

Adaptive Dynamic Programming: Solltrajektorienfolge- regelung und Konvergenzbedingungen

Zur Erlangung des akademischen Grades eines
DOKTOR-INGENIEURS
von der KIT-Fakultät für
Elektrotechnik und Informationstechnik
des Karlsruher Instituts für Technologie (KIT)
genehmigte

DISSERTATION

von
Florian Köpf, M.Sc.
geb. in Leonberg

Tag der mündlichen Prüfung:	13. Januar 2022
Hauptreferent:	Prof. Dr.-Ing. Sören Hohmann
Korreferent:	Prof. Dr.-Ing. Daniel Görge

Danksagung

Diese Dissertation entstand während meiner Zeit als wissenschaftlicher Mitarbeiter am Institut für Regelungs- und Steuerungssysteme (IRS) am Karlsruher Institut für Technologie (KIT). Für die dort gebotene Gelegenheit zur fachlichen und persönlichen Weiterentwicklung möchte ich mich an dieser Stelle herzlich bedanken.

Besonderer Dank gilt meinem Hauptreferenten Prof. Dr.-Ing. Sören Hohmann für seine ausgezeichnete Betreuung. Die zahlreichen anregenden Diskussionen und angenehmen Gespräche haben mir stets neue Denkanstöße geliefert und das mir entgegengebrachte Vertrauen und die gebotenen Freiheiten haben wesentlich zum Gelingen dieser Arbeit beigetragen. Herzlich danke ich auch Prof. Dr.-Ing. Daniel Görge für die freundliche Übernahme des Korreferats, spannende Gespräche und das Interesse an der vorliegenden Arbeit.

Der gemeinsame Austausch mit meinen Kolleginnen und Kollegen am IRS und die gegenseitige Unterstützung haben für eine bemerkenswerte Arbeitsatmosphäre gesorgt. Neben unzähligen fachlichen Diskussionen haben auch die persönlichen Gespräche und gemeinsamen Erlebnisse ein inspirierendes Umfeld erschaffen. Besonderer Dank hierfür gilt Jairo Inga und Michael Flad, die mich als Gruppenleiter über fachliche Belange hinaus bei meiner Tätigkeit als wissenschaftlicher Mitarbeiter unterstützt haben. Esther Bischoff, Christian Braun, Philipp Karg, Simon Rothfuß, Julian Schneider und abermals Jairo Inga, mit denen teilweise gemeinsame Publikationen entstanden sind, haben wertvolle Anmerkungen zur vorliegenden Arbeit geliefert. Lukas Kölsch, Simon Rothfuß und Manuel Schwartz haben daneben den Arbeitsalltag in der „Legebatterie“ bereichert, während Mathias Kluwe und Beate Stassen durch ihre Hilfsbereitschaft und herzliche Art wichtige Säulen des Instituts darstellen. Dankbar bin ich auch Luca Puccetti, der mich im Rahmen eines gemeinsamen Forschungsprojekts in Aschheim auf die Teststrecke mitgenommen hat und der immer für fachliche Diskussionen zu begeistern war, sowie Sean Kille für seine Unterstützung am Ball-auf-Platte-System.

Ganz herzlich bedanke ich mich zudem bei allen Studierenden, deren Abschlussarbeiten ich im Rahmen meiner Tätigkeit am IRS betreuen durfte. Die Zusammenarbeit hat mir dabei viel Freude bereitet und der hierbei geführte fachliche Austausch lieferte wertvolle Beiträge zum Erfolg der vorliegenden Dissertation.

Für den gebotenen Ausgleich während der letzten Jahre danke ich meinen Freunden und meiner Familie von Herzen. Meine Eltern Karin und Jürgen sind immer für mich da und unterstützen mich jederzeit uneingeschränkt. Katharina hat mich unermüdlich zwischen all den Formeln und Beweisen geerdet und ihr gewissenhaftes Lektorat der vorliegenden Arbeit hat den letzten Feinschliff gegeben. Danke, dass du für mich da bist.

Karlsruhe, im Januar 2022

The outcome of any serious research can only be to make two questions grow where only one grew before.

Thorstein Bunde Veblen

Kurzfassung

Adaptive Dynamic Programming (ADP) steht als vielversprechendes und zukunftsorientiertes regelungstechnisches Werkzeug im Fokus der aktuellen Forschung. Allerdings existieren hierfür bislang weder flexibel einsetzbare, mit dem ADP-Mechanismus kompatible Solltrajektorien- darstellungen noch theoretische Untersuchungen hinsichtlich einer geeigneten System- anregung zur Sicherstellung der Konvergenz.

Die vorliegende Arbeit schließt diese Lücken: Zum einen werden erstmals zeitdiskrete und zeitkontinuierliche Methoden präsentiert und analysiert, die flexible Solltrajektorien- darstellungen in ADP-Ansätze integrieren. Die explizite Abhängigkeit der vorgestellten, neuartigen Value- bzw. Q-Function und des darauf basierenden gelernten Regelgesetzes von Trajektorien- parametern, die den aktuellen Sollverlauf repräsentieren, ermöglicht eine variable Vorgabe der Solltrajektorie zur Laufzeit. Zum anderen werden erstmalig theoretische Bedingungen an den Systemzustand hergeleitet, die sicherstellen, dass eine für die Konvergenz der Adap- tion zentrale Anregungseigenschaft erfüllt ist. Verbleibende Freiheitsgrade erlauben zudem die Berücksichtigung anwendungsspezifischer Anforderungen bei der Systemanregung. Die theoretischen Aussagen werden in Simulationen bestätigt.

Erste reale Anwendungen der vorgestellten adaptiven optimalen Trajektorienfolgeregelungs- methoden offenbaren schließlich das Potenzial dieser Ansätze. Flexible und effiziente Regler, die aufgrund der Berücksichtigung des Solltrajektorienverlaufs vorausschauend agieren, kön- nen ohne aufwendige Modellbildung aus realen Messdaten erlernt werden und sind zudem bisherigen Ansätzen bezüglich ihrer Performanz überlegen.

Inhaltsverzeichnis

Danksagung	i
Kurzfassung	iii
Abbildungsverzeichnis	ix
Tabellenverzeichnis	xiii
Abkürzungen und Symbole	xv
1 Einleitung	1
2 Stand der Wissenschaft und Forschungslücke	5
2.1 Einführung der Notation und Grundbegriffe des ADP	5
2.1.1 Zeitdiskrete ADP-Grundgleichungen	6
2.1.2 Zeitkontinuierliche ADP-Grundgleichungen	8
2.1.3 Funktionsapproximatoren	10
2.1.4 ADP-Lösungsansätze	12
2.2 ADP-basierte Solltrajektorienfolgeregler	17
2.2.1 ADP-Ansätze unter Nutzung der dynamischen Inversion	18
2.2.2 Globale Solltrajektorienvorgabe durch eine Exosystemdynamik ohne externe Eingriffsmöglichkeit	20
2.2.3 Stationäre Sollzustandsvorgabe	22
2.2.4 ADP in realen regelungstechnischen Anwendungen	23
2.2.5 Fazit	24
2.3 Anregung ADP-basierter Regelungsansätze	25
2.4 Wissenschaftliche Fragestellungen und Beiträge der Arbeit	30
2.4.1 ADP-kompatible, flexible Solltrajektorien-darstellung	30
2.4.2 Konvergenz ADP-basierter Regelungsansätze	31
3 Zeitdiskrete ADP-basierte Solltrajektorienfolgeregelung	33
3.1 Definition ADP-kompatibler zeitdiskreter Trajektorien	33
3.2 Zeitdiskrete ADP-kompatible parametrisierte Referenztrajektorie	39
3.2.1 Allgemeine Problemstellung	40
3.2.2 Q-Function mit parametrierter Referenzdarstellung	41
3.2.3 Funktionsapproximation und Policy Iteration der erweiterten Q-Function	43

3.2.4	Linear-quadratische optimale Trajektorienfolgeregelung mit parametrierter Referenz	47
3.2.5	Implementierung und simulativer Vergleich des PRADP	56
3.3	Zeitdiskrete ADP-kompatible Referenztrajektorie auf einem endlichem Vorausschauhorizont	66
3.3.1	Problemdefinition	67
3.3.2	Solltrajektorienabhängige Q-Function	70
3.3.3	Modellfreies Erlernen der Q-Function	74
3.3.4	Ergebnisse	85
3.3.5	Diskussion der Simulationsergebnisse	93
3.4	Resümee zur zeitdiskreten ADP-basierten Solltrajektorienfolgeregelung	94
4	Zeitkontinuierliche ADP-basierte Solltrajektorienfolgeregelung	97
4.1	Definition ADP-kompatibler zeitkontinuierlicher Trajektorien	97
4.2	Zeitkontinuierliche ADP-kompatible parametrisierte Referenztrajektorie	99
4.2.1	Solltrajektorienendarstellung	99
4.2.2	Trajektorienfolgeregelung mit global diskontiertem Gütemaß	104
4.2.3	Trajektorienfolgeregelung mit gedämpfter Referenzdynamik	108
4.2.4	ADP-Umsetzung	114
4.2.5	Simulationsergebnisse	117
4.3	Zusammenfassung	126
5	Konvergenzbedingungen zeitkontinuierlicher adaptiver Optimalregler	129
5.1	Eingangsaffines Differenzialspiel mit unbekanntem Gegenspielern	129
5.2	Policy Iteration für Nicht-Nullsummen-Differenzialspiele	131
5.3	Funktionsapproximation und Anregungsbedingung	133
5.4	Hinreichende Bedingungen zur Erfüllung der PE-Eigenschaft in ADP-basierten Differenzialspielen	138
5.4.1	Hilfsmenge $\Omega^{(1)}$	142
5.4.2	Hilfsmenge $\Omega^{(2)}$	144
5.4.3	Frequenzbedingungen Ω	148
5.5	Signal zur Überprüfung der Erfüllung der PE-Eigenschaft	153
5.6	Anregungssignale für ADP-basierte Differenzialspiele	155
5.7	Simulationsbeispiel zur Anregung von ADP-basierten Differenzialspielen	162
5.7.1	Beispielproblem	163
5.7.2	Konstruktion geeigneter Anregungssignale	164
5.7.3	Simulationsergebnisse	166
5.8	Diskussion	173
6	Reale Anwendung ADP-basierter Solltrajektorienfolgeregler	175
6.1	Modellfreie, adaptive Längsregelung eines realen Fahrzeugs	175
6.1.1	Problemstellung	177
6.1.2	Modellfreier ADP-Solltrajektorienfolgeregler mit Zustandsrekonstruktion	178

6.1.3	ADP-Solltrajektorienfolgeregelung im Realfahrzeug	184
6.1.4	Diskussion	190
6.2	Modellfreier Trajektorienfolgeregler für ein reales Ball-auf-Platte-System	191
6.2.1	Ball-auf-Platte-System und Problembeschreibung	193
6.2.2	ADP-Solltrajektorienfolgeregler für ein Ball-auf-Platte-System	195
6.2.3	Ergebnisse	201
6.2.4	Diskussion	207
7	Zusammenfassung	211
A	Anhang zu Kapitel 3	I
A.1	Beweisskizze zu Satz 3.4	I
A.2	Beweis zu Lemma 3.8	IV
A.3	Beweis zu Lemma 3.9	V
A.4	Ergänzungen zum linearen Einspurmodell	VI
A.5	ADP-Solltrajektorienregler für das lineare Einspurmodell mit $\gamma = 1$	VII
B	Anhang zu Kapitel 4	IX
B.1	Beweis zu Lemma 4.1	IX
B.2	Beweis zu Lemma 4.2	IX
B.3	Zustandstransformation für reellwertige $\zeta(t)$ und D	XI
B.4	Äquivalenz zwischen Problem 4.1 und Problem 4.2	XI
B.5	Beweis zu Lemma 4.3	XII
B.6	Beweis zu Lemma 4.4	XIII
B.7	Beweis zu Lemma 4.6	XIII
B.8	Value Iteration nach Bian und Jiang [BJ16a]	XIV
C	Anhang zu Kapitel 5	XVII
C.1	Beweisskizze zu Lemma 5.2	XVII
C.2	Beweis zu Lemma 5.9	XVIII
C.3	M_i und $v_{\text{freq}}(t)$ für das verwendete Beispielsystem	XVIII
D	Anhang zu Kapitel 6	XXI
D.1	Exkurs: Reinforcement Learning mit unvollständiger Zustandsinformation	XXI
D.2	Actor-Critic-Grundlagen	XXI
D.3	Wahl der Hyperparameter für das Online-Training	XXIII
D.4	Ergänzende Messdaten des realen Ball-auf-Platte-Systems	XXIV
	Literaturverzeichnis	XXXI

Abbildungsverzeichnis

1.1	Schematische Übersicht der Hauptbeiträge der vorliegenden Arbeit.	4
2.1	Klassifikation von ADP-Methoden.	13
2.2	Literaturansatz ohne äußere Beeinflussbarkeit der Solltrajektorie.	22
2.3	Literaturansatz unter Vorgabe eines stationären Sollzustands.	23
2.4	Schematische Darstellung des gewünschten generalisierenden ADP-Solltrajektorienfolgeregelungsansatzes.	25
3.1	ADP-kompatible Approximation einer Solltrajektorie.	39
3.2	Beispielhafter Ausschnitt des Solltrajektorienverlaufs und lokale Approximation durch kubische Polynome.	59
3.3	Ablaufschema des PRADP-Algorithmus.	61
3.4	Ergebnis PRADP für Szenario 1.	63
3.5	Ergebnis PRADP für Szenario 2.	64
3.6	Ergebnis PRADP für Szenario 3.	65
3.7	Ergebnis der ADP-Solltrajektorienregelung für System 1.	89
3.8	Ergebnis der ADP-Solltrajektorienregelung für System 2 ($\gamma = 0,9$).	90
3.9	Gewichtsfehlerverlauf während des Lernvorgangs für System 1.	90
3.10	Gewichtsfehlerverlauf während der ersten 30 Iterationen des Lernvorgangs für System 2 ($\gamma = 0,9$).	91
3.11	Trajektorienfolgeergebnisse für ADP-Regler, die bei unterschiedlich starkem Messrauschen trainiert und validiert wurden.	92
3.12	Detailansicht von Abbildung 3.7 (System 1), um das prädiktive Verhalten des vorgestellten ADP-Ansatzes zu visualisieren.	93
4.1	Verlauf der Systemtrajektorie für das betrachtete Beispielsystem bei globaler Diskontierung.	108
4.2	Verlauf der Systemtrajektorie für das betrachtete Beispielsystem bei teilweiser Dämpfung.	113
4.3	Struktur des vorgestellten zeitkontinuierlichen ADP-Trajektorienfolgeregelungsansatzes.	116
4.4	Reglergewichte und Gewichtsfehlernorm für das Beispiel der stationären Sollvorgabe.	118
4.5	Sollzustandsverlauf, Zustandsverlauf und Stellgröße für das Beispiel des harmonischen Oszillators als Solltrajektorienparametrierung im Vergleich zur stationären Sollzustandsvorgabe.	120

4.6	Reglergewichte und Gewichtsfehlnorm für das Beispiel des harmonischen Oszillators als Solltrajektorienparametrierung.	120
4.7	Rang der Datenmatrizen D_ϕ und D_ψ für das Beispiel des harmonischen Oszillators als Solltrajektorienparametrierung.	121
4.8	Resultierende Kosten für das Beispiel des harmonischen Oszillators als Solltrajektorienparametrierung im Vergleich zur stationären Sollzustandsvorgabe.	121
4.9	Sollzustandsverlauf, Zustandsverlauf und Stellgröße für das Beispiel der Polynomvorgabe im Vergleich zur stationären Sollzustandsvorgabe.	123
4.10	Reglergewichte und Gewichtsfehlnorm für das Beispiel der Polynomvorgabe.	123
4.11	Rang der Datenmatrizen D_ϕ und D_ψ für das Beispiel der Polynomvorgabe.	124
4.12	Resultierende Kosten für das Beispiel der Polynomvorgabe.	124
4.13	Sollzustandsverlauf, Zustandsverlauf und Stellgröße für das Beispiel der Polynomvorgabe unter Einfluss von Messrauschen.	125
4.14	Gewichtsfehlnorm für das Beispiel der Polynomvorgabe unter Einfluss von Messrauschen.	126
5.1	Strukturdiagramm Differenzialspiel mit globaler Anregung.	156
5.2	Ablaufdiagramm zum Design eines Anregungssignals u_{ex}	158
5.3	Grafische Veranschaulichung der Frequenzmenge Ω	165
5.4	Verlauf der Critic-Gewichte unter Verwendung des Anregungssignals $u_{ex,1,1}$ bzw. weißen Rauschens mit normierten Lernraten.	169
5.5	Critic-Fehlnorm für die Anregungssignale $u_{ex,1,j}$ mit $j \in \{4, 5, 6\}$ und $u_{ex,1,w}$ mit normierten Lernraten.	170
5.6	Eigenwertsignal $\lambda_{\min,2}(t)$ zur Überprüfung der PE-Eigenschaft nach Lemma 5.11.	170
5.7	Eigenwertsignal $\lambda_{\min,1}(t)$ zur Überprüfung der PE-Eigenschaft nach Lemma 5.10.	171
5.8	Anregungssignal $u_{ex,1,1}$ und $u_{ex,1,w}$	171
5.9	Verlauf der Zustandsgröße $x(t)$ für die Anregungssignale $u_{ex,1,1}$ und $u_{ex,1,w}$	172
5.10	Verlauf der Critic-Gewichtsfehlnorm unter Verwendung des Anregungssignals $u_{ex,1,1}$ bzw. weißen Rauschens, welche die gleiche mittlere Anregungssignalleistung aufweisen, bei identischen Lernraten.	172
6.1	Versuchsfahrzeug (<i>BMW 740Li</i>), das für die modellfreie, adaptive Längsregelung mittels eines Online-ADP-Ansatzes verwendet wurde.	176
6.2	Struktur der gegebenen Problemstellung zur Längsregelung eines realen Fahrzeugs.	178
6.3	Verwendete Netzwerkarchitektur für die Approximation der Q-Function.	181
6.4	Übersicht über den Gesamtalgorithmus des ADP-basierten Geschwindigkeitsreglers.	182
6.5	Beispielhafter Ausschnitt der Geschwindigkeit y_k sowie der Stellgröße u_k während des Trainingsvorgangs.	186
6.6	Verteilung der während des Trainingsvorgangs verwendeten verrauschten Referenzparameter.	187
6.7	Reglerparameter $\theta^{[l]}$ während des Trainingsvorgangs.	187
6.8	Sollgeschwindigkeitsverlauf und gemessene Geschwindigkeit während der Auswertungsfahrt.	188
6.9	Solltrajektorienparameter während der Auswertungsfahrt.	189

6.10	Stellgröße während des Auswertungsmanövers.	189
6.11	Approximierte Value Function während des Auswertungsmanövers.	190
6.12	Ball-auf-Platte-System für die modellfreie, ADP-basierte Solltrajektorienfolgeregelung.	192
6.13	Architektur des Ball-auf-Platte-Systems.	194
6.14	Mithilfe der mechanischen Drehregler durch einen Menschen aufgezeichnete, geglättete Messdaten des Ball-auf-Platte-Systems.	200
6.15	Ablaufschema des Trainingsvorgangs für den modellfreien ADP-Solltrajektorienfolgeregler für das Ball-auf-Platte-System.	202
6.16	Gewichtsvektor $\hat{w}^{[l]}$ über die Iterationen l des LSPI-Algorithmus.	203
6.17	Vorgabe einer stationären Sollposition für den gelernten, ADP-basierten Regler sowie den modellbasierten Vergleichsregler.	204
6.18	Vorgabe eines polynomiellen Sollverlaufs mit $d = 2$ für den gelernten, ADP-basierten Regler sowie den modellbasierten Vergleichsregler im Vergleich zur Vorgabe einer stationären Sollposition für den gelernten, ADP-basierten Regler.	205
6.19	Vergleich der gelernten Solltrajektorienfolgeregler bei Vorgabe eines polynomiellen Sollverlaufs mit einer stationären Sollpositionsvorgabe für eine Validierungstrajektorie.	206
6.20	Vergleich eines gelernten Reglers mit und ohne lernbare Offsetkorrektur bei Vorgabe einer stationären Sollposition.	207
6.21	Gleichzeitige Vorgabe eines Solltrajektorienverlaufs in beiden Plattendimensionen und resultierende Ballposition für den gelernten Regler und den modellbasierten Vergleichsregler.	208
A.1	Geometrische Zusammenhänge des linearen Einspurmodells.	VII
A.2	Ergebnis der ADP-Solltrajektorienregelung für System 2 ($\gamma = 1$).	VII
A.3	Gewichtsfehlerverlauf während des Lernvorgangs für System 2 ($\gamma = 1$).	VIII
D.1	Zustände und Stellgröße bei Vorgabe einer stationären Sollposition für den gelernten, ADP-basierten Regler sowie den modellbasierten Vergleichsregler.	XXIV
D.2	Zustände und Stellgröße bei Vorgabe eines polynomiellen Sollverlaufs für den gelernten, ADP-basierten Regler sowie den modellbasierten Vergleichsregler im Vergleich zur Vorgabe einer stationären Sollposition für den gelernten, ADP-basierten Regler.	XXV
D.3	Vergleich der gelernten Solltrajektorienfolgeregler bei Vorgabe eines polynomiellen Sollverlaufs und einer stationären Sollposition für eine Validierungstrajektorie.	XXVI
D.4	Vergleich eines gelernten Reglers mit und ohne lernbare Offsetkorrektur bei Vorgabe einer stationären Sollposition.	XXVII
D.5	Gleichzeitige Vorgabe eines Sollpositionsverlaufs in beiden Plattendimensionen für den gelernten Regler (X -Richtung).	XXVIII
D.6	Gleichzeitige Vorgabe eines Sollpositionsverlaufs in beiden Plattendimensionen für den gelernten Regler (Y -Richtung).	XXIX

Tabellenverzeichnis

2.1	Übersicht über ADP-Trajektorienfolgeregelungsmethoden, die auf dem Konzept der dynamischen Inversion basieren.	19
2.2	Übersicht über ADP-Trajektorienfolgeregelungsmethoden, die eine globale Vorgabe der Referenztrajektorie annehmen.	21
2.3	Anregungsheuristiken verschiedener ADP-Ansätze.	29
3.1	Beispielwerte für das hinreichende Stabilitätskriterium nach Satz 3.3.	56
3.2	Verringerung der Anzahl der zu schätzenden Gewichte durch Analyse der Struktur von \mathbf{H}	77
3.3	Trajektorienfolgefehler und Gewichtsfehler der betrachteten Simulationsbeispiele.	89
3.4	Abweichungen der gelernten Regler vom optimalen Regler für unterschiedlich starkes Messrauschen.	92
4.1	Parameterwahl für die harmonische Referenzzustandsdarstellung.	119
4.2	Parameterwahl für die Polynomvorgabe.	122
5.1	Matrizen \mathbf{C}_z der Frequenzbedingungen Ω für das Simulationsbeispiel nach Abschnitt 5.7.	165
5.2	Normierte Lernraten, Konvergenzzeiten und Zeitkonstanten für die Anregungssignale $u_{\text{ex},1,j}(t)$, $j \in \{1, 2, 3\}$, und $u_{\text{ex},1,w}(t)$	168
D.1	Wahl der Hyperparameter für das Online-Training.	XXIII

Abkürzungen und Symbole

Abkürzungen

Abkürzung	Beschreibung
ADP	Adaptive Dynamic Programming
BINF	Byrnes-Isidori-Normalform
CAN	Controller Area Network
DPG	Deterministic Policy Gradient
FIR	Finite Impulse Response
GUI	Graphical User Interface
HDP	Heuristic Dynamic Programming
HJB	Hamilton-Jacobi-Bellman
IRL	Integral Reinforcement Learning
LQ	linear-quadratisch
LSPI	Least-Squares Policy Iteration
o. B. d. A.	ohne Beschränkung der Allgemeinheit
PE	persistently excited
PI	Policy Iteration
PRADP	Parametrized Reference Adaptive Dynamic Programming
QEP	Quadrature Encoder Pulse
RL	Reinforcement Learning
SR	sufficiently rich
TD	temporal difference
u. d. N.	unter der Nebenbedingung
USB	Universal Serial Bus
VI	Value Iteration

Lateinische Buchstaben

Symbol	Beschreibung
<i>A</i>	autonomer Systemanteil bei linearem System
<i>a</i>	Amplituden von Sinusfunktionen bzw. Beschleunigung
<i>B</i>	Eingangsdynamikmatrix bei linearem System

Symbol	Beschreibung
b	Frequenzvorfaktor
C_z	Matrizen zur Definition der Frequenzbedingungen Ω
c	Amplituden Kosinusfunktionen
c_{norm}	Normierungsfaktor
D	Dynamikmatrix Referenzparameter bzw. Entkopplungsmatrix
d	Frequenzvorfaktor bzw. Polynomgrad
e	Einheitsvektor bzw. Abweichung vom Sollzustand
e	Restterm
$e_{\hat{w}}$	Schwellwert für Abbruchbedingung des LSPI-Algorithmus
$F(\cdot)$	allgemeine Systemdynamik
F_{ref}	Dynamikmatrix des Sollzustands x_r
$f(\cdot)$	autonomer Systemanteil
$f_{x_r}(\cdot)$	Dynamik des Sollzustands x_r
$f_{x_r, x}(\cdot)$	Abbildung zwischen Systemzustand x und Sollzustand x_r
$f_{\zeta}(\cdot)$	Dynamik des Referenzparameters ζ
$f_{x_r, \zeta}(\cdot)$	Abbildung zwischen Referenzparameter ζ und Sollzustand x_r
f	Potenzen der Basisfunktionsmonome
\bar{f}	Amplituden der Basisfunktionselemente
$g(\cdot)$	Eingangsfunktion der Systemdynamik
g	Amplituden Sinusfunktionen
\bar{g}	Amplituden Kosinusfunktionen
$H(\cdot)$	Hamilton-Funktion
h	Dimension des Gewichtsvektors w
\bar{h}	Laufindex
h_{FIR}	Länge des FIR-Filters
I, I_n	Einheitsmatrix bzw. Einheitsmatrix der Dimension $n \times n$
$I^{[X]}, I^{[Y]}$	Motorstrom in X - bzw. Y -Richtung
i	Laufindex
$J(\cdot)$	Gütefunktional
j	Laufindex
j	imaginäre Einheit
K	Reglermatrix
K	Summenobergrenze
k	diskreter Zeitschritt bzw. Laufindex
L	Summenobergrenze
l	Iterationsindex/Trainingsschritt bzw. Laufindex
M	Matrix
M	Sampleanzahl Training
M_B	Batchgröße Training
m	Anzahl an Frequenzen bzw. Index
N	Anzahl der Spieler/Regler
N_P	Anzahl verschiedener Frequenztuple \mathcal{P}_i

Symbol	Beschreibung
n	Systemordnung
n_h	Vorausschauhorizontlänge der Solltrajektorie
n_z	Dimension von z
n_ζ	Dimension des Referenzparameters ζ
o	Laufindex
P	Riccati-Matrix
P^f	Frequenzkombinationen
P_{LS}	Projektionsmatrix der (gewichteten) Least-Squares-Regression
p	Anzahl der Eingangsgrößen
Q	Matrix zur quadratischen Bestrafung des Systemzustands
$Q(\cdot)$	Q-Funktion
Q_y	Faktor zur quadratischen Bestrafung des Ausgangsfehlers
$q(\cdot)$	Bestrafungsfunktion Systemzustand
R	Matrix zur quadratischen Bestrafung der Stellgröße
r	Einschrittkostenvektor
$r(\cdot)$	Einschrittkosten
S	Projektionsmatrix
s	Projektionsvektor
$s^{[X]}, s^{[Y]}$	Ballposition in X - bzw. Y -Richtung
s_r	Soll-Ballposition in ADP-kompatibler Form
$s_{r,soll}$	Soll-Ballposition
T^f	trigonometrischer Typ zu Frequenzkombinationen P^f
T	Zeitpunkt
T_d	Konstante für zeitdiskrete PE-Definition
T_{IRL}	Integrationsdauer des Integral-Reinforcement-Learning-Intervalls
T_{konv}	Konvergenzzeit der Policy Iteration
T_{sim}	Simulationsdauer
T_t	Schwellwert
t	Zeit
Δt	Abtastzeit
u	Stellgröße
u_{ex}	Anregungssignal
$u_{ex,1,w}$	Anregungssignal weißes Rauschen
$V(\cdot)$	Value Function
v_{freq}	Vektor aus Sinus- und Kosinusfunktionen
$v^{[X]}, v^{[Y]}$	Ballgeschwindigkeit in X - bzw. Y -Richtung
W_p	Gewichtungsmatrix der gewichteten Least-Squares-Regression
$W_{\bar{h}}$	Frequenzmenge
w	Gewichtsvektor
w_{FIR}	Gewichtsvektor FIR-Filter
X	Hilfsmatrix
X	Plattendimension in X -Richtung

Symbol	Beschreibung
\boldsymbol{x}	Systemzustand
\boldsymbol{x}_0	Anfangszustand
\boldsymbol{x}_r	Sollzustand in ADP-kompatibler Form
$\boldsymbol{x}_{r,\text{soll}}$	Sollzustand
$\boldsymbol{x}_{\text{PE}}$	Anregungstrajektorie
$\tilde{\boldsymbol{x}}$	erweiterter Zustand
$\tilde{\boldsymbol{x}}^\diamond$	erweiterter Zustand mit Anregungsrauschen auf Referenzparameter
$\tilde{\boldsymbol{x}}^{\diamond\diamond}$	erweiterter Zustand mit kompatibelem Referenzparameter
$\bar{\boldsymbol{x}}$	Hilfzustand
\boldsymbol{Y}	Hilfsmatrix
Y	Plattendimension in Y -Richtung
\boldsymbol{y}	realer oder fiktiver Systemausgang
$\bar{\boldsymbol{y}}, \tilde{\boldsymbol{y}}, \tilde{\boldsymbol{y}}$	Hilfsvektoren
y	Ausgang bzw. laterale Position
y_r	Sollausgang in ADP-kompatibler Form
$y_{r,\text{soll}}$	Sollausgang
\boldsymbol{Z}	Solltrajektorienparametermatrix bzw. Hilfsmatrix
\boldsymbol{z}_B	Systemzustand in BINF
$z_{k,i}$	Solltrajektorienparameter für Zustand i im Zeitschritt k
$\bar{\boldsymbol{z}}$	Hilfsvektor
z	Laufindex

Griechische Buchstaben

Symbol	Beschreibung
$\boldsymbol{\alpha}$	Vektor
α, α_I	Grad der PE
$\alpha^{[X]}, \alpha^{[Y]}$	Plattenwinkel in X - bzw. Y -Richtung
α_t	Schwellwert
β	Faktor für gewichtete Least-Squares-Schätzung bzw. Schwimmwinkel
$\boldsymbol{\gamma}(\cdot)$	funktionaler Zusammenhang der Spalten der Eingangsmatrix $\boldsymbol{g}(\boldsymbol{x})$
γ	Diskontierungsfaktor
$\boldsymbol{\delta}$	TD-Fehlervektor
δ	relativer Grad/Differenzordnung bzw. TD-Fehler
δ_{LR}	Lenkradwinkel
$\delta_{\text{LR},v}$	Winkelgeschwindigkeit des Lenkrads
δ_s	Lenkwinkel
ϵ	Approximationsfehler
$\bar{\epsilon}$	obere Schranke des Approximationsfehlers
ϵ_1	beschränktes Fehlersignal

Symbol	Beschreibung
ε_2	transienter Fehler
$\bar{\varepsilon}_\sigma$	obere Schranke von ε_1
ζ	Referenzparametervektor
ζ	Frequenzvariable
η	Lernrate
θ	Actor-Gewicht
ι	Platzhalter
κ	Laufindex
$\lambda_{\min,1}$	Eigenwertsignal (minimaler Eigenwert von Ξ_1)
$\lambda_{\min,2}$	Eigenwertsignal (minimaler Eigenwert von Ξ_2)
λ_{LM}	Regularisierung des Levenberg-Marquardt-Algorithmus
$\mu(\cdot)$	Regelgesetz
ν	Skalierungsfaktor bzw. Messrauschen
Ξ_1	PE-Matrix im Intervall $[t, t + T]$
Ξ_2	PE-Matrix im Intervall $[t_0, t]$
$\rho(\cdot)$	Basisfunktionen Referenzapproximation
ρ	Konvergenzrate
σ	für Adaption relevantes Signal
σ_s	skalares Signal
τ	Zeit/Integrationsvariable
Φ	Hilfsmatrix für den LSPI-Algorithmus (Minimierung quadrierter TD-Fehler)
$\bar{\Phi}$	Hilfsmatrix für den LSPI-Algorithmus (mit Fixpunktforderung)
$\phi(\cdot)$	Basisfunktionsvektor
Ψ_1	funktionaler Zusammenhang zwischen Zustand und flachem Ausgang sowie dessen Ableitungen
Ψ_2	funktionaler Zusammenhang zwischen Stellgröße und flachem Ausgang sowie dessen Ableitungen
Ψ	Menge der zulässigen Regelgesetze
$\psi(\cdot)$	Basisfunktionsvektor
ψ	Gierwinkel
ψ_{\ddagger}	Gierrate
Ω	Menge an Frequenzbedingungen
ω	Frequenzvektor
$\omega, \underline{\omega}$	Frequenz
$\omega^{[X]}, \omega^{[Y]}$	Winkelgeschwindigkeit Platte in X - bzw. Y -Richtung

Kalligraphische und andere Symbole

Symbol	Beschreibung
$\mathbf{0}_{n \times m}$	Nullmatrix der Dimension $n \times m$
\mathcal{B}	Bellman-Operator
$\mathcal{C}^l(\mathcal{M})$	Klasse der auf einer Menge \mathcal{M} l -fach stetig differenzierbaren Funktionen
\mathcal{D}	Menge der Dimensionen des Ball-auf-Platte-Systems
\mathbb{N}	Menge der natürlichen Zahlen
\mathcal{N}	Menge an Spielern/Reglern
$\mathcal{N}(\mu; \sigma^2)$	Normalverteilung mit Mittelwert μ und Standardabweichung σ
\mathcal{P}	Frequenztuplel
\mathbb{R}	Menge der reellen Zahlen
$\mathbb{R}_{\geq 0}$	Menge der nicht-negativen reellen Zahlen
$\mathcal{S}\{\cdot\}$	Transformation eines Signals
\mathcal{T}	Datentuplel
\mathcal{X}	Menge der Systemzustände
\mathcal{Z}	Definitionsbereich des Referenzparametervektors ζ

Indizes, Exponenten und Operatoren

Symbol	Beschreibung
$\text{Bild}(\cdot)$	Bild einer Matrix
$\det(\cdot)$	Determinante einer Matrix
$\text{diag}(\cdot)$	Diagonalmatrix
$\text{Dim}(\cdot)$	Dimension einer Matrix
$\text{eig}(\cdot)$	Eigenwerte einer Matrix
$\text{Im}\{\cdot\}$	Imaginärteil
$L_f^j h(\mathbf{x})$	j -te Lie-Ableitung von $h(\mathbf{x})$ bezüglich $\mathbf{f}(\mathbf{x})$
$\text{mat}(\cdot)$	Konvertierung eines Vektors in eine Matrix, wobei $\text{mat}(\text{vec}(\mathbf{M}), n, p) = \mathbf{M} \in \mathbb{R}^{n \times p}$
mod	Modulo
$\text{Rang}(\cdot)$	Rang einer Matrix
$\text{Re}\{\cdot\}$	Realteil
$\text{sgn}(\cdot)$	Signumfunktion
$\text{std}(\cdot)$	Standardabweichung
$\text{vec}(\cdot)$	Vektorisierung einer Matrix (vertikale Konkatenation der Spalten)
$\text{vecr}(\cdot)$	Vektorisierung einer symmetrischen Matrix \mathbf{M} , sodass $\mathbf{x}^\top \mathbf{M} \mathbf{x} = (\mathbf{x} \otimes_r \mathbf{x})^\top \text{vecr}(\mathbf{M})$
\otimes	Kronecker-Produkt
\otimes_r	reduziertes Kronecker-Produkt mit nicht-redundanten Elementen

Symbol	Beschreibung
$\lambda_{\text{Re}}^+(\cdot)$	größter Eigenwert einer reellwertigen Matrix
$\lambda_{\text{min}}(\cdot)$	kleinster Eigenwert einer Matrix
$\nabla_{\mathbf{x}}$	Gradient bezüglich \mathbf{x} , sodass $\nabla_{\mathbf{x}} V(\mathbf{x}) := \left[\frac{\partial V(\mathbf{x})}{\partial x_1} \quad \dots \quad \frac{\partial V(\mathbf{x})}{\partial x_n} \right]^T$
∂	partielle Ableitung
\succ	$M_1 \succ M_2$: Matrix $M_1 - M_2$ positiv definit
\succeq	$M_1 \succeq M_2$: Matrix $M_1 - M_2$ positiv semidefinit
\forall	Allquantor
\exists	Existenzquantor
$\ \cdot\ _2$	Spektralnorm einer Matrix bzw. Euklidische Norm eines Vektors
$\ \cdot\ _\infty$	Maximumsnorm
\cdot	Platzhalter
$[\cdot]_{n_1:n_2}$	Teilvektor: $[v]_{n_1:n_2} = [v_{n_1} \quad \dots \quad v_{n_2}]^T$ mit $n_1, n_2 \in \mathbb{N}_{>0}, n_1 < n_2$
$\hat{\square}$	geschätzte Größen
$\check{\square}$	aus exakter Systeminversion resultierende Größen
$\dot{\square}$	zeitliche Ableitung
\square^*	Optimalität bzw. Nash-Lösung
\square^*	Komplexe Konjugation
$\bar{\square}$	Kennzeichnung von Hilfsgrößen bzw. Komplementärmengen
$\tilde{\square}$	Fehlersignal
$\tilde{\square}$	Kennzeichnung von Größen eines erweiterten Systems bzw. Hilfsgrößen
\square^T	Transposition einer Matrix oder eines Vektors
$\square^{(j)}(t)$	j -te zeitliche Ableitung einer zeitkontinuierlichen Größe
$\square^{(\kappa)}$	um κ Zeitschritte verschobene zeitdiskrete Größe
$\square^{\{k\}}$	Kennzeichnung von Größen aus lokaler Perspektive im Zeitschritt k
$\square^{[l]}$	Kennzeichnung von Größen der l -ten Iteration
$\square^{\langle \cdot \rangle}$	Indizierung zur Kennzeichnung von Abhängigkeiten von einer Größe

1 Einleitung

Lernbasierte, adaptive Optimalregelungsmethoden rücken zunehmend in den Fokus der aktuellen Forschung [WHL17], [SB18], [KHL⁺12]. Während modellbasierte optimale Regelungsansätze (vgl. beispielsweise [Gee07], [SMC⁺11], [TBBH10]) nur anwendbar sind, wenn ein Systemmodell vorhanden ist [LV09], [VLV13, S. 1], liegt der Fokus von selbstlernenden Optimalreglern auf Anwendungen, bei denen kein oder nur unvollständiges Systemwissen vorliegt, eine Modellbildung aufwendig ist oder unbekannte Parameter schwierig zu ermitteln sind [HJK20]. Adaptive, lernende Ansätze erlauben somit trotz vorhandener Modellunsicherheiten den Entwurf leistungsfähiger, flexibler Regler [Wer99], [Wer92], [SBW92], [Tao03, S. 10], [HLM03]. Flexible Reglerentwurfsmethoden [CF14], [LLSX14] und damit einhergehende Fertigungsprozesse [MMP⁺10] versprechen schließlich Produkte und Anwendungen, die an kundenspezifische Wünsche und individuelle Nutzungsanforderungen angepasst werden können. Derartige Anforderungen sind insbesondere im Kontext der Industrie 4.0 verstärkt zu erwarten [Pla19]. Im Vergleich zu nicht-optimierungsbasierten adaptiven Regelungsansätzen (vgl. beispielsweise [ÄW95], [Tao03], [IS96]) steht bei den in der vorliegenden Arbeit betrachteten Methoden die Minimierung eines zugrunde liegenden Gütefunktional bei der Adaption von Reglergewichten im Zentrum (vgl. [KHL⁺12], [LVV12], [VLV13, S. 4]). Dabei können gewünschte Designziele, wie beispielsweise Kosten, Präzision, Zeitverhalten, Energieverbrauch, Komfort, Sicherheit und individuelle Bedürfnisse, durch Gütemaße berücksichtigt werden [LVS12], [AM89], [Ber95].

Reglerparameter basierend auf Simulations- oder Messdaten sowie einem Belohnungssignal zu adaptieren, um ein optimales Regelgesetz zu erlernen, entspricht dem Grundprinzip des Reinforcement Learning (RL) [SB18]. Mithilfe von RL ist es in den letzten Jahren gelungen, tiefe neuronale Netze zu trainieren und komplexe Aufgaben zu lösen. So konnten Meister des Brettspiels Go erstmalig durch einen Computer besiegt [SSS⁺17], [SHM⁺16], Greifbewegungen für Roboter unter Nutzung von Kameradaten gelernt [KIP⁺18] und Bewegungen einer robotischen Hand erlernt werden [ABC⁺20]. Jedoch musste für diese Erfolge ein erheblicher Trainingsaufwand betrieben werden, insbesondere wurden sehr große Mengen an Trainingsdaten benötigt. Beispielsweise spielte der Go-Computer ALPHAGO ZERO [SSS⁺17] im Laufe des Trainingsprozesses fast fünf Millionen komplette Partien gegen sich selbst, Kalashnikov et al. [KIP⁺18] sammelten Trainingsdaten aus 580 000 Greifversuchen über einen Zeitraum von vier Monaten und betrieben dazu gleichzeitig sieben Roboterarme, während Andrychowicz et al. [ABC⁺20] mithilfe einer Simulationsumgebung künstlich generierte Trainingsdaten nutzten, die, abhängig von der Variabilität der Umgebung, drei bis 100 Jahren an realen Erfahrungsdaten entsprachen.

Aus einer anwendungsorientierten und regelungstechnischen Perspektive sind solche außerordentlichen Datenmengen zum Training in bestimmten Anwendungen jedoch nicht verfügbar oder unpraktikabel [DAMH19], [BBdE10, S. 8]. Beispiele sind Messdaten realer technischer Systeme, Mensch-Maschine-Systeme oder biomedizinische Anwendungen (vgl. [LL04], [HMT⁺21], [GSK19]). Aus diesem Grund wird in der vorliegenden Arbeit das Konzept des sogenannten Adaptive Dynamic Programming (ADP)¹ [Wer92], [BT96], [MCLS02], [LV09], [LVV12] betrachtet, das eine Kombination aus adaptiver und optimaler Regelung darstellt [KHL⁺12] und Mechanismen des Reinforcement Learning nutzt. Die Einbeziehung gegebenenfalls vorhandenen Vorwissens über das System und die zugrunde liegende Problemstellung, wie beispielsweise die Systemordnung oder die Struktur der gesuchten Lösung, kann hierbei die Komplexität des Trainingsvorgangs und somit die Menge der benötigten Trainingsdaten entscheidend reduzieren (vgl. [Gör17], [BBdE10, Abschnitt 3.7.3 und Abschnitt 5.4]). ADP stellt zudem im Gegensatz zum klassischen, indirekten Vorgehen, zunächst ein Systemmodell zu identifizieren und anschließend eine modellbasierte Optimierung vorzunehmen, einen ganzheitlichen (vgl. [SB18, S. 3]), direkten Ansatz dar [LV09, S. 41], [SBW92], [JJ12], [FWS⁺18].

Dieser regelungstechnischen und anwendungsorientierten Perspektive folgend werden in der vorliegenden Arbeit zwei zentrale und bislang nur unzureichend gelöste Probleme ADP-basierter Regelungsansätze behandelt:

1. ADP-basierte Solltrajektorienfolgeregerler und
2. Konvergenzbedingungen für eine erfolgreiche Adaption.

Die Relevanz der ersten Problemstellung resultiert dabei aus der Anforderung vieler technischer Anwendungen, wie beispielsweise Fahrzeugen, Robotern oder verfahrenstechnischen Anlagen, dass Systemgrößen einer vorgegebenen Referenz- bzw. Solltrajektorie folgen sollen (vgl. z. B. [BK18], [FOSH14], [SCN⁺04], [LLHW16], [MBTL12], [DCP96]). Viele Literaturbeiträge zu ADP-basierten Methoden beschränken sich jedoch entweder auf den Fall, den Systemzustand optimal bezüglich eines Gütemaßes zu null zu regeln, oder betrachten Referenztrajektorien unter sehr einschränkenden Annahmen. In Anwendungen, wie beispielsweise der Robotik oder dem hochautomatisierten oder autonomen Fahren, die im Allgemeinen eine möglichst flexible und zur Laufzeit veränderbare Solltrajektorienvorgabe erfordern (vgl. [van97]), sind bestehende ADP-Ansätze daher ungeeignet. Eine wesentliche Herausforderung beim Entwurf von ADP-Zustandstrajektorienfolgereglern ist eine Repräsentation der Solltrajektorie, die in den ADP-Formalismus integriert werden kann. Eine hierfür geeignete Darstellung des Sollzustandsverlaufs, die einerseits eine flexible Solltrajektorienvorgabe ermöglicht und andererseits mit einem vertretbaren Lernaufwand verbunden ist, ist dabei nicht trivial, speziell aus einer anwendungsorientierten Perspektive jedoch unverzichtbar. Bislang existieren

¹ Weitere Bezeichnungen sind *Approximate Dynamic Programming* [Wer13], *Neuro-Dynamic Programming* [BT96], *Heuristic Dynamic Programming* [Wer77] oder *Incremental Dynamic Programming* [Wat89]. In der Literatur werden die Begriffe ADP und RL gelegentlich synonym verwendet [Gos09]. Eine präzise Abgrenzung ist im Allgemeinen nicht möglich [Wer13], die Bezeichnung ADP deutet jedoch meist auf einen eher regelungstechnisch orientierten Zweig des RL hin (vgl. [Gör17], [BBT⁺18], [LV09]).

somit weder eine allgemeine theoretische Definition, unter welchen Bedingungen eine Solltrajektorienendarstellung kompatibel für die Verwendung im ADP-Kontext ist, noch geeignete Ansätze ADP-basierter Solltrajektorienfolgeregler, welche die Einbeziehung einer von außen vorgebbaren Repräsentation variabler Solltrajektorienverläufe ermöglichen. Die vorliegende Arbeit schließt diese Lücke. Neben der theoretischen Analyse, Simulation und Diskussion der präsentierten Ansätze wird zudem die Anwendbarkeit von ADP-basierten Solltrajektorienfolgereglern anhand zweier realer Anwendungsbeispiele untersucht.

Ein zweiter wesentlicher Aspekt ADP-basierter Regler ist die Frage, unter welchen Bedingungen die Konvergenz des Lernprozesses gewährleistet ist. Diese Frage hängt eng mit dem Begriff der Exploration aus dem Bereich des RL zusammen. Dabei muss stets ein Kompromiss zwischen dem Ausprobieren neuer Stellgrößen und Zustände, die möglicherweise bislang unbekannte, bessere Optionen zutage bringen, und dem Ausnutzen bisheriger Erfahrungen eingegangen werden (Exploration-Exploitation-Dilemma) [SB18, S. 3]. Anschaulich betrachtet müssen die Trainingsdaten kausale Zusammenhänge, wie das Systemverhalten, den Einfluss potenzieller weiterer Akteure und die resultierenden Kosten, angemessen abbilden. Für eine erfolgreiche Konvergenz von ADP-Reglern ist somit eine ausreichende Anregung des Systems bzw. der für die Adaption relevanten Signale erforderlich. Speziell im Kontext des ADP existieren bislang jedoch kaum theoretische Erkenntnisse, wie eine geeignete Anregung entworfen werden kann [JKBD15] und welche Bedingungen die Systemzustände erfüllen müssen, um Konvergenz des Adaptionsprozesses zu erzielen. Der zweite zentrale Beitrag der vorliegenden Arbeit ist daher durch hinreichende Bedingungen an den Systemzustand gegeben, die gewährleisten, dass die für die Konvergenz benötigte Anregung des Systems erfüllt ist. Die theoretischen Ergebnisse werden anhand von Simulationen analysiert und diskutiert.

Gliederung der Arbeit

Die vorliegende Arbeit ist wie folgt gegliedert: **Kapitel 2** gibt einen Überblick über den relevanten Stand der Wissenschaft und führt benötigte Notationen und Begriffe ein. Dabei werden Forschungslücken herausgearbeitet, Forschungsfragen konkretisiert und ein Überblick über die Beiträge der vorliegenden Arbeit gegeben. Die Kapitel 3–6 bilden anschließend den Kern der Arbeit. Eine schematische Übersicht dieser Hauptinhalte ist in Abbildung 1.1 dargestellt.

Dabei werden in **Kapitel 3** und **Kapitel 4** flexible Referenztrajektorien Darstellungen, die mit dem ADP-Formalismus kompatibel sind, vorgestellt. In zeitdiskreter Formulierung (Kapitel 3) werden einerseits parametrisierte Solltrajektorienverläufe für den Entwurf eines ADP-basierten Trajektorienfolgereglers entwickelt, andererseits wird eine Methode zur direkten Verwendung der Sollzustände auf einem endlichen Vorausschauhorizont vorgestellt. In zeitkontinuierlicher Darstellung (Kapitel 4) werden ebenfalls parametrisierte, ADP-kompatible Sollzustandsverläufe präsentiert und in einen ADP-Ansatz integriert. Die Eigenschaften der unterschiedlichen Methoden werden analysiert und diskutiert. Zudem werden jeweils Simulationsergebnisse präsentiert.

Kapitel 5 liefert formale theoretische Aussagen zur Anregung ADP-basierter Regelungsansätze. Für die verallgemeinerte Problemstellung eines zeitkontinuierlichen eingangsaffinen Nicht-Nullsummen-Differenzialspiels werden hinreichende Bedingungen an den Systemzustand präsentiert, welche die Erfüllung der Anregungsbedingung und somit Konvergenz der betrachteten ADP-Methode garantieren. Illustrative Simulationsergebnisse werden präsentiert und diskutiert.

In **Kapitel 6** werden adaptive optimale Solltrajektorienfolgeregerler anhand zweier realer Anwendungsbeispiele präsentiert. Zunächst wird ein adaptiver Geschwindigkeitsregler, der einem vorgegebenen Geschwindigkeitsprofil möglichst kostenoptimal folgen soll, in einem realen Fahrzeug angewandt. Dieser wird online, d. h. während sich das Fahrzeug auf der Teststrecke befindet, adaptiert. Schließlich wird als zweites reales Anwendungsbeispiel ein adaptiver optimaler Trajektorienfolgeregerler für ein Ball-auf-Platte-System vorgestellt, der anhand von aufgezeichneten Messdaten und ohne Verwendung eines konkreten Systemmodells erlernt wird.

Kapitel 7 fasst schließlich die Hauptkenntnisse der Arbeit zusammen und bewertet diese.

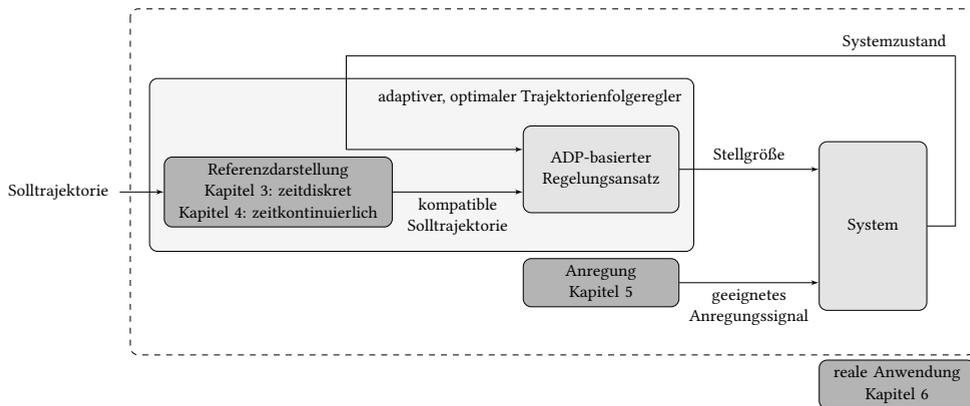


Abbildung 1.1: Schematische Übersicht der Hauptbeiträge der vorliegenden Arbeit.

2 Stand der Wissenschaft und Forschungslücke

In Abschnitt 2.1 dieses Kapitels werden zunächst zentrale Grundgleichungen und Grundbegriffe des ADP eingeführt, um eine Einordnung und Diskussion bestehender Methoden zu ermöglichen. Anschließend wird in Abschnitt 2.2 der Stand der Wissenschaft hinsichtlich bestehender ADP-basierter Solltrajektorienfolgeregelungsmethoden diskutiert. Weiterhin wird in Abschnitt 2.3 beleuchtet, wie die Anregung von Systemzuständen im ADP-Kontext in der Literatur bislang erfolgt. Nachdem die Forschungslücken herausgearbeitet wurden, werden in Abschnitt 2.4 die wissenschaftlichen Fragestellungen der Arbeit präzisiert und eine Übersicht über die Beiträge präsentiert.

2.1 Einführung der Notation und Grundbegriffe des ADP

ADP-basierte Regelungsansätze, die mit RL-Methoden und gemessenen Systemtrajektorien Optimalregelungsprobleme lösen sollen (vgl. [LV09, S. 39]), stellen ein vergleichsweise junges Gebiet der Regelungstechnik dar (vgl. [VLV13, S. 2], [SB18, Abschnitt 1.7], [Wer13]). Zwar postulieren Mendel und McLaren [MM70] bereits 1970 erste trial-and-error-basierte Regelungsansätze, die Ideen des RL aufgreifen, Werbos [Wer77] beschreibt 1977 das Konzept des *Heuristic Dynamic Programming* und Barto et al. [BSA83] gelangen durch Diskretisierung des Zustandsraums² simulative Trainingserfolge, dennoch beginnt die eigentliche Entwicklung moderner ADP-Ansätze erst mit der Arbeit von Watkins [Wat89] (vgl. [SB18, S. 14]). Eine ausführliche Übersicht über die Ursprünge und Grundlagen von RL- und ADP-Methoden ist beispielsweise in [LVV12], [SB18, Kapitel 1], [LV09], [WZL09], [LWW⁺17, Kapitel 1] und [Wer13] gegeben.

ADP- und RL-basierte Ansätze lassen sich in zwei wesentliche Klassen unterteilen. Neben Monte-Carlo-Methoden [SB18, Kapitel 5] spielen insbesondere auf dem sogenannten *Temporal-Difference-Fehler* (TD-Fehler) [Sut88], [SB18, Kapitel 6] beruhende Methoden eine Rolle. Letztere nutzen einen skalaren Prädiktionsfehler, der auf der Bellman-Gleichung bzw. der Hamilton-Jacobi-Bellman-Gleichung basiert, zur Adaption. Monte-Carlo-Methoden betrachten üblicherweise episodische Aufgaben (vgl. [SB18, S. 91]) mit klar definierten Endzuständen, welche viele Male wiederholt werden. Hierbei müssen Lernalgorithmen stets das Ende einer Episode abwarten, wohingegen Algorithmen, die den TD-Fehler nutzen, den Vorteil aufweisen, mit

² Dieser unter dem Begriff des *tile coding* bekannte Mechanismus kann jedoch zahlreiche theoretische und praktische Probleme verursachen (vgl. [van12, Abschnitt 2.1.2]).

jedem Zeitschritt ein Datentupel zu erhalten und unabhängig eines episodischen Charakters Adaptionen vornehmen zu können [LVV12, S. 87 f.], [SB18, Abschnitt 6.2], [VLV13, S. 29 f.]. Ein weiteres Argument, Anpassungen auf Basis eines skalaren Vorhersagefehlers, wie dem TD-Fehler, vorzunehmen, ist, dass dieses Grundprinzip des RL über den rein technischen Ursprung hinausgeht³. Im Folgenden wird der Fokus dieser Arbeit auf die Klasse der TD-Methoden gelegt.

Um den für die vorliegende Arbeit relevanten Stand der Wissenschaft angemessen einordnen zu können, werden in den nachfolgenden Abschnitten zunächst im ADP-Kontext wichtige Grundgleichungen in zeitdiskreter und zeitkontinuierlicher Darstellung eingeführt. Anschließend wird die Notwendigkeit der Verwendung von Funktionsapproximatoren bei wertkontinuierlichen Zustands- und Stellgrößenräumen diskutiert. Hierbei dient ein Funktionsapproximator beispielsweise der Beschreibung des funktionalen Zusammenhangs zwischen dem Systemzustand und den daraus resultierenden Kosten. Schließlich werden grundlegende ADP-Lösungsansätze vorgestellt.

2.1.1 Zeitdiskrete ADP-Grundgleichungen

In zeitdiskreter Darstellung werde zunächst ein System

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k) + \mathbf{g}(\mathbf{x}_k)\mathbf{u}_k, \quad (2.1)$$

$\mathbf{x} \in \mathbb{R}^n$, $\mathbf{u} \in \mathbb{R}^p$, $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times p}$, $\mathbf{f}(\mathbf{0}) = \mathbf{0}$, mit dem Zeitindex k und dem zu minimierenden⁴ Gütefunktional

$$J(\mathbf{x}_0, \boldsymbol{\mu}) = \sum_{\kappa=0}^{\infty} \gamma^{\kappa} (q(\mathbf{x}_{\kappa}) + \boldsymbol{\mu}(\mathbf{x}_{\kappa})^{\top} \mathbf{R} \boldsymbol{\mu}(\mathbf{x}_{\kappa})) =: \sum_{\kappa=0}^{\infty} \gamma^{\kappa} r(\mathbf{x}_{\kappa}, \boldsymbol{\mu}(\mathbf{x}_{\kappa})) \quad (2.2)$$

betrachtet (vgl. beispielsweise [ATLAK08], [LV09], [LVS12, Kapitel 11.5] und [WZL09]). Hierbei stellt $\boldsymbol{\mu} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ eine Zustandsrückführung dar. Weiterhin sei $q : \mathbb{R}^n \rightarrow \mathbb{R}$ eine positiv definite Funktion (vgl. [NA05, S. 53]) und \mathbf{R} symmetrisch und positiv definit. Zudem sei das

³ So können Parallelen zu verhaltenspsychologischen Modellen und der Neurobiologie gezogen werden. Bentham [Ben89] definiert bereits 1789 das Nützlichkeitsprinzip, nach welchem Handlungen bevorzugt werden, die den Gesamtnutzen vergrößern. 1911 spricht Thorndike [Tho11], das Lernverhalten von Tieren untersuchend, vom „Law of Effect“. Demnach werden Aktionen, die zu einer Belohnung geführt haben, in künftigen, ähnlichen Situationen häufiger wiederholt als mögliche Alternativen. Fortgesetzt wird dieser behavioristische Ansatz durch Watsons Grundsatz, Verhalten lasse sich durch Reiz und Reaktion beschreiben [Wat13]. Zudem findet laut Rescorla und Wagner [RW72], die auf Pavlovs klassischer Konditionierung [Pav27] aufbauen, gerade dann ein Lernprozess statt, wenn sich Prädiktion und tatsächliche Beobachtung unterscheiden. Somit ist ein Vorhersagefehler treibende Kraft des Lernens. Weiterhin lassen sich Parallelen zwischen dem TD-Learning und Lernvorgängen im menschlichen Gehirn, die auf dem Neurotransmitter Dopamin beruhen, ziehen. So lässt sich die Dopaminaktivität im Mittelhirn, ähnlich wie der TD-Fehler, als skalares Prädiktionsfehlersignal interpretieren, welches Lernvorgänge auf Neuronenebene beeinflusst [SDM97], [Gli11], [Niv09], [HDB94]. Letztlich lässt sich durch ADP-Methoden auch menschliches Verhalten beim Erlernen sensomotorischer Bewegungen nachbilden [JJ14a], [Bia17].

⁴ In dieser Arbeit wird die in der optimalen Regelung übliche Konvention der Minimierung von Kosten verwendet, während im Bereich des RL eine Maximierung von Belohnungen üblich ist (vgl. [LVV12]).

System stabilisierbar auf der kompakten Menge $\mathcal{X} \subset \mathbb{R}^n$, die den Ursprung enthält. Durch den sogenannten Diskontierungsfaktor γ mit $0 < \gamma \leq 1$ lässt sich parametrieren, ob die Einschrittkosten $r(\cdot)$ zu jedem Zeitschritt gleich gewichtet werden ($\gamma = 1$), oder inwiefern weiter in der Zukunft liegende Kosten schwächer ins Gewicht fallen ($\gamma < 1$). Weiterhin stellt der Diskontierungsfaktor ein Werkzeug dar, um auch bei nicht-verschwindenden Einschrittkosten $r(\cdot)$ einen endlichen Wert des Gütefunktional erhalten zu können (vgl. Abschnitt 2.2).

Eine wichtige Größe, die eng mit dem Gütefunktional $J(\mathbf{x}_0, \boldsymbol{\mu})$ in (2.2) verknüpft ist, stellt die sogenannte Zustands-Nutzenfunktion⁵, die auch als *Value Function* bezeichnet wird, dar. Die Value Function ist für zulässige⁶, d. h. stabilisierende, stetige⁷ Regler $\boldsymbol{\mu}(\mathbf{x}_k) \in \mathcal{C}^0(\mathcal{X})$, die zudem zu einem endlichen Gütefunktional führen (vgl. [ATLAK08, Definition 2]), durch

$$V^\mu(\mathbf{x}_k) = \sum_{\kappa=k}^{\infty} \gamma^{\kappa-k} r(\mathbf{x}_\kappa, \boldsymbol{\mu}(\mathbf{x}_\kappa)) = r(\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k)) + \gamma V^\mu(\mathbf{x}_{k+1}), \quad (2.3a)$$

$$V^\mu(\mathbf{0}) = 0, \quad (2.3b)$$

gegeben. Die durch (2.3b) beschriebene Anfangsbedingung folgt aufgrund von $\mathbf{f}(\mathbf{0}) = \mathbf{0}$ und $q(\mathbf{x}_k)$ positiv definit. In Übereinstimmung mit Bellmans Optimalitätsprinzip [Bel57a] resultiert für

$$\boldsymbol{\mu}^*(\mathbf{x}_k) = \arg \min_{\boldsymbol{\mu}} V^\mu(\mathbf{x}_k) \quad (2.4)$$

aus (2.3a) die Bellman-Gleichung

$$V^*(\mathbf{x}_k) := V^{\boldsymbol{\mu}^*}(\mathbf{x}_k) = \min_{\boldsymbol{\mu}(\mathbf{x}_k)} r(\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k)) + \gamma V^*(\mathbf{x}_{k+1}). \quad (2.5)$$

Das optimale Regelgesetz ergibt sich zu

$$\boldsymbol{\mu}^*(\mathbf{x}_k) = -\frac{\gamma}{2} \mathbf{R}^{-1} \mathbf{g}^\top(\mathbf{x}_k) \nabla_{\mathbf{x}_{k+1}} V^*(\mathbf{x}_{k+1}) \quad (2.6)$$

(vgl. [LVV12]). Neben der Value Function $V^\mu(\mathbf{x}_k)$ spielt die sogenannte Zustands-Aktions-Nutzenfunktion⁸, die in Anlehnung an klassisches Q-Learning [Wat89], [WD92] auch als *Q-Function*⁹ bezeichnet wird, eine wichtige Rolle in der ADP-Literatur (vgl. beispielsweise [KLM⁺14], [LLW⁺17], [LLHW16], [ATLAK07], [LVV12]). Die Q-Function ist durch

$$\begin{aligned} Q^\mu(\mathbf{x}_k, \mathbf{u}_k) &:= r(\mathbf{x}_k, \mathbf{u}_k) + \sum_{\kappa=k+1}^{\infty} \gamma^{\kappa-k} r(\mathbf{x}_\kappa, \boldsymbol{\mu}(\mathbf{x}_\kappa)) \\ &= r(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q^\mu(\mathbf{x}_{k+1}, \boldsymbol{\mu}(\mathbf{x}_{k+1})) \\ &= r(\mathbf{x}_k, \mathbf{u}_k) + \gamma V^\mu(\mathbf{x}_{k+1}) \end{aligned} \quad (2.7)$$

⁵ (engl.): *state value function*. Per Definition gilt $V^\mu(\mathbf{x}_0) = J(\mathbf{x}_0, \boldsymbol{\mu})$.

⁶ (engl.): *admissible*.

⁷ Die Notation $\mathcal{C}^l(\mathcal{X})$ kennzeichnet die Klasse der auf der Menge \mathcal{X} l -fach stetig differenzierbaren Funktionen.

⁸ (engl.): *state-action value function*.

⁹ Diese Bezeichnung geht auf den Begriff *quality function* zurück [LVV12].

definiert. Sie beschreibt die (diskontierten) Langzeitkosten, wenn sich das System im Zustand \mathbf{x}_k befindet, im Zeitschritt k die Stellgröße \mathbf{u}_k und in allen weiteren Zeitschritten das Regelgesetz $\boldsymbol{\mu}(\mathbf{x}_k)$ angewandt wird. Mit $Q^*(\mathbf{x}_k, \mathbf{u}_k) := Q^\mu(\mathbf{x}_k, \mathbf{u}_k)$ ergibt sich das optimale Regelgesetz zu

$$\boldsymbol{\mu}^*(\mathbf{x}_k) = \mathbf{u}^*(\mathbf{x}_k) = \arg \min_{\mathbf{u}_k} Q^*(\mathbf{x}_k, \mathbf{u}_k) \quad (2.8)$$

[LLW⁺17], [BBdE10, S. 18], [LVV12], zudem gilt

$$V^*(\mathbf{x}_k) = \min_{\mathbf{u}_k} Q^*(\mathbf{x}_k, \mathbf{u}_k). \quad (2.9)$$

Ein wesentlicher Unterschied von (2.8) im Vergleich zu (2.6) ist, dass die Eingangsdynamik $\mathbf{g}(\mathbf{x}_k)$ darin nicht präsent ist. Dies ist auf die explizite Abhängigkeit der Q-Function von \mathbf{u}_k zurückzuführen (vgl. [LVS12, S. 502]). Weiterhin hängt (2.8) nicht von \mathbf{x}_{k+1} , und somit implizit von der Stellgröße \mathbf{u}_k (vgl. [LVS12, S. 497]), sondern ausschließlich von \mathbf{x}_k ab.

2.1.2 Zeitkontinuierliche ADP-Grundgleichungen

Basierend auf [VL09], [VL10], [LVV12], [WHL17] und [LVS12, Kapitel 10.1] wird ein eingangsaflines System mit der Dynamik

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t)) + \mathbf{g}(\mathbf{x}(t))\mathbf{u}(t), \quad (2.10a)$$

$$\mathbf{x}(0) = \mathbf{x}_0, \quad (2.10b)$$

$\mathbf{x} \in \mathbb{R}^n$, $\mathbf{u} \in \mathbb{R}^p$, $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times p}$, \mathbf{f} und \mathbf{g} Lipschitz-stetig auf der kompakten Menge $\mathcal{X} \subset \mathbb{R}^n$, die den Ursprung enthält, $\mathbf{f}(\mathbf{0}) = \mathbf{0}$, (2.10) stabilisierbar auf \mathcal{X} , und ein zu minimierendes Gütefunktional der Form

$$J(\mathbf{x}_0, \boldsymbol{\mu}) = \int_0^\infty q(\mathbf{x}) + \boldsymbol{\mu}^\top(\mathbf{x})\mathbf{R}\boldsymbol{\mu}(\mathbf{x}) \, d\tau =: \int_0^\infty r(\mathbf{x}, \boldsymbol{\mu}) \, d\tau \quad (2.11)$$

betrachtet¹⁰. Dabei sei $q : \mathbb{R}^n \rightarrow \mathbb{R}$ abermals eine positiv definite Funktion und $\boldsymbol{\mu} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ eine Zustandsrückführung. Zudem gelte $\mathbf{R} = \mathbf{R}^\top \succ \mathbf{0}$. Diese Form des Systems (2.10) und Gütefunktionals (2.11) ist in der ADP-Literatur weitverbreitet und schließt insbesondere für $\mathbf{f}(\mathbf{x}(t)) = \mathbf{A}\mathbf{x}(t)$ und $\mathbf{g}(\mathbf{x}(t)) = \mathbf{B}$ sowie $q(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q}\mathbf{x}$ mit $\mathbf{Q} \succeq \mathbf{0}$ die Problemklasse der linear-quadratischen (LQ-)Optimierungsprobleme ein.

¹⁰ Ähnlich wie bei der zuvor in Abschnitt 2.1.1 vorgestellten zeitdiskreten Formulierung lassen sich auch die zeitkontinuierlichen Gleichungen auf ein diskontiertes Gütemaß verallgemeinern. Mit $\gamma_c \in \mathbb{R}_{\geq 0}$ wird die später in (2.12) definierte Value Function dann zu $V^\mu(\mathbf{x}(t)) = \int_t^\infty e^{-\gamma_c(\tau-t)} r(\mathbf{x}, \boldsymbol{\mu}) \, d\tau = \int_t^{t+T_{\text{IRL}}} e^{-\gamma_c(\tau-t)} r(\mathbf{x}, \boldsymbol{\mu}) \, d\tau + e^{-\gamma_c(T_{\text{IRL}}-t)} V^\mu(\mathbf{x}(t+T_{\text{IRL}}))$ und die Lyapunov-Gleichung nach (2.13) zu $0 = r(\mathbf{x}, \boldsymbol{\mu}) + (\nabla_{\mathbf{x}} V^\mu(\mathbf{x}))^\top (\mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\boldsymbol{\mu}(\mathbf{x})) - \gamma_c V^\mu(\mathbf{x})$, $V^\mu(\mathbf{0}) = 0$ (vgl. [Doy00]). Aus Gründen der Übersichtlichkeit ist in Abschnitt 2.1.2 jedoch der Fall $\gamma_c = 0$ dargestellt.

Die Value Function wird für zulässige Regelgesetze¹¹ [BSW97, Definition 1] durch

$$V^\mu(\mathbf{x}) := V^\mu(\mathbf{x}(t)) = \int_t^\infty r(\mathbf{x}, \boldsymbol{\mu}) \, d\tau \quad (2.12)$$

definiert [VPAKL09], [VL09]. Somit sind durch die Value Function $V^\mu(\mathbf{x}(t))$ die Gesamtkosten gegeben, die entstehen, wenn sich das System zum Zeitpunkt t im Zustand $\mathbf{x}(t)$ befindet und das Regelgesetz $\boldsymbol{\mu}(\mathbf{x})$ verwendet wird.

Die infinitesimale Version von $V^\mu(\mathbf{x})$ basierend auf (2.12) ist durch die sogenannte Lyapunov-Gleichung (vgl. [BSW97], [VL10])

$$0 = r(\mathbf{x}, \boldsymbol{\mu}) + (\nabla_{\mathbf{x}} V^\mu(\mathbf{x}))^\top (\mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\boldsymbol{\mu}(\mathbf{x})), \quad V^\mu(\mathbf{0}) = 0, \quad (2.13)$$

gegeben¹². Die mit dem Optimierungsproblem assoziierte Hamilton-Funktion¹³ lautet

$$H(\mathbf{x}, \nabla_{\mathbf{x}} V(\mathbf{x}), \boldsymbol{\mu}(\mathbf{x})) = r(\mathbf{x}, \boldsymbol{\mu}) + (\nabla_{\mathbf{x}} V(\mathbf{x}))^\top (\mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\boldsymbol{\mu}(\mathbf{x})) \quad (2.14)$$

(vgl. [LLW14], [WHL17]). Die optimale Value Function

$$V^*(\mathbf{x}) =: V^{\mu^*}(\mathbf{x}) = \min_{\boldsymbol{\mu}} V^\mu(\mathbf{x}) \quad (2.15)$$

erfüllt unter Annahme der Existenz des Minimums die Hamilton-Jacobi-Bellman-Gleichung (HJB-Gleichung)

$$0 = \min_{\boldsymbol{\mu}} H(\mathbf{x}, \nabla_{\mathbf{x}} V^*(\mathbf{x}), \boldsymbol{\mu}(\mathbf{x})). \quad (2.16)$$

Diese entspricht der Lyapunov-Gleichung (2.13) für die optimale Value Function $V^*(\mathbf{x})$ und das optimale Regelgesetz $\boldsymbol{\mu}^*$. Zudem gilt nach (2.13) für jedes zulässige Regelgesetz $\boldsymbol{\mu}(\mathbf{x})$ mit der zugehörigen Value Function $V^\mu(\mathbf{x})$ mit $V^\mu(\mathbf{0}) = 0$

$$0 = H(\mathbf{x}, \nabla_{\mathbf{x}} V^\mu(\mathbf{x}), \boldsymbol{\mu}(\mathbf{x})). \quad (2.17)$$

Unter der Annahme, dass das Minimum auf der rechten Seite von (2.16) existiert und eindeutig ist, ergibt sich das optimale Regelgesetz zu

$$\boldsymbol{\mu}^*(\mathbf{x}) = \arg \min_{\boldsymbol{\mu}(\mathbf{x})} H(\mathbf{x}, \nabla_{\mathbf{x}} V^*(\mathbf{x}), \boldsymbol{\mu}(\mathbf{x})) = -\frac{1}{2} \mathbf{R}^{-1} \mathbf{g}^\top(\mathbf{x}) \nabla_{\mathbf{x}} V^*(\mathbf{x}) \quad (2.18)$$

¹¹ (engl.): *admissible policies*. Diese sind durch $\boldsymbol{\mu}(\mathbf{0}) = \mathbf{0}$, $\boldsymbol{\mu}(\mathbf{x})$ stetig und stabilisierend auf \mathcal{X} und $V^\mu(\mathbf{x})$ endlich definiert. Letzteres wird hier zusätzlich zu einem stabilisierenden Regelgesetz gefordert, da ein stabilisierender Regler allein nicht gewährleistet, dass das Kostenintegral endlich bleibt [BSW97, S. 2162].

¹² Da im Fall des hier betrachteten unendlichen Optimierungshorizonts keine explizite Abhängigkeit der Value Function von der Zeit existiert, gilt $\frac{\partial V^\mu(\mathbf{x})}{\partial t} = 0$. Zudem wird angenommen, dass $V^\mu(\mathbf{x})$ stetig differenzierbar ist [VL10], d. h. $V^\mu(\mathbf{x}) \in C^1(\mathcal{X})$.

¹³ Hierbei ist $H(\cdot)$ als Funktion mit beliebigem $V(\mathbf{x})$ und $\boldsymbol{\mu}(\mathbf{x})$ zu verstehen und es muss im Allgemeinen nicht $V(\mathbf{x}) = V^\mu(\mathbf{x})$ gelten.

(vgl. [VL09], [VL10]). Das Einsetzen von (2.18) in (2.16) und die Verwendung von (2.11) liefern die HJB-Gleichung bezüglich $\nabla_{\mathbf{x}} V^*$ (vgl. [VL10]):

$$0 = q(\mathbf{x}) + \nabla_{\mathbf{x}} V^{*\top} \mathbf{f}(\mathbf{x}) - \frac{1}{4} \nabla_{\mathbf{x}} V^{*\top}(\mathbf{x}) \mathbf{g}(\mathbf{x}) \mathbf{R}^{-1} \mathbf{g}^\top(\mathbf{x}) \nabla_{\mathbf{x}} V^*(\mathbf{x}), \quad V^*(\mathbf{x}) = 0. \quad (2.19)$$

Falls es nun gelingt, V^* zu finden, sodass (2.19) erfüllt ist, ergibt sich nach (2.18) das optimale Regelgesetz $\boldsymbol{\mu}^*(\mathbf{x})$.

Eine Alternative zur direkten Verwendung der Lyapunov-Gleichung (2.13) ergibt sich durch Umformulierung von (2.12) zu

$$V^\mu(\mathbf{x}) := V^\mu(\mathbf{x}(t)) = \int_t^\infty r(\mathbf{x}, \boldsymbol{\mu}) \, d\tau \quad (2.20a)$$

$$= \int_t^{t+T_{\text{IRL}}} r(\mathbf{x}, \boldsymbol{\mu}) \, d\tau + V^\mu(\mathbf{x}(t + T_{\text{IRL}})) \quad (2.20b)$$

(vgl. [VPAKL09], [VL09]), wobei $T_{\text{IRL}} > 0$ eine beliebige Intervalllänge darstellt. Gleichung (2.20b) spielt eine zentrale Rolle für ADP-Methoden, die sich unter dem Begriff des *Integral Reinforcement Learning* (IRL) zusammenfassen lassen (vgl. beispielsweise [LVV12], [SLW17], [JJ14c], [BJ16a], [VPAKL09], [VL09], [ML14b], [LS17], [LPC15]). Ein IRL-basierter Ansatz [BJ16a] wird in Kapitel 4 der vorliegenden Arbeit genutzt, während die Lyapunov-Gleichung (2.13) in Kapitel 5 Verwendung findet.

Das Lösen der HJB-Gleichung (2.16) bzw. (2.19) ist im Allgemeinen schwierig [BSW97], [VL09], [VL10], zudem wird die vollständige Kenntnis der Systemdynamik in Form von $\mathbf{f}(\mathbf{x})$ und $\mathbf{g}(\mathbf{x})$ benötigt. Ebenso lässt sich aus (2.20b) nicht unmittelbar auf V^* schließen. An dieser Stelle setzen ADP-Methoden an, um basierend auf gemessenen Zustands- und Stellgrößentrjektorien sowie dem Kostensignal $r(\cdot)$ Regelgesetze im Sinne des Gütemaßes J nach (2.11) zu adaptieren. Die wichtigsten Lösungsansätze hierzu werden in Abschnitt 2.1.4 überblicksartig vorgestellt.

2.1.3 Funktionsapproximatoren

Die in dieser Arbeit betrachteten regelungstechnischen Problemstellungen weisen wertkontinuierliche Zustands- und Stellgrößenräume auf. Daher sind klassische RL-Algorithmen ungeeignet, die Markov-Entscheidungsprozesse mit endlichen Zustands- und Stellgrößenräumen betrachten [KLM96] und mit tabellarischen Darstellungen der Value Function V bzw. Q-Function Q arbeiten (z. B. Q-Learning [Wat89], klassisches Temporal-Difference Learning [Sut88], tabellarische Policy-Iteration- [SB18, Kapitel 4.3] oder Value-Iteration-Algorithmen [SB18, Kapitel 4.4]). Insbesondere bei wertkontinuierlichen Zustands- und Stellgrößenräumen ist daher die Verwendung von geeigneten Funktionsapproximatoren entscheidend (vgl. beispielsweise [BBdE10, Kapitel 3.2], [LVV12, S. 90], [van12], [BBT⁺18], [LP03], [LV09], [LLHW16]). Unabhängig davon, ob zeitkontinuierliche oder zeitdiskrete Systeme betrachtet werden, wird eine Value Function $V(\mathbf{x})$ durch h stetig differenzierbare Basisfunktionen

$\phi(\mathbf{x}) \in \mathcal{C}^1$ bzw. eine Q-Function $Q(\mathbf{x}, \mathbf{u})$ mithilfe h stetig differenzierbarer Basisfunktionen $\phi(\mathbf{x}, \mathbf{u}) \in \mathcal{C}^1$ approximiert. Zusammenfassen der jeweiligen Basisfunktionen ergibt die Vektoren $\phi(\mathbf{x}) \in \mathbb{R}^h$ bzw. $\phi(\mathbf{x}, \mathbf{u}) \in \mathbb{R}^h$. Am Beispiel der Value Function $V^\mu(\mathbf{x})$ resultiert bei Verwendung linearer Funktionsapproximatoren¹⁴

$$V^\mu(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + \epsilon(\mathbf{x}) \quad (2.21)$$

mit dem Gewichtsvektor $\mathbf{w} \in \mathbb{R}^h$ und dem Approximationsfehler $\epsilon(\mathbf{x})$ (vgl. [VL10], [LVV12], [VL09], [WHL17])¹⁵. Der Funktionsapproximator einer Value Function oder Q-Function wird in der Literatur als *Critic* bezeichnet, da dieser eine Bewertung eines Regelgesetzes vornimmt (vgl. [BSA83], [LVV12]). Demgegenüber wird ein Funktionsapproximator, der ein Regelgesetz beschreibt, als *Actor* bezeichnet. Zahlreiche ADP-Methoden verwenden lediglich einen Critic (z. B. [WLL14], [LYWW15], [LLW14], [VL09]), beispielsweise, wenn ein analytischer Zusammenhang zwischen der Value Function V oder Q-Function Q und dem (geschätzten) optimalen Regelgesetz genutzt werden kann (vgl. (2.6) und (2.18)). Verfahren, die sowohl V bzw. Q als auch das Regelgesetz μ durch einen eigenen Funktionsapproximator parametrieren, werden als *Actor-Critic-Methoden* bezeichnet (z. B. [VL10], [BKJ⁺13], [LW14], [KL15])¹⁶. Zusammenfassend wird durch Funktionsapproximatoren die Verwendung von wertkontinuierlichen Zustands- und Stellgrößenräumen ermöglicht. Die Suche nach einer optimalen Value- bzw. Q-Function und einem optimalen Regelgesetz wird somit auf eine Parametersuche übertragen.

Schließlich sei noch anzumerken, dass die Formulierung des zu minimierenden Gütemaßes J über einen unendlichen Zeithorizont eine zentrale Bedeutung für die erfolgreiche Verwendung von Funktionsapproximatoren einnimmt, wie die nachfolgende Bemerkung konstatiert.

¹⁴ Wenngleich grundsätzlich auch nichtlineare Funktionsapproximatoren verwendet werden können, so sind im Allgemeinen lineare Funktionsapproximatoren, d. h. eine linear gewichtete Summation ggf. nichtlinearer Basisfunktionen auf theoretischer Ebene besser verstanden und zudem aus praktischer Sicht einfach handhabbar (vgl. [van12], [BT96], [BBT⁺18], [BBdE10, Kapitel 3.3.1]).

¹⁵ Unter der Annahme, dass die Value Function $V^\mu(\mathbf{x})$ auf einer kompakten Menge \mathcal{X} stetig differenzierbar ist, lässt sich $\forall \mathbf{x} \in \mathcal{X}$ der Approximationsfehler $\epsilon(\mathbf{x})$ bei geeignet gewählten Basisfunktionen $\phi(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^h$ mit steigender Anzahl h an Basisfunktionen beliebig verringern. Für polynomielle Basisfunktionen basiert dies beispielsweise auf dem Satz von Weierstraß [Wei85], jedoch existieren auch Verallgemeinerungen (vgl. [HSW90], [BSW97], [AKL05], [VL10]). Da eine sehr große Anzahl an Basisfunktionen und somit zu lernenden Gewichten vermieden werden soll, sei jedoch anzumerken, dass ohne Vorwissen über die Systemdynamik eine praktikable Wahl an Basisfunktionen bei nichtlinearen Systemen grundsätzlich eine bislang ungelöste Herausforderung darstellt (vgl. [BKJ⁺13, Remark 5]).

¹⁶ In der Literatur wird häufig von einer *Actor-Critic-Struktur* gesprochen, selbst, wenn kein expliziter, zusätzlicher Funktionsapproximator zur Beschreibung der Stellgröße verwendet wird (vgl. [LVV12]). Dies ist darin begründet, dass der Regler μ , der bei allen ADP-Methoden vorhanden ist, in Anlehnung an RL-Methoden als *Actor* bezeichnet werden kann. Zur Abgrenzung wird in der vorliegenden Arbeit der Oberbegriff *verallgemeinerte Critic-Struktur* definiert, um sowohl Methoden mit zusätzlichem, gesondertem Actor-Funktionsapproximator, als auch solche, die direkt aus dem Critic-Gewicht das Regelgesetz μ bestimmen, zu bezeichnen. Demgegenüber wird explizit dann von *Actor-Critic-Methoden* gesprochen, wenn neben dem Critic-Approximator ein zusätzlicher Funktionsapproximator für das Regelgesetz μ verwendet wird.

Bemerkung 2.1

Für den Entwurf von ADP-Methoden mithilfe von Funktionsapproximatoren ist das Auftauchen derselben Value Function $V^\mu(\cdot)$ bzw. Q-Function $Q^\mu(\cdot)$ auf beiden Seiten der zentralen Gleichungen (2.3a), (2.7) und (2.20b) fundamental. Für diese Eigenschaft ist nach [LVV12, S. 80] maßgeblich die Verwendung eines unendlichen Optimierungshorizontes in (2.2) bzw. (2.11) verantwortlich.

2.1.4 ADP-Lösungsansätze

An dieser Stelle lässt sich der Begriff des ADP weiter schärfen. Im Folgenden stellt ADP in Anlehnung an [LWW⁺17] und [WHL17] eine Kombination aus dynamischer Programmierung, der Verwendung von Funktionsapproximatoren und einer verallgemeinerten Critic-Struktur¹⁷ dar, um mithilfe von Systemtrajektorien Optimalregelungsprobleme zu lösen (vgl. [LV09, S. 39]).

Hierfür werden in diesem Abschnitt grundlegende ADP-Mechanismen eingeführt, die im weiteren Verlauf der vorliegenden Arbeit benötigt werden. Neben der Anpassung des Regelgesetzes $\mu(x)$, entweder direkt aus der Value Function V bzw. der Q-Function Q oder über die Anpassung von Actor-Gewichten (vgl. Abschnitt 2.1.3), nimmt insbesondere die Adaption der Critic-Gewichte w eine zentrale Rolle ein. Wird bei Anwendung eines Regelgesetzes μ die aktuelle Schätzung des Critic-Gewichts angepasst, um die mit diesem Regelgesetz verbundenen Langzeitkosten in Form von V^μ bzw. Q^μ zu beschreiben, so wird dies als *Policy-Evaluation-Schritt* bezeichnet. Eine Anpassung des Regelgesetzes μ basierend auf der aktuellen Schätzung von V^μ bzw. Q^μ , mit dem Ziel, ein bezüglich des Gütefunktional J verbessertes Regelgesetz zu erhalten, wird hingegen *Policy-Improvement-Schritt* genannt [LVV12]. Je nachdem, ob diese Schritte alternierend oder parallel durchgeführt und ob sie vollständig oder nur teilweise ausgeführt werden, lassen sich ADP-Methoden klassifizieren. Die drei für die vorliegende Arbeit wesentlichen Klassen sind *Policy-Iteration-Algorithmen*, *Value-Iteration-Algorithmen* und *Actor-Critic-Methoden*. Diese werden in den Abschnitten 2.1.4.1–2.1.4.3 thematisiert. Alle drei Klassen sind der sehr allgemeinen Definition der sogenannten *Generalized-Policy-Iteration-Methoden* zuzuordnen, die, unabhängig der Details und Granularität, auf einer irgendwie gearteten Interaktion zwischen Policy-Evaluation- und Policy-Improvement-Prozessen beruhen (vgl. [SB18, Abschnitt 4.6]). Nach der Betrachtung dieser drei zentralen ADP-Klassen wird in Abschnitt 2.1.4.4 schließlich der Unterschied zwischen sogenannten *On-Policy*- und *Off-Policy-ADP-Ansätzen* charakterisiert. Abbildung 2.1 gibt eine Übersicht über die nachfolgend präsentierten ADP-Klassen und deren wichtigste Eigenschaften.

¹⁷ Dies schließt sowohl Actor-Critic-Methoden als auch Ansätze ohne expliziten, zusätzlichen Actor-Funktionsapproximator ein (vgl. Abschnitt 2.1.3). Ansätze, die versuchen, ohne einen Critic-Approximator das Regelgesetz anzupassen (siehe beispielsweise [RPS07], [Mun06]), stehen hingegen nicht im Fokus der vorliegenden Arbeit.

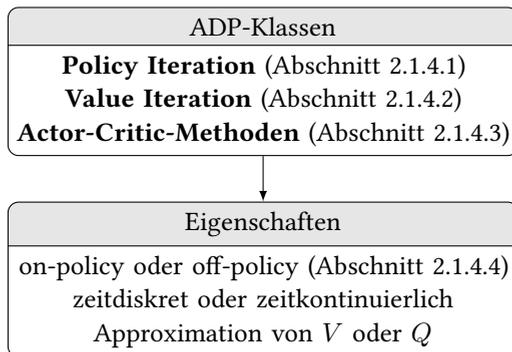


Abbildung 2.1: Klassifikation von ADP-Methoden.

2.1.4.1 Policy Iteration

Policy-Iteration-Algorithmen (PI-Algorithmen) basieren darauf, Regelgesetze auszuwerten, indem deren zugehörige Value Functions bzw. Q-Functions ermittelt werden. Diese Value- bzw. Q-Functions werden dann dazu verwendet, um neue, bezüglich des Gütefunktional J verbesserte, Regelgesetze zu finden. Dieser Prozess findet iterativ statt [BBdE10, Kapitel 2.4], [SB18, Kapitel 4.3]. PI-Algorithmen sind somit dadurch gekennzeichnet, dass der Policy-Evaluation-Schritt vollständig durchgeführt wird, also stets die Value Function V^μ (bzw. Q-Function Q^μ) zu einem Regelgesetz μ bestimmt wird, bevor das Regelgesetz im Policy-Improvement-Schritt angepasst wird. Am Beispiel einer zeitdiskreten Problemstellung, bei der eine Value Function V^μ approximiert wird, basiert diese iterative Prozedur auf (2.3a) und (2.6). Der prinzipielle Ablauf ist in Algorithmus 2.1 skizziert (vgl. [LVV12], [LV09]). Diese Iteration kann entweder fortgeführt oder durch ein Abbruchkriterium, beispielsweise bei Konvergenz von $V^{[l]}$, beendet werden. Die praktische Auswertung der Schritte 1 und 2 hängt letztlich von der konkreten Problemstellung und den gewählten Lösungsansätzen ab. Für den Policy-Evaluation-Schritt kann beispielsweise ein Gradientenabstieg, ein Batch-Least-Squares- oder ein rekursiver Least-Squares-Ansatz genutzt werden [LVV12]. Zur Lösung des Policy-Improvement-Schrittes kann z. B. direkt die analytische Lösung [LVV12] oder ein Gradientenabstiegsverfahren [LW14] verwendet werden.

Policy-Iteration-Algorithmen existieren neben der in Algorithmus 2.1 gezeigten zeitdiskreten Formulierung, die eine Value Function V lernt (vgl. [LVV12], [LV09], [Hey16]), in zahlreichen Varianten. Beispielsweise kann basierend auf (2.7) und (2.8) für eine zeitdiskrete Systemdarstellung eine Policy-Iteration definiert werden, die eine Q-Function lernt [LP03], [ATLAK07], [BYB94], [LV09]. Erste Ansätze dezentraler und verteilter Q-Learning-Methoden, die ebenfalls auf einer PI basieren, wurden von Görjes [Gör19] präsentiert. PI-Algorithmen zur Lösung zeitkontinuierlicher Problemstellungen, die eine Value Function V lernen, müssen entweder im Policy-Evaluation-Schritt die durch (2.13) gegebene Lyapunov-Gleichung lösen [BSW97] oder nutzen die in (2.20b) gegebene IRL-Darstellung [VPAKL09], [LVV12, S. 97]. Analog zu

Algorithmus 2.1 Zeitdiskrete Policy Iteration

1: **Initialisiere** Iterationsindex $l := 0$, zulässiges initiales Regelgesetz $\boldsymbol{\mu}^{[0]}(\mathbf{x}_k)$

Schritt 1 (Policy Evaluation):

2: Finde $V^{[l+1]}$, sodass gilt:

$$V^{[l+1]}(\mathbf{x}_k) = r(\mathbf{x}_k, \boldsymbol{\mu}^{[l]}(\mathbf{x}_k)) + \gamma V^{[l+1]}(\mathbf{x}_{k+1}). \quad (2.22)$$

Schritt 2 (Policy Improvement):

3: Aktualisiere das Regelgesetz

$$\boldsymbol{\mu}^{[l+1]}(\mathbf{x}_k) = \arg \min_{\boldsymbol{\mu}(\cdot)} \left(r(\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k)) + \gamma V^{[l+1]}(\mathbf{x}_{k+1}) \right). \quad (2.23)$$

4: Setze $l := l + 1$. Gehe zu Schritt 1.

IRL-Methoden [VPAKL09] kann eine PI auch so formuliert werden, dass eine Q-Function einer zeitkontinuierlichen Problemstellung iterativ gelernt wird [LPC12].

Die Forderung eines initial zulässigen Regelgesetzes $\boldsymbol{\mu}^{[0]}(\cdot)$ bei PI-Algorithmen kann entweder bei vorhandenem Teilwissen über die Systemdynamik mittels robuster Regelungsansätze erfüllt werden [ML14a, Remark 11], [KLNSK15, Remark 9], [YDZY20] oder vor Anwendung des PI-Algorithmus überprüft werden (vgl. [LW14, Theorem 3.3]). Wenngleich die Forderung nach einem initial zulässigen Regelgesetz eine gewisse Einschränkung¹⁸ von PI-Algorithmen darstellt, so resultiert für nicht-diskontierte Problemstellungen (d. h. für $\gamma = 1$) nach jedem Policy-Improvement-Schritt erneut ein zulässiges Regelgesetz [BSW97, Lemma 9], [LW14]. Auch die Verwendung einer durch einen Funktionsapproximator nicht exakt beschriebenen Value Function V führt in diesem Fall zu einem zulässigen Regelgesetz, wenn die Funktionsapproximation hinreichend genau erfolgt [BSW97, Theorem 26]. Insbesondere ist somit eine exakte Lösung des Policy-Evaluation-Schrittes nicht notwendig, um Stabilität des geschlossenen Regelkreises zu erreichen.

2.1.4.2 Value Iteration

Value-Iteration-Algorithmen (VI-Algorithmen) suchen basierend auf der Bellman-Gleichung oder der HJB-Gleichung iterativ eine optimale Value Function V^* bzw. Q-Function Q^* sowie den optimalen Regler $\boldsymbol{\mu}^*$, führen dabei jedoch keine vollständige Policy Evaluation durch [SB18, Kapitel 4.4], [BBdE10, Kapitel 2.3]. Eine zeitdiskrete Value Iteration, die eine Value Function V lernt, ist in Algorithmus 2.2 gegeben (vgl. [LV09]). Die Schritte *Policy Evaluation*

¹⁸ Hinsichtlich realer Anwendungen erscheint diese Forderung jedoch wenig einschränkend, da ohnehin die Beschränktheit der Systemzustände während der Datenakquise gewährleistet sein muss und ein Betrieb eines Realsystems mit einem instabilen Regelgesetz zu vermeiden ist.

und *Policy Improvement* lassen sich zu einer Gleichung kombinieren (vgl. [LVV12, (30)], [Ber17, (5)]), sodass

$$V^{[l+1]}(\mathbf{x}_k) = \min_{\boldsymbol{\mu}(\cdot)} \left(r(\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k)) + \gamma V^{[l]}(\mathbf{x}_{k+1}) \right) \quad (2.24)$$

resultiert. Im Unterschied zur PI muss das initiale Regelgesetz hierbei nicht zulässig (und damit insbesondere nicht stabilisierend) sein und im Gegensatz zu (2.22) wird auf der rechten Seite von (2.25) die vorherige Schätzung $V^{[l]}(\cdot)$ der Value Function verwendet. Das bedeutet, dass im Policy-Evaluation-Schritt nicht die zum aktuellen Regelgesetz $\boldsymbol{\mu}^{[l]}(\mathbf{x}_k)$ gehörende Value Function ermittelt wird, sondern lediglich ein einzelner Update-Schritt in Richtung dieser Value Function¹⁹ vorgenommen wird. Auch hier existieren Formulierungen, die für zeitdiskrete Probleme eine Q-Funktion [LVV12, S. 95] oder in zeitkontinuierlichen IRL-Ansätzen eine Value Function [LVV12, S. 97] lernen.

VI-Algorithmen weisen den Vorteil auf, prinzipiell nicht mit einem zulässigen Regelgesetz initialisiert werden zu müssen. Zudem erfordert eine einzelne Iteration einer VI, abhängig von der konkreten Implementierung, ggf. einen geringeren Berechnungsaufwand als eine Iteration der PI [KLM96, S. 250]²⁰. Jedoch kann bei VI-basierten Verfahren, im Gegensatz zu PI-Algorithmen, im Allgemeinen keine Aussage über die Stabilität der Regler $\boldsymbol{\mu}^{[l+1]}$ während der

Algorithmus 2.2 Zeitdiskrete Value Iteration

1: **Initialisiere** Iterationsindex $l := 0$, initiales Regelgesetz $\boldsymbol{\mu}^{[0]}(\mathbf{x}_k)$, $V^{[0]} \geq 0$

Schritt 1 (Policy Evaluation):

2: Finde $V^{[l+1]}$, sodass gilt:

$$V^{[l+1]}(\mathbf{x}_k) = r(\mathbf{x}_k, \boldsymbol{\mu}^{[l]}(\mathbf{x}_k)) + \gamma V^{[l]}(\mathbf{x}_{k+1}). \quad (2.25)$$

Schritt 2 (Policy Improvement):

3: Aktualisiere das Regelgesetz

$$\boldsymbol{\mu}^{[l+1]}(\mathbf{x}_k) = \arg \min_{\boldsymbol{\mu}(\cdot)} \left(r(\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k)) + \gamma V^{[l+1]}(\mathbf{x}_{k+1}) \right). \quad (2.26)$$

4: Setze $l := l + 1$. Gehe zu Schritt 1.

¹⁹ Diese Aussage wird insbesondere dadurch gestützt, dass die Bellman-Gleichung eine Fixpunktgleichung ist und die durch $V^{[i+1]}(\mathbf{x}_k) = r(\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k)) + \gamma V^{[i]}(\mathbf{x}_{k+1})$ definierte Iteration mit fixiertem Regelgesetz $\boldsymbol{\mu}(\mathbf{x}_k)$ eine Kontraktion darstellt (vgl. [LV09], [Ber20b]). Für ein zulässiges $\boldsymbol{\mu}(\cdot)$ konvergiert diese Iteration für $i \rightarrow \infty$ gegen die Lösung des Policy-Evaluation-Schrittes der zeitdiskreten Policy Iteration (2.22) [LV09]. Somit kann (2.25) als einzelner Schritt der für i definierten Iteration mit $\boldsymbol{\mu}(\cdot) = \boldsymbol{\mu}^{[l]}(\cdot)$ interpretiert werden. Dies deckt sich auch mit der Definition der PI und VI für endliche Zustands- und Stellgrößenräume in [SB18, Kapitel 4.3] bzw. [SB18, Kapitel 4.4].

²⁰ Bei der Anwendung einer PI oder VI auf klassische Markov-Entscheidungsprozesse mit endlichen Zustands- und Aktionsräumen trifft die Aussage zum Komplexitätsunterschied der Berechnung einer einzelnen Iteration z. B. zu (vgl. [KLM96, S. 251]). Bei Verwendung linearer Funktionsapproximatoren und einer blockweisen Least-Squares-Schätzung des Policy-Evaluation-Schrittes kann die Komplexität einer einzelnen Value Iteration ebenfalls reduziert werden, falls die auftretende Pseudoinverse vorab berechnet werden kann (vgl. [BJ16b, Remark 4.1]).

einzelnen Iterationen l getroffen werden, weshalb erst das finale (konvergierte) Regelgesetz für $l \rightarrow \infty$ zur Regelung verwendet werden sollte [LW14]. Zudem konvergieren PI-Algorithmen meist in weniger Iterationen l , da die Policy Evaluation eine komplette Bewertung des aktuellen Regelgesetzes vornimmt [LVV12], [Hey16].

2.1.4.3 Actor-Critic-Methoden

Actor-Critic-Methoden (vgl. [KT03], [GBLB12], [SB18, Kapitel 13]) adaptieren die Parameter von Actor- und Critic-Funktionsapproximatoren (vgl. Abschnitt 2.1.3). Diese Actor-Critic-Methoden sind häufig durch PI- oder VI-Algorithmen motiviert und erlauben eine unterschiedliche Granularität der Policy-Evaluation- und Policy-Improvement-Prozesse. Die Adaption des approximierten optimalen Regelgesetzes μ findet, ähnlich wie bei einer PI oder VI, üblicherweise basierend auf einer geschätzten Value Function oder Q-Function statt. Der Hauptunterschied hierbei ist jedoch, dass bei Actor-Critic-Methoden zumeist gradientenbasierte Methoden verwendet werden, um das Regelgesetz in die durch den Critic bestärkte Richtung anzupassen, anstatt direkt eine vollständige Minimierung basierend auf der aktuellen Schätzung des Critics durchzuführen (vgl. (2.23) und (2.26)).

Beispielsweise präsentieren Vamvoudakis und Lewis [VL10] eine zeitkontinuierliche Actor-Critic-Methode, bei welcher der Critic basierend auf dem quadratischen Fehler der Hamiltonfunktion (vgl. (2.14) mit Verwendung eines Funktionsapproximators) mithilfe eines normierten Gradientenabstiegs und der Actor mit einem modifizierten normierten Gradientenabstieg angepasst wird. Da die Methode von Vamvoudakis und Lewis [VL10] Kenntnis der Eingangsdynamik $g(\mathbf{x})$ erfordert, erweitern Bhasin et al. [BKJ⁺13], die den Critic mithilfe einer zeitkontinuierlichen Least-Squares-Formulierung und den Actor mit einem Gradientenabstieg mit anschließender Projektion adaptieren, die Actor-Critic-Struktur um ein zusätzliches neuronales Netz, das die Systemdynamik identifiziert. Silver et al. [SLH⁺14] präsentieren eine zeitdiskrete Actor-Critic-Formulierung, bei der die Critic-Parameter einer approximierten Q-Function basierend auf dem Gradienten des quadratischen TD-Fehlers (vgl. (2.7) unter Verwendung eines Funktionsapproximators) adaptiert werden. Die Actor-Gewichte θ , die das Regelgesetz μ approximieren, werden, beispielsweise mithilfe eines Gradientenabstiegs oder anderer gradientenbasierter Optimierungsverfahren, in Richtung des Gradienten des Gütemaßes J bezüglich des Actor-Gewichtes θ angepasst. Da dieser Mechanismus später in der Anwendung in Abschnitt 6.1 zum Einsatz kommt, ist eine kurze Einführung in Anhang D.2 gegeben. Darüber hinaus existieren unzählige weitere Actor-Critic-Ansätze. Beispielhaft sei die Arbeit von Li et al. [LGM20] zu nennen, in der ein einzelnes neuronales Netz genutzt wird, um sowohl Gewichte der Value Function²¹ als auch Gewichte eines optimiertes Regelgesetzes zu erlernen. Dieser Ansatz kann somit als Verschmelzung der Actor-Critic-Struktur zu einem kombinierten neuronalen Netz interpretiert werden.

²¹ Konkret ermöglicht die zusätzliche Nutzung der sogenannten *Advantage Function* (vgl. [WSH⁺16]), aus der geschätzten Value Function die Q-Function zu schätzen.

2.1.4.4 On-Policy- und Off-Policy-Ansätze

Ein wichtiges Klassifikationsmerkmal von RL- bzw. ADP-Algorithmen ist, ob diese *On-Policy-Algorithmen* oder *Off-Policy-Algorithmen* darstellen. Dabei sind On-Policy-Methoden dadurch charakterisiert, dass die approximierten Value Function V bzw. Q-Function Q die (ggf. diskontierten) Langzeitkosten des während des Adaptionsvorgangs tatsächlich verwendeten Regelgesetzes repräsentieren. Im Gegensatz dazu lernen Off-Policy-Ansätze eine approximierten Value Function V bzw. Q-Function Q , die zu einer sogenannten *Target Policy* gehört. Diese Target Policy weicht im Allgemeinen von der *Behavior Policy* ab, die auf das System angewandt wird, um Messdaten zu generieren (vgl. [KLJ17], [SLH⁺14], [LCL⁺19]). Konkret bedeutet das auch, dass das optimale Regelgesetz bei Off-Policy-Ansätzen nicht ausgeführt werden muss, um dieses zu erlernen [van12, S. 17]. Beispiele für On-Policy-Verfahren sind On-Policy-TD-Learning²² [SB18, Kapitel 6.4] oder die On-Policy-IRL-Methoden von Vrabie et al. [VPAKL09] oder Modares und Lewis [ML14a]. Die Klasse der Off-Policy-Algorithmen schließt beispielsweise klassisches Q-Learning [Wat89], [WD92], Off-Policy-IRL-Algorithmen [JJ12], [JJ14c], [SLW17] den LSPI-Algorithmus [LP03] und Off-Policy-Actor-Critic-Methoden [SLH⁺14] ein.

Wenn es gelingt, ADP-Methoden off-policy zu entwerfen, bringt dies zwei wesentliche Vorteile mit sich. Zum einen können aufgrund der Off-Policy-Charakteristik Daten beim Training wiederverwendet werden, was die Dateneffizienz signifikant erhöht und essenziell für Anwendungen, beispielsweise in der Robotik, ist [GHLL17]. Diese gezielte Mehrfachverwendung von Datentupeln ist im Bereich des RL unter dem Begriff des *Experience Replay* bekannt, welches maßgeblich für den Erfolg von RL mit tiefen neuronalen Netzen verantwortlich ist [MKS⁺15], [MKS⁺13]. Im Gegensatz dazu müssen bei On-Policy-Algorithmen nach jeder Adaption des Regelgesetzes alle bisherigen Trainingsdaten verworfen und neu aufgezeichnet werden. Dies stellt eine datenineffiziente und unpraktikable Einschränkung von On-Policy-Methoden dar [KIP⁺18]. Zum anderen führen On-Policy-Algorithmen bei der Verwendung von Explorationsrauschen, das zur Systemanregung benötigt wird (vgl. Kapitel 5), zu einem Offset in der Schätzung der Critic-Gewichte und somit auch dem approximierten Regelgesetz [LYD17, Theorem 1], [LCL⁺19, Lemma 3], [YDZY20, Lemma 1], [KLJ17, Remark 2]. Diesen Offset weisen Off-Policy-Verfahren hingegen nicht auf [LYD17, Theorem 3], [LCL⁺19, Theorem 3], [YDZY20, Theorem 3], [KLJ17, Theorem 3]. Insbesondere tragen Off-Policy-Ansätze folglich dazu bei, das sogenannte *Exploration-Exploitation-Dilemma* [SB18, S. 3] in den Griff zu bekommen, da die Behavior Policy eine umfassende Exploration ermöglicht, während dennoch eine Adaption in Richtung der optimalen Target Policy stattfinden kann (vgl. [LCL⁺19], [LHP⁺16]).

2.2 ADP-basierte Solltrajektorienfolgereger

Zahlreiche zeitdiskrete (z. B. [LLW⁺17], [WZL16], [WLLS17], [WHQ20]) und zeitkontinuierliche (z. B. [MCLS02], [KWD16], [Vam17], [Wan20]) ADP-Ansätze betrachten lediglich

²² Auch unter dem Namen *SARSA* bekannt.

die Regelung des Systemzustands auf eine konstante Ruhelage $\mathbf{x} = \mathbf{0}$ ²³ und keine Regelung auf Solltrajektorienverläufe $\mathbf{x}_{r,k}$ bzw. $\mathbf{x}_r(t)$. Da jedoch die Solltrajektorie die Kosten r und somit auch die Value Function V (vgl. (2.3a) und (2.20)) bzw. Q-Function Q (vgl. (2.7)) beeinflusst, sind diese Regler, die für eine konstante Ruhelage trainiert wurden, nicht für andere Sollzustände geeignet und somit insbesondere nicht auf den Trajektorienfolgeregelungsfall übertragbar.

Des Weiteren weisen Ansätze, die keine explizite funktionale Abhängigkeit der Value Function V oder Q-Function Q von der Referenztrajektorie verwenden, sondern die Form $V^\mu(\mathbf{x})$ bzw. $Q^\mu(\mathbf{x}, \mathbf{u})$ aufweisen, keine Übertragbarkeit auf eine andere als die während des Trainings verwendete (einzelne, fest vorgegebene) Trajektorie auf. Als Beispiel sei die Arbeit von Yu et al. [YSH⁺17] zu nennen, bei welcher der Sollzustand kein Argument der Q-Function darstellt, die Kosten r aber dennoch von einem zeitveränderlichen aktuellen Sollzustand abhängen. Yu et al. verwenden in ihrem Actor-Critic-Ansatz während des in zahlreichen Episoden wiederholten Trainingsvorgangs immer wieder dieselbe Solltrajektorie, weshalb diese mittels der Kosten r implizit Einfluss auf die gelernte Q-Function $Q^\mu(\mathbf{x}, \mathbf{u})$ hat. Jedoch trainieren sie ihren Regler damit auf genau diese während des Trainings verwendete Referenztrajektorie, da sowohl $Q(\cdot)$ als auch das gelernte Regelgesetz μ implizit von dieser Referenz abhängen, aber kein explizites Verständnis über die Abhängigkeit von einer während des Trainingsvorgangs nicht gesehenen Referenztrajektorie aufweisen (vgl. auch (3.4)). Die Verwendung einer anderen Solltrajektorie erfordert somit einen erneuten Trainingsvorgang (vgl. [YSH⁺17, Abschnitt 5]).

Im Folgenden werden in den Abschnitten 2.2.1–2.2.3 unterschiedliche Klassen ADP-basierter Solltrajektorienfolgeregelungsansätze aus der Literatur vorgestellt und in Abschnitt 2.2.4 eine kurze Übersicht der Anwendung von ADP-Methoden auf reale regelungstechnische Anwendungen gegeben, um die erste zentrale Forschungslücke herauszuarbeiten.

2.2.1 ADP-Ansätze unter Nutzung der dynamischen Inversion

Zahlreiche ADP-basierte Methoden nutzen einen auf einer Zustandslinearisierung (vgl. beispielsweise [Kha02, Kapitel 13]) basierenden Vorsteuerentwurf, der sich in der Literatur unter dem Begriff einer *dynamischen Inversion* findet [EBHS07]. Diese Konzepte bestimmen üblicherweise a priori und unter vollständiger Kenntnis oder nach vorheriger Identifikation der Systemdynamik einen inversionsbasierten Vorsteuerterm und nutzen erst dann ADP-Verfahren, um lediglich den Zustandsrückführungsterm des Reglers zu adaptieren. Somit wird genau genommen kein datenbasierter ADP-Solltrajektorienfolgeregler gelernt, da der Vorsteuerterm modellbasiert berechnet wird. Tabelle 2.1 gibt eine Übersicht über diese Methoden. Die Spalte *globale Diskontierung* γ gibt hierbei Aufschluss darüber, ob ein Diskontierungsfaktor γ

²³ Wenngleich diese Methoden prinzipiell auch für einen konstanten Sollzustand ungleich null trainiert werden könnten, so ist dieser Sollzustand zur Laufzeit nicht veränderbar und auch kein expliziter Parameter des gelernten Regelgesetzes μ . Vielmehr wird während des Trainings implizit der Einfluss des einmalig fest gewählten Sollzustands auf V bzw. Q gelernt. Dies kann letztlich als Festwertregelung (vgl. [Lun20a, S. 362]) mit der zusätzlichen wesentlichen Einschränkung, dass die Führungsgröße nicht verändert werden kann, d. h. ohne äußere Eingriffsmöglichkeit, interpretiert werden.

im Gütemaß verwendet wird (vgl. (2.2)), um die Value- bzw. Q-Function endlich zu halten, oder welche anderen Mechanismen hierfür genutzt werden.

Zhang et al. [ZWL08] stellen einen zeitdiskreten Ansatz vor, der eine Value Iteration nutzt. Huang und Liu [HL14] präsentieren ebenfalls eine zeitdiskrete Methode. Neben einem neuronalen Netz zur Approximation der Systemdynamik verwenden sie einen Actor-Critic-Ansatz und passen sowohl die Critic-Gewichte als auch den Funktionsapproximator zur Beschreibung des Zustandsrückführungsterms mithilfe eines Gradientenabstiegs an. Auch die zeitkontinuierliche ADP-Trajektorienfolgeregelungsmethode von Kamalapurkar et al. [KDBD15] verwendet eine derartige dynamische Inversion sowie ein Actor-Critic-Verfahren.

Ein Nachteil all dieser Methoden ist jedoch die im Allgemeinen nicht erfüllte Annahme, dass die Eingangsmatrix $\mathbf{g}(\mathbf{x})$ bzw. \mathbf{B} invertierbar²⁴ sein muss. Zudem muss die gesamte Systemdynamik entweder bekannt sein oder identifiziert werden (vgl. [KLM⁺14], [KLL15]). Mu et al. [MSWS17] fordern bei ihrem Ansatz für zeitkontinuierliche, eingangsaffine Systeme ebenfalls, dass die durch $\mathbf{f}(\mathbf{x})$ und $\mathbf{g}(\mathbf{x})$ beschriebene Systemdynamik bekannt ist, ihr Ansatz erlaubt jedoch immerhin eine beschränkte und zustandsabhängige Unsicherheit in der Systemdynamik. Dierks und Jagannathan [DJ09] nutzen in ihrem gradientenbasierten Actor-Critic-Ansatz

Veröffentlichung	Systemmodell erforderlich (oder identifiziert)	zeitkontinuierlich (k) oder zeitdiskret (d)	Systemdynamik	Value Function V oder Q-Function Q	Grundalgorithmus	globale Diskontierung γ	dynamische Inversion
[ZWL08]	ja	d	eingangsaffin	V	VI	nein ^c	ja
[HL14]	ja	d	eingangsaffin	V	AC ^a	nein ^c	ja
[KDBD15]	ja	k	eingangsaffin	V	AC ^b	nein ^c	ja
[MSWS17]	ja ^d	k	eingangsaffin	V	PI	nein ^c	ja
[DJ09]	teilw.	d	eingangsaffin	V	AC ^a	nein ^c	ja

^a Actor-Critic mit (ggf. normiertem) Gradientenabstieg.

^b Actor-Critic, Critic mit zeitkontinuierlichem Least-Squares-Update mit Vergessensfaktor.

^c Gütemaß bestraft Abweichung der Stellgröße vom Vorsteuerungsanteil.

^d Jedoch Unsicherheit im Modell erlaubt.

Tabelle 2.1: Übersicht über ADP-Trajektorienfolgeregelungsmethoden, die auf dem Konzept der dynamischen Inversion basieren.

²⁴ Im Sinne der Existenz einer Matrix $\tilde{\mathbf{g}}(\mathbf{x})$, sodass $\mathbf{g}(\mathbf{x})\tilde{\mathbf{g}}(\mathbf{x}) = \mathbf{I}$ gilt (vgl. beispielsweise [DJ09], [DJ10]).

neben Funktionsapproximatoren zur Schätzung des Critics und des Zustandsrückführungsterms noch einen dritten Funktionsapproximator²⁵, um auch den Vorsteuerterm zu adaptieren. Dennoch muss auch bei diesem Verfahren $g(x)$ bekannt und invertierbar sein. Abschließend sei noch angemerkt, dass bei all diesen Konzepten, die den Ansatz der Zustandslinearisierung nutzen, im Gütefunktional $J(\cdot)$ die quadratische Abweichung der Stellgröße vom Vorsteuerterm und somit nur die transiente Stellgröße und nicht die eigentliche Stellenergie bestraft wird. Ein geschlossener optimierungsbasierter Ansatz liegt somit nicht vor.

2.2.2 Globale Solltrajektorienvorgabe durch eine Exosystemdynamik ohne externe Eingriffsmöglichkeit

Eine andere in der ADP-Literatur weitverbreitete Annahme ist, dass die Referenztrajektorie global einer (häufig unbekannt) Exosystemdynamik folgt [ML13], [ML14a], [ML14b], [ZZXS17], [KL15], [LLHW16], [KLM⁺14], [LYD17], [Vam16], [KLNSK15], [QZL13], [GJ16], [GJ15], [BA18], [KLL15]. Basierend auf dieser Annahme wird der Systemzustand um einen Exosystemzustand erweitert, sodass auf bestehende ADP-Mechanismen zurückgegriffen werden kann. Eine Übersicht über ADP-Trajektorienfolgeregelungsmethoden, die eine solche globale Vorgabe der Referenztrajektorie annehmen, ist in Tabelle 2.2 gegeben. Diese Ansätze werden im Folgenden diskutiert.

Modares und Lewis [ML13] stellen für zeitkontinuierliche linear-quadratische optimale Trajektorienfolgeregelungsprobleme einen Policy-Iteration-Algorithmus vor, der jedoch die vollständige Kenntnis der erweiterten Systemdynamik benötigt. Darüber hinaus existieren sowohl zeitkontinuierliche [ML14a], [ML14b], [ZZXS17] als auch zeitdiskrete [KL15] Ansätze, welche die Value Function V approximieren und zwar keine Kenntnis über die interne Systemdynamik $f(x)$ bzw. A benötigen, die Eingangsmatrix $g(x)$ bzw. B jedoch bekannt sein muss. Diese Einschränkung einer bekannten Eingangsmatrix weisen die Ansätze von Luo et al. [LLHW16] (zeitdiskrete, nichtlineare Systemdynamik), Kiumarsi et al. [KLM⁺14] und Li et al. [LYD17] (jeweils zeitdiskrete, lineare Systemdynamik), sowie Vamvoudakis [Vam16] (zeitkontinuierliche, lineare Systemdynamik), die jeweils eine Q-Function lernen, nicht auf. Kiumarsi et al. [KLNSK15] betrachten eine lineare, zeitdiskrete Systemdynamik und lernen eine ADP-basierte optimale Ausgangsrückführung, indem der Systemzustand aus einer endlichen Anzahl vergangener Werte der Stellgrößen, Ausgangsgrößen und Sollzustände rekonstruiert wird und anschließend alternativ entweder eine Policy Iteration oder eine Value Iteration verwendet wird. Qin et al. [QZL13] stellen für zeitkontinuierliche linear-quadratische Ausgangsfolgeregelungsprobleme eine IRL-basierte Policy Iteration vor, die keine Kenntnis der Werte der Systemmatrizen A und B benötigt. Hierbei wird jedoch einschränkend angenommen, dass die Anzahl der Exosystemzustände identisch zur Anzahl der Ausgangsgrößen des zu regelnden Systems ist. Gao und Jiang [GJ16], [GJ15] verwenden eine auf dem Prinzip des IRL basierende Policy Iteration, um linear-quadratische Trajektorienfolgeregelungsprobleme mit zeitkontinuierlicher Systemdynamik und unbekannt Systemmatrizen A und B zu lösen. Bernhard und Adamy [BA18] betrachten Trajektorienfolgeregelungsprobleme mit linearer,

²⁵ Dieser schätzt letztlich implizit die interne Dynamik $f(\cdot)$.

Veröffentlichung	Systemmodell erforderlich (oder identifiziert)	zeitkontinuierlich (k) oder zeitdiskret (d)	Systemdynamik	Value Function V oder Q -Function Q	Grundalgorithmus	globale Diskontierung γ	dynamische Inversion
[ML13]	ja	k	linear	V	PI	ja	nein
[ML14a]	teilw.	k	linear	V	PI	ja	nein
[ML14b]	teilw.	d	eingangsaffin	V	AC ^a	ja	nein
[ZZXS17]	teilw.	k	eingangsaffin	V	PI	ja	nein
[KL15]	teilw.	d	eingangsaffin	V	AC ^a	ja	nein
[LLHW16]	nein	d	nichtlinear	Q	PI	ja	nein
[KLM ⁺ 14]	nein	d	linear	Q	PI	ja	nein
[LYD17]	nein	d	linear	Q	PI	ja	nein
[Vam16]	nein	k	linear	Q	AC ^a	nein ^b	nein
[KLNSK15]	nein	d	linear	V	PI/VI	ja	nein
[QZL13]	nein	k	linear	V	PI	nein ^b	nein
[GJ16], [GJ15]	nein	k	linear	V	PI	nein ^c	nein
[BA18]	nein	k	linear	V	PI	nein ^d	nein
[KLL15]	teilw.	d	eingangsaffin	V	PI/VI	ja	nein

^a Actor-Critic mit (ggf. normiertem) Gradientenabstieg.

^b Exosystem muss stabil sein.

^c Gütemaß bestraft Abweichung der Stellgröße vom Vorsteuerungsanteil.

^d Unendliche Kosten erlaubt.

Tabelle 2.2: Übersicht über ADP-Trajektorienfolgeregelungsmethoden, die eine globale Vorgabe der Referenztrajektorie annehmen.

unbekannter, zeitkontinuierlicher Systemdynamik und unendlichen Kosten²⁶. Kiumarsi et al. [KLL15] behandeln zeitdiskrete, eingangsaffine, zeitveränderliche Systeme. Hierbei werden für unterschiedliche Systemdynamiken unterschiedliche Value Functions gelernt, anschließend wird mittels sogenannter *adaptive self-organizing maps* automatisiert zwischen diesen Repräsentationen gewechselt.

²⁶ Hierbei wird die Definition der überholenden Optimalität (vgl. [Ber20a, Definition 2.3], (engl.): *overtaking optimality*) verwendet.

Alle im vorliegenden Abschnitt diskutierten Arbeiten weisen jedoch einen entscheidenden Nachteil auf: die Annahme, dass der Exosystemzustand ausgehend von dessen Initialzustand propagiert wird und zur Laufzeit nicht durch eine externe Eingriffsmöglichkeit verändert werden kann²⁷. Die Referenztrajektorie ist daher global vorgegeben und nicht von außen beeinflussbar²⁸. Daher können beispielsweise Nutzereingaben, Straßenverläufe oder gewünschte Systemzustandsverläufe, die aus einer übergeordneten Ebene resultieren, nicht berücksichtigt werden. Als illustratives Beispiel, um die damit verbundene Problematik zu veranschaulichen, werde ein Fahrzeug betrachtet, das einem beliebigen, realen Straßenverlauf folgen soll. Um mit diesem Fahrzeug beispielsweise von Karlsruhe nach Kaiserslautern zu fahren, müsste unter Verwendung der in diesem Abschnitt diskutierten Methoden der gesamte Straßenverlauf auf dieser Strecke durch die Exosystemdynamik und dessen Anfangszustand beschrieben sein. Während dies bereits eine nicht handhabbare Annahme darstellt, müsste zudem für jede andere Fahrtstrecke ein geeignetes Exosystem gefunden und der Regler erneut trainiert werden. Damit sind mit den bestehenden Methoden insbesondere auch keine beliebigen Solltrajektorien als Eingabe der Regelungsmethoden möglich. Diese Eigenschaft ist in Abbildung 2.2 veranschaulicht und stellt erhebliche Einschränkungen für die allgemeine Anwendbarkeit dieser adaptiven optimalen Trajektorienfolgeregelungen dar.

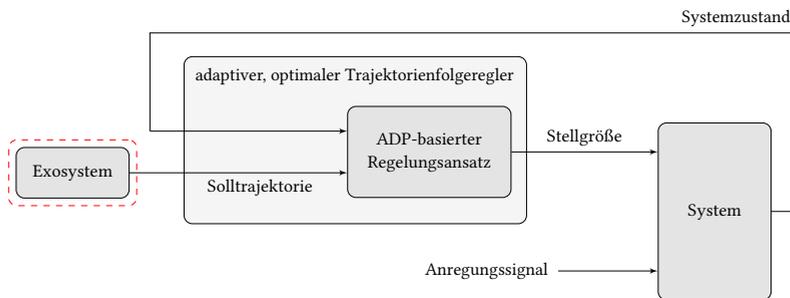


Abbildung 2.2: Schematische Darstellung von Literaturansätzen, die annehmen, die Referenztrajektorie folge global einer Exosystemdynamik. Eine flexible Beeinflussbarkeit des Solltrajektorienverlaufs von außen fehlt hierbei.

2.2.3 Stationäre Sollzustandsvorgabe

Eine Alternative zur Annahme, die Referenztrajektorie werde global durch eine Exosystemdynamik erzeugt, stellt die Erweiterung des Systemzustands um den aktuellen Sollzustand oder die Abweichung vom aktuellen Sollzustand [HSSH17], [PRH19], [PKRH20], [WXL⁺14], bzw. die Erweiterung um einen projizierten aktuellen Sollzustand [NKJS04] dar. Shi et al. [SSW18] führen diese Idee fort, indem der Systemzustand nicht nur um den Sollzustand im

²⁷ Andernfalls würde die Markov-Eigenschaft verletzt und die Trainingsdaten wären nicht mit der Bellman-Gleichung (2.5) bzw. HJB-Gleichung (2.16) verträglich (vgl. [van12], Abschnitt 3.1 und Abschnitt 4.1).

²⁸ Dies kann als Servoregelung (vgl. [Lun20a, S. 398], [AM89, S. 89]) ohne äußere Eingriffsmöglichkeit interpretiert werden, bei der das Exosystem die Führungsgröße generiert.

aktuellen Zeitschritt, sondern auch um den Sollzustand des direkt nachfolgenden Zeitschritts ergänzt wird. Zwar lassen die ADP-basierten Methoden, welche den aktuellen (und ggf. nachfolgenden) Sollzustand in der Value Function bzw. Q-Function explizit berücksichtigen, die Vorgabe beliebiger stationärer Sollzustände zu (vgl. Abbildung 2.3), jedoch wird der zeitliche Verlauf der Solltrajektorie hierbei nicht (oder im Fall von Shi et al. [SSW18] nur sehr eingeschränkt) berücksichtigt. Vielmehr repräsentiert die Value Function $V^\mu(\mathbf{x}_k, \mathbf{x}_{r,k})$ nur die unter dem Regelgesetz μ entstehenden Langzeitkosten, wenn sich das System im Zustand \mathbf{x}_k befindet und auf den stationären Endwert $\mathbf{x}_{r,k}$ eingeregelt werden soll, was jedoch nicht dem nach (2.11) zu minimierenden Gütefunktional entspricht. Der Regler, der eine Festwertregelung (vgl. [Lun20a, S. 362]) darstellt, reagiert damit lediglich auf Abweichungen zum aktuellen Sollzustand, berücksichtigt jedoch nicht den zukünftigen Verlauf der Referenztrajektorie. Dies kann zu einem zeitlichen Versatz zwischen dem Istzustand und dem Sollzustand führen.

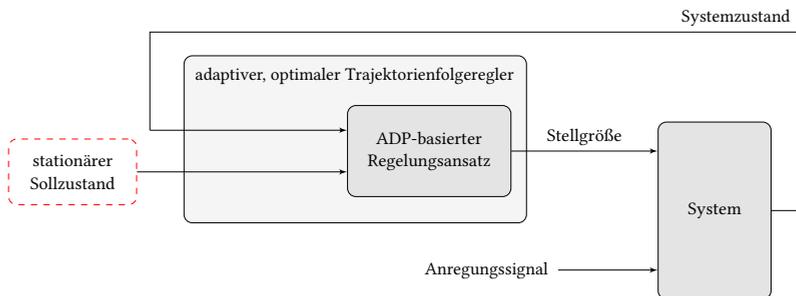


Abbildung 2.3: Schematische Darstellung ADP-basierter Optimalregelungskonzepte aus der Literatur, welche die Vorgabe eines stationären Sollzustands verwenden.

2.2.4 ADP in realen regelungstechnischen Anwendungen

Nach diesen Einblicken in unterschiedliche Klassen ADP-basierter Solltrajektorienfolgeregler wird im Folgenden ein kurzer Überblick hinsichtlich bisheriger realer ADP-Anwendungen gegeben. Obwohl zahlreiche Publikationen zu ADP-basiertem Reglerentwurf und zu ADP-basierten Zustandsfolgereglern existieren, betrachten die meisten dieser Veröffentlichungen lediglich Simulationsergebnisse, jedoch keine reale Anwendung der Methoden im Experiment. Insofern ist die empirische Basis dünn und vernachlässigt praktische Herausforderungen, z. B. aufgrund von Mess- und Stellgliedernauigkeiten. Als prominente Beispiele seien [LV09], [WLMZ18], [JJ14c], [BKJ⁺13], [ML14a], [LLHW16] und [KLM⁺14] genannt. Gleiches gilt für das in Abschnitt 6.1 betrachtete Anwendungsbeispiel ADP- oder RL-basierter Geschwindigkeitsregler für Fahrzeuge. Auch hier existieren zumeist nur Simulationsergebnisse (siehe [PRH19], [BK18], [NCH08], [DCd11] und [KBJ⁺21, Table II]) oder Laborergebnisse anhand eines Fahrsimulators [PT12]. Huang et al. [HXH⁺19] arbeiten zwar mit einem Realfahrzeug, trainieren ihren Geschwindigkeitsregler jedoch basierend auf einem stark vereinfachten, datenbasierten Modell der Longitudinaldynamik. Aus diesem Modell werden Stichproben

gezogen, anhand derer ein Längsregler trainiert wird. Zudem findet keine Berücksichtigung des zukünftigen Solltrajektorienverlaufs statt, der Adaptionsprozess erfolgt nicht anhand von Realdaten während der Fahrt und dynamische Vorgänge werden ignoriert, d. h. interne, nicht messbare Systemgrößen werden nicht rekonstruiert (vgl. [PRH19]).

Auch die Betrachtung anderer Anwendungsbeispiele liefert ein ähnliches Bild. Li et al. [LGM20] demonstrieren beispielsweise mithilfe von Simulationsergebnissen, wie ADP genutzt werden kann, um das Energiemanagement in hybridelektrischen Fahrzeugen zu optimieren. Li und Görges [LG19], [LG20] stellen Ansätze für modellfreie, ADP-basierte Abstandsregeltempomaten vor und präsentieren Simulationsergebnisse. Shi et al. [SSW18] untersuchen einen RL-basierten Regler für ein Unterwassergefährt, wobei eine stationäre Sollzustandsvorgabe (vgl. Abschnitt 2.2.3) verwendet wird. Die dortigen Ergebnisse beschränken sich auf Simulationen, eine reale Anwendung findet nicht statt. Hwangbo et al. [HSSH17] stellen einen RL-basierten Regler vor, der neuronale Netze trainiert und die Position eines Quadrocopters regelt. Dem Regler wird dabei lediglich eine Sollposition vorgegeben (Festwertregelung), an welcher sich der Quadrocopter stabilisieren soll. Eine Berücksichtigung des Verlaufs der Solltrajektorie findet nicht statt. Auch die von Ng et al. [NKJS04] präsentierte RL-basierte Regelung eines Modellhelikopters verwendet als Vorgabe nur die aktuelle Abweichung von einer Wunschposition, anstatt den Sollverlauf zu übergeben. Sowohl Hwangbo et al. [HSSH17] als auch Ng et al. [NKJS04] nutzen zwar ihre durch einen RL-basierten Ansatz trainierten Regler nach dem Trainingsvorgang an realen Systemen, jedoch basiert das Erlernen des Regelgesetzes jeweils auf Reinforcement Learning anhand von Daten, die aus einem Simulationsmodell erzeugt wurden und nicht direkt mittels Messdaten. Somit muss ein geeignetes Systemmodell für das Training vorliegen, also zunächst modelliert oder identifiziert werden. Es bleibt ungeklärt, inwiefern eine Adaption des Reglers direkt basierend auf realen Messdaten möglich ist. Die in Kapitel 1 erwähnten Arbeiten [KIP⁺18], [ABC⁺20] erfordern eine enorme Menge an Real- bzw. Simulationsdaten für einen erfolgreichen Trainingsvorgang. Zudem wird in [KIP⁺18] im Fall eines erfolgreichen Greifversuchs des Roboters eine binäre Belohnung ausgeschüttet und in [ABC⁺20] nur die Abweichung von einer stationären Soll-Orientierung eines in einer Roboterhand befindlichen Würfels bestraft und kein Solltrajektorienverlauf explizit vorgegeben.

Insbesondere eine reale Anwendung eines modellfreien ADP-Solltrajektorienreglers, der nicht nur den aktuellen Sollzustand, sondern einen flexiblen Sollzustandsverlauf berücksichtigt, und dabei strukturelles Vorwissen über das zugrunde liegende Problem ausnutzt, ist somit bis zum jetzigen Zeitpunkt in der Literatur nicht vorhanden. Aufgrund von Unterschieden zwischen Simulationen und Realanwendungen (vgl. [DAMH19], [CHL19], [Bro92]) ist aus einer anwendungsorientierten regelungstechnischen Perspektive eine Validierung der direkten Anwendung der ADP-Methoden auf reale Systeme jedoch unverzichtbar.

2.2.5 Fazit

Bislang existierende ADP-Solltrajektorienfolgeregelungsansätze lassen sich in drei verschiedene Gruppen klassifizieren. Inversionsbasierte Methoden nach Abschnitt 2.2.1 erfordern

Kenntnis über die Systemdynamik und bestrafen nicht die eigentliche Stellenergie. Daher werden sie nachfolgend nicht weiter für den ADP-Kontext betrachtet. Ansätze, die annehmen, dass die Referenztrajektorie einer globalen Exosystemdynamik folgt (vgl. Abschnitt 2.2.2), weisen bislang keine äußere Beeinflussungsmöglichkeit dieses Solltrajektorienverlaufs auf (vgl. Abbildung 2.2). Die Vorgabe eines stationären Endwertes nach Abschnitt 2.2.3 erlaubt zwar eine Sollwertvorgabe von außen (vgl. Abbildung 2.3), der weitere zeitliche Verlauf der Referenztrajektorie wird hierbei jedoch nicht berücksichtigt. Abschnitt 2.2.4 unterstreicht zudem, dass bislang nur wenige reale Anwendungen ADP-basierter Regelungsansätze betrachtet wurden und diese insbesondere keine flexible Solltrajektorienvorgabe erlauben.

Somit kann konstatiert werden, dass ein generalisierender Ansatz (vgl. Abbildung 2.4), der aus der externen Vorgabe einer Solltrajektorie eine für die Verwendung durch ADP-Ansätze kompatible²⁹ Approximation generiert, die zugleich Informationen über den zeitlichen Verlauf der Referenztrajektorie beinhaltet, bis dato nicht existiert.

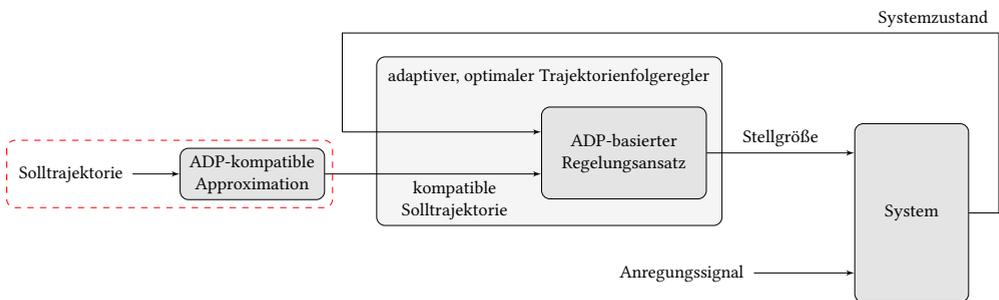


Abbildung 2.4: Schematische Darstellung eines generalisierenden Ansatzes, der die externe Vorgabe eines Solltrajektorienverlaufs erlaubt und diese Referenztrajektorie auf eine mit dem ADP-Formalismus kompatible Weise approximiert. Eine solche Methode, die vielfältige, flexible Solltrajektorien in den ADP-Mechanismus integrieren kann, existiert bislang nicht.

2.3 Anregung ADP-basierter Regelungsansätze

Für Konvergenzaussagen von Systemidentifikationsmethoden sowie adaptiven und ADP-basierten Algorithmen wird stets die Annahme gefordert, dass das betrachtete System ausreichend angeregt ist. Am Beispiel des ADP lässt sich anschaulich plausibilisieren, weshalb eine geeignete Anregung des Systems erforderlich ist. So lassen sich beispielsweise bei einem System, das sich dauerhaft in einer stabilen Ruhelage $x = 0$ befindet, keine umfassenden Informationen über dessen dynamisches Verhalten gewinnen. Sämtliche Regelgesetze mit $\mu(0) = 0$ erscheinen dann aus Sicht eines zu minimierenden Gütefunktions gleichermaßen

²⁹ Eine ADP-kompatible Solltrajektorienarstellung erlaubt die Integration eines Sollzustandsverlaufs in ADP-basierte Regelungsansätze. Eine genaue Definition erfolgt durch Definition 3.1 für den zeitdiskreten Fall und durch Definition 4.1 für den zeitkontinuierlichen Fall.

optimal, da der Fall $\mathbf{x} \neq \mathbf{0}$ nicht beurteilt werden kann. Erst durch die Betrachtung unterschiedlicher Zustands- und Stellgrößenkombinationen, d. h. durch eine geeignete *Exploration* (vgl. [SB18, S. 3]) können Informationen darüber gewonnen werden, welche Regelgesetze hinsichtlich des Gütefunktionalis zu bevorzugen sind.

Eine Formalisierung derartiger Anregungsannahmen wird unter dem Begriff der sogenannten *Persistent-Excitation-Bedingung*³⁰ zusammengefasst. Dieser Begriff entstand in den 1960er-Jahren zunächst im Kontext der Systemidentifikation [ÅB66] (vgl. [NA05, S. 238], [BU16]), spielt jedoch auch bei adaptiven Regelungsansätzen [NA05, S. 239], [IS96] und ADP-Ansätzen [VLV13, S. 62] eine zentrale Rolle bei der Betrachtung von Konvergenzeigenschaften.

Konkret wird für zeitkontinuierliche Signale $\boldsymbol{\sigma}(t) \in \mathbb{R}^h$ die Einhaltung einer PE-Bedingung

$$\int_t^{t+T} \boldsymbol{\sigma}(\tau) \boldsymbol{\sigma}^\top(\tau) \, d\tau \succeq \alpha \mathbf{I}, \quad (2.27)$$

$\alpha, T \in \mathbb{R}_{>0}$, $\forall t \geq t_0$, gefordert (vgl. [NA05, Definition 6.2], [PSA17, Definition 4]). Falls das Signal $\boldsymbol{\sigma}(t)$ diese PE-Bedingung erfüllt, so wird es im Folgenden als *PE-Signal* bezeichnet. Die Relevanz der PE-Bedingung (2.27) offenbart sich anhand der häufig auftretenden Differenzialgleichung

$$\frac{d\tilde{\mathbf{w}}(t)}{dt} = -\boldsymbol{\sigma}(t) \boldsymbol{\sigma}^\top(t) \tilde{\mathbf{w}}(t), \quad (2.28)$$

wobei $\tilde{\mathbf{w}}(t) \in \mathbb{R}^h$ den Fehler eines zu schätzenden Parameters $\mathbf{w} \in \mathbb{R}^h$ darstellt. Für die gleichmäßige asymptotische Stabilität der Ruhelage von (2.28) ist es notwendig und hinreichend, dass $\boldsymbol{\sigma}(t)$ ein PE-Signal ist [NA05, S. 246]. Anschaulich interpretiert genügt aufgrund von $\text{Rang}(\boldsymbol{\sigma}(t) \boldsymbol{\sigma}^\top(t)) \leq 1$ die in $\boldsymbol{\sigma}(t)$ zu einem einzelnen Zeitpunkt t enthaltene Information nicht, um die h Unbekannten des Parametervektors \mathbf{w} zu bestimmen. Die PE-Bedingung (2.27) fordert daher, dass das Integral der positiv semidefiniten Matrix $\boldsymbol{\sigma}(\tau) \boldsymbol{\sigma}^\top(\tau)$ über ein endliches Intervall $[t, t+T]$, $\forall t \geq t_0$, positiv definit ist (vgl. [NA05, S. 247]). Eine zweite anschauliche Interpretation ergibt sich wie folgt: Nach [NA05, S. 247] ist das Integral des Betrags der Projektion eines PE-Signals $\boldsymbol{\sigma}(t)$ entlang jedes beliebigen Einheitsvektors des \mathbb{R}^h über ein endliches Intervall $[t, t+T]$, $\forall t \geq t_0$, ungleich null. Dies bedeutet, dass in diesem endlichen Intervall das Signal $\boldsymbol{\sigma}(t)$ eine Basis des \mathbb{R}^h durchlaufen muss (vgl. [Kar19, S. 42]).

Die Erfüllung der PE-Bedingung für den Regressor $\boldsymbol{\sigma}(t)$ ist für Konvergenzaussagen bei der zeitkontinuierlichen Parameteradaptation [Ngu18, S. 132], [JALG18], [Pra17], der adaptiven Regelung [NA05, Abschnitt 6.5.2], [BS86], [AG08], [LK98], [NMH⁺15] und bei ADP-Ansätzen [VL10], [VL11], [Vam17], [JKBD15], [ML14b], [VVL09b], [TCTH19], [BKJ⁺13], [LLW14], [LYWW15], [ZCZL11], [SWL19, Assumption 6.4], [ZZXS17], [ZDJ14] zentral.

Das zeitdiskrete Pendant der durch (2.27) formulierten PE-Bedingung ist durch

$$\sum_{i=k}^{k+T_d} \boldsymbol{\sigma}_i \boldsymbol{\sigma}_i^\top \succeq \alpha \mathbf{I} \succ \mathbf{0} \quad (2.29)$$

³⁰ In seltenen Fällen ist auch die deutschsprachige Bezeichnung *fortwährende Erregung* [BU16] zu finden.

mit $T_d \in \mathbb{N}$, $T_d > 0$, gegeben [Bit84], [Mar84], [GM86]. Diese Bedingung findet sich unter anderem bei rekursiven zeitdiskreten ADP-Algorithmen [ATLAK07], [LLHW16], [BYB94], [ZQJL14], aber auch bei zeitkontinuierlichen ADP-Ansätzen, die mithilfe von IRL-Methoden gelöst werden [BJ16a], [FC16].

Zu erwähnen seien an dieser Stelle außerdem Arbeiten, die zur PE-Bedingung äquivalente [PSA17] und alternative [CJ10], [KKD14], [KWD13] Anregungsbedingungen vorstellen. Die unter dem Begriff des *concurrent learning* aufzufindenden Methoden der adaptiven Regelung [CJ10] bzw. ADP-Literatur [KKD14], [KWD13] fordern hierbei die Erfüllung einer Rangbedingung an zuvor abgespeicherte Daten. Zwar sind diese Rangbedingungen ggf. in der Anwendung nachträglich komfortabel überprüfbar, jedoch bleibt auch bei diesen Ansätzen ungeklärt, wie Systeme angeregt werden sollten, um geeignete Daten zu erzeugen. Somit wird auch hier das Vorhandensein geeignet angeregter Daten a priori vorausgesetzt (vgl. [CJ10, Condition 1], [KKD14, Assumption 1]).

Da in Kapitel 5 eine Betrachtung der zeitkontinuierlichen PE-Bedingung nach (2.27) stattfindet, konzentriert sich die nachfolgende Übersicht auf diesen Fall. Für die Erfüllung der PE-Bedingung spielt dabei insbesondere der Frequenzgehalt von σ eine wesentliche Rolle (vgl. [SB89, Abschnitt 4.3], [LK98], [NA87]). So ist beispielsweise ein skalares Signal $\sigma_s \in \mathbb{R}$ genau dann ein PE-Signal, wenn dessen Leistungsdichtespektrum mindestens eine Frequenzlinie aufweist [NA05, S. 253], [Kar19, Lemma 4.2]. Für den Fall $\sigma \in \mathbb{R}^h$ sind Aussagen für den Zusammenhang zwischen dem Frequenzspektrum von σ und der Erfüllung der PE-Bedingung in [BS83, Lemma 3.4] gegeben. Bei der Analyse von Signalen ist zudem die Frage, unter welchen mathematischen Operationen die PE-Eigenschaft eines Signals erhalten bleibt, relevant. Die bestehende Literatur zu Systemidentifikationsmethoden und adaptiver Regelung offenbart hierbei, dass für *lineare* algebraische und dynamische Transformationen von Signalen theoretische Aussagen hinsichtlich der PE-Eigenschaft existieren (vgl. [NA05, S. 249 ff.], [BS86]). Ein Beispiel für eine solche algebraische Signaltransformation ist, dass bei linearen Abbildungen eines Signals durch Matrizen $M \in \mathbb{R}^{m_1 \times m_2}$ mit Maximalrang und $m_1 \leq m_2$ dessen PE-Eigenschaft erhalten bleibt [NA05, Lemma 6.1]. Als Beispiel für eine dynamische Transformation lässt sich nennen, dass der Systemzustand eines linearen, Eingangs-Ausgangs-stabilen, steuerbaren dynamischen Systems n -ter Ordnung mit skalarer Eingangsgröße genau dann PE ist, wenn das Leistungsdichtespektrum des skalaren Eingangssignals mindestens n Frequenzlinien aufweist [BS86], [NA05, S. 255 ff.]. Aussagen für mehrdimensionale Eingangsgrößen sind in [MB90] und [GM86] zu finden.

Bei ADP-basierten Ansätzen treten jedoch selbst bei scheinbar einfachen Problemstellungen Nichtlinearitäten in den Basisfunktionen $\phi(\mathbf{x}(t))$ der linearen Funktionsapproximatoren (vgl. Abschnitt 2.1.3) auf. Ein Beispiel ist durch LQ-Optimierungsprobleme gegeben (vgl. [PLB15, S. 303]). Da diese Nichtlinearitäten den Frequenzgehalt der betrachteten Signale verändern können, und sich aus dem Frequenzgehalt des Systemzustands $\mathbf{x}(t)$ somit keine direkten Rückschlüsse auf den Frequenzgehalt der Basisfunktion $\phi(\mathbf{x}(t))$ (vgl. (2.21)) bzw. deren zeitliche Ableitungen $\frac{d\phi(\mathbf{x}(t))}{dt}$ ³¹ ziehen lassen, sind existierende theoretische Erkenntnisse nicht ohne

³¹ Gerade für das Signal $\frac{d\phi(\mathbf{x}(t))}{dt}$ ist jedoch die Erfüllung der PE-Bedingung von entscheidender Bedeutung für die Konvergenz des Critic-Gewichts w (vgl. [VL10, Theorem 1]).

Weiteres auf den ADP-Kontext übertragbar. Abhängig von den konkreten Nichtlinearitäten können beispielweise zusätzliche Frequenzen auftreten und sich positiv auf die Erfüllung der PE-Eigenschaft auswirken oder Frequenzen können sich gegenseitig aufheben und somit die PE-Bedingung verletzen (vgl. [LK98], [LK99]). Eine sorgfältige Analyse des Einflusses der vorhandenen Nichtlinearitäten auf die PE-Eigenschaft ist daher insbesondere im ADP-Kontext unumgänglich.

Unter den wenigen in der Literatur verfügbaren Ansätzen, welche die PE-Bedingung im Zusammenhang mit Nichtlinearitäten analysieren, sind die Arbeiten von Lin und Kanellakopoulos [LK98] (System in strenger Ausgangs-Rückkopplungsform), [LK99] (System in strenger parametrischer Rückkopplungsform) sowie von Adetola und Guay [AG06] (zustandslinearisierbare Systeme mit linearer Parameterabhängigkeit), [AG08] (nichtlineare Systeme mit linearer Parameterabhängigkeit) zu nennen, die jedoch dem Bereich der adaptiven Regelung zuzuordnen sind. Allgemeine Untersuchungen und theoretische Analysen hinsichtlich einer geeigneten Anregung im ADP-Kontext existieren bislang nicht [JZLH19], [JKBD15], [KLL15], [BKJ]⁺13].

Obwohl eine geeignete Anregung insbesondere auch bei ADP-Methoden stets erforderlich ist, wird in vorhandenen Publikationen üblicherweise die Erfüllung einer auf die jeweils konkrete Formulierung zugeschnittenen PE-Bedingung als gegeben vorausgesetzt. Eine methodische Analyse findet bislang nicht statt. In der Hoffnung, die zur Konvergenz der jeweiligen Methode erforderliche PE-Bedingung zu erfüllen, sind in der Literatur unterschiedliche heuristische Strategien zu finden. So wird häufig Explorationsrauschen in Form von Zufallssignalen auf die Stellgröße addiert, wie beispielsweise Gaußsches weißes Rauschen oder, seltener, auf einem Intervall gleichverteiltes Rauschen. Häufig ist dieses Anregungssignal jedoch nicht näher spezifiziert und die Erfüllung der Anregungsbedingung wird nicht weiter untersucht. Eine ähnlich verbreitete Methode ist die Addition eines Anregungssignals, das aus einer Summation und Multiplikation von Sinus- und Kosinusfunktionen resultiert³², auf die Stellgröße. Eine Untersuchung, wie diese Sinus- und Kosinussignale zu wählen sind, sodass eine geeignete Anregung erreicht wird, findet jedoch bislang nicht statt. Schließlich schlagen manche Autorinnen und Autoren als dritte Strategie vor, das System für verschiedene Zustände auszuwerten, beispielsweise durch eine wiederholte Neuinitialisierung des Anfangszustands. Auch hier findet keine theoretische Untersuchung statt und es bleibt ungeklärt, wie eine solche Neuinitialisierung konkret erfolgen müsste, um die Erfüllung der PE-Bedingung zu gewährleisten. Beispiele für ADP-Ansätze, die diese drei Heuristiken zur Anregung verwenden, sind Tabelle 2.3 zu entnehmen.

Bei der Betrachtung der Anregungssignale für ADP-Methoden ist schließlich noch hervorzuheben, dass mögliche anwendungsspezifische Anforderungen bislang, insbesondere bei Verwendung weißen Rauschens und anderer Zufallssignale zur Anregung, nicht berücksichtigt werden. Beispielsweise weisen reale physikalische und biologische Systeme üblicherweise Tiefpassverhalten auf [EP10], hochfrequente Stellsignale werden also gedämpft und sind somit

³² Da sich Produkte von Sinus- und Kosinusfunktionen stets durch Summen aus Sinus- und Kosinusfunktionen darstellen lassen, wird im Folgenden nur noch von Summationen gesprochen.

additives Zufallssignal auf Stellgröße:

[ATLAK07], [BK18], [DJ09], [DJ10], [DJ11], [FY16], [JKBD15], [KL15], [KLM⁺14], [KLNSK15], [LCL⁺19], [LLW14], [LP03], [LPC12], [LVV12], [QZLY19], [Vam17], [VMKL17], [WLMZ18], [WLZZ16], [WY18], [ZCL13], [ZDJ15], [ZQJL14], [ZZWZ16]

additiver Summenterm aus Sinus- und Kosinusfunktionen auf Stellgröße:

[BA18], [BKJ⁺13], [FC16], [GJ16], [JJ14b], [JJ14c], [JJ12], [JKBD15], [KDBD15], [KLJ17], [KLL15], [LYD17], [LYWW15], [LNY⁺15], [MRLP16], [ML14b], [QZLY19], [TCTH19], [VMKL17], [VMH16], [Vam16], [Vam15], [WLL14], [YHL13], [YLLL16]

wiederholte Neuinitialisierung des Systemzustands:

[DLL⁺19], [LLW14], [VL09], [VPAKL09], [WL12]

Tabelle 2.3: Anregungsheuristiken verschiedener ADP-Ansätze.

potenziell ungeeignet zur Anregung. Zudem ist zu beachten, dass die Verwendung ungeeigneten Explorationsrauschens eine erhöhte mechanische Belastung technischer Systeme mit sich bringt und zur schnelleren Abnutzung oder sogar Beschädigung des Systems führen kann (vgl. [KBP13], [dKTB18]). So beansprucht die Verwendung weißen Rauschens durch ruckartige Stelleingriffe potenziell die Aktuatorik. Schließlich ist eine solch ruckartige Anregung für Adaptionsverfahren, die sich online zur Laufzeit anpassen sollen, meist unerwünscht und beispielsweise auch in einer möglichen zukünftigen Anwendung im Mensch-Maschine-Kontext weder komfortabel noch zur Anregung unbekannter Reaktionen des Menschen geeignet³³.

Zusammenfassend kann konstatiert werden, dass zahlreiche Publikationen zu ADP-basierten Verfahren die Erfüllung geeigneter Anregungsbedingungen voraussetzen, bis dato jedoch noch keine verifizierbare Methode existiert, welche die Einhaltung der PE-Bedingung im nichtlinearen Kontext sicherstellt (vgl. [JZLH19], [JKBD15], [KLL15], [BKJ⁺13]). Dies ist insbesondere auf den Mangel an allgemeinen theoretischen Aussagen zur PE-Bedingung unter nichtlinearen Signaltransformationen, die im Bereich des ADP aufgrund nichtlinearer Gütefunktionale selbst bei LQ-Optimierungsproblemen auftreten, zurückzuführen. Im Hinblick auf eine spätere Übertragbarkeit auf reale Anwendungsbeispiele sind zudem problemspezifische

³³ Im Kontext der Mensch-Maschine-Interaktion, bei der sich Mensch und Maschine gegenseitig bei einer Regelungsaufgabe unterstützen, kann der Mensch aufbauend auf Arbeiten zur Modellierung des sensomotorischen Systems [Sco04], [Tod04], [JJ14a] als Optimalregler modelliert werden. Im Zusammenspiel mit einer ebenfalls optimierungstheoretisch formulierten Automation ergibt sich ein Differenzialspiel (vgl. auch Definition 5.1) [FFH17], [NC15]. Da üblicherweise sowohl das Regelgesetz als auch das Gütefunktional des menschlichen Handlungspartners unbekannt ist, muss entweder das Regelgesetz (vgl. beispielsweise [KNFH20]) oder das Gütefunktional identifiziert werden (vgl. [KIR⁺17], [IKFH17], [IBKH20], [RIK⁺17]), oder die Automation adaptiert sich mithilfe eines ADP-basierten Ansatzes aus der Interaktion mit dem Menschen. Da jedoch das menschliche neuromuskuläre System Tiefpasscharakteristik aufweist [KAV⁺14], sind hochdynamische Anregungssignale hierbei ungeeignet, um adäquate Informationen zur Adaption der Automation zu erhalten. Eine sinnvoll gewählte Anregung ist also insbesondere auch für zukünftige adaptive Mensch-Maschine-Konzepte erforderlich.

Anforderungen zu beachten, wie beispielsweise Freiheitsgrade beim Anregungsdesign zur Berücksichtigung von Tiefpasseigenschaften.

2.4 Wissenschaftliche Fragestellungen und Beiträge der Arbeit

Ausgehend von den Forschungslücken, die in den vorigen Abschnitten aufgezeigt wurden, werden im Folgenden zwei konkrete Fragestellungen abgeleitet, zu denen die vorliegende Arbeit Beiträge liefert. Außerdem wird eine Übersicht über diese wissenschaftlichen Beiträge gegeben.

2.4.1 ADP-kompatible, flexible Solltrajektorienendarstellung

ADP-basierte Regelungsansätze sind üblicherweise auf der Bellman-Gleichung bzw. Hamilton-Jacobi-Bellman-Gleichung begründet und für einen unendlichen Optimierungshorizont formuliert. Eine wesentliche Herausforderung beim Entwurf von ADP-Solltrajektorienfolgereglern ist eine mit der Bellman- bzw. HJB-Gleichung kompatible Repräsentation der Solltrajektorie. Bisher existieren aber weder eine allgemeine theoretische Definition ADP-kompatibler Solltrajektorien Darstellungen noch geeignete Ansätze ADP-basierter Solltrajektorienfolgeregler, welche die Einbeziehung einer von außen vorgebbaren Darstellung variabler Solltrajektorienverläufe ermöglichen. Daraus ergibt sich die folgende Fragestellung:

Forschungsfrage 1:

Wie können Solltrajektorienverläufe dargestellt werden, um in einen modellfreien ADP-Formalismus integrierbar zu sein, der einen approximierten optimalen Solltrajektorienfolgeregler mithilfe von Messdaten eines dynamischen Systems und ohne Verwendung eines Systemmodells erlernt?

Beitrag 1:

In dieser Arbeit wird erstmalig der Begriff der ADP-kompatiblen Solltrajektorienendarstellung eingeführt und formal definiert. Weiterhin werden neuartige und mit dem ADP-Formalismus kompatible Darstellungsformen flexibler, von außen vorgegebbarer Referenztrajektorien für Sollzustände präsentiert sowie deren Eigenschaften analysiert. Insbesondere beinhaltet die vorliegende Arbeit:

- Theoretische Beiträge in zeitdiskreter Darstellung (Kapitel 3):
 - Eine neuartige, mit dem ADP-Formalismus kompatible, flexible, parametrisierte Darstellungsform der Referenztrajektorie, die insbesondere die Beschreibung des Solltrajektorienverlaufs in einer lokalen Umgebung des aktuellen Zeitpunkts erlaubt. Des Weiteren werden Existenz und Eindeutigkeit des optimalen Regelgesetzes für den LQ-Solltrajektorienfolgeregelungsfall analysiert. Zudem wird ein neues Stabilitätskriterium präsentiert. Schließlich wird ein Vorgehen vorgestellt, das es ermöglicht, ADP-basierte Regler, die diese parametrisierte Referenztrajektorie nutzen, anhand von Messdaten zu trainieren.
 - Die im ADP-Kontext erstmalige direkte Verwendung der Sollzustände auf einem endlichen, gleitenden Vorausschauhorizont. Hierbei wird für den LQ-Fall die Existenz und Eindeutigkeit sowie die exakte Form der optimalen Lösung analysiert. Letzteres ist essenziell, um effiziente Funktionsapproximatoren mit einer möglichst geringen Anzahl zu schätzender Gewichte definieren und mithilfe von Messdaten trainieren zu können. Außerdem wird die Konvergenz der verwendeten ADP-Methode für den Solltrajektorienfolgeregelungsfall bewiesen.
- Theoretische Beiträge in zeitkontinuierlicher Darstellung (Kapitel 4):
 - Eine neuartige, mit dem ADP-Formalismus kompatible, flexible, parametrisierte Repräsentation der Referenztrajektorie, die sich aus der Superposition der Lösungen gewöhnlicher Differenzialgleichungen ergibt. Existenz, Eindeutigkeit und Stabilität der optimalen Lösung werden analysiert. Schließlich wird ein Konzept vorgestellt, das diese Solltrajektorien-darstellung in eine ADP-Methode integriert.
- Reale Anwendungsbeispiele eines zuvor entwickelten ADP-basierten Solltrajektorienfolgereglers (Kapitel 6):
 - Ein Geschwindigkeitsregler für ein Realfahrzeug, dessen Abhängigkeit von einer parametrisierten Beschreibung des Sollgeschwindigkeitsverlaufs direkt aus Messdaten erlernt wird. Insbesondere findet die Adaption hierbei online statt, d. h. die Reglergewichte werden während der Fahrt im geschlossenen Regelkreis angepasst.
 - Ein ADP-basierter Solltrajektorienfolgeregler für ein reales Ball-auf-Platte-System. Dieser wird einerseits mit einem ADP-basierten Sollzustandsregler (Festwertregelung) und andererseits mit einem modellbasierten optimalen Reglerentwurf verglichen.

2.4.2 Konvergenz ADP-basierter Regelungsansätze

Bei sämtlichen adaptiven und ADP-basierten Regelungsverfahren wird eine ausreichende Anregung des Systems bzw. der für die Adaption relevanten Signale gefordert, um Konvergenz der zu schätzenden Parameter gegen die jeweiligen Zielgrößen zu gewährleisten. Die Signale, auf denen diese Adaption im ADP-Kontext beruht und für welche die Erfüllung der

PE-Eigenschaft gefordert wird, resultieren aus nichtlinearen Transformationen der Systemzustände (vgl. Abschnitt 2.3), woraus die folgende Forschungsfrage resultiert:

Forschungsfrage 2:

Welche Bedingungen an die dynamischen Zustände eines Systems gewährleisten, dass die PE-Eigenschaft eingehalten wird und ADP-Methoden konvergieren?

Beitrag 2:

In Kapitel 5 dieser Arbeit wird eine weitverbreitete PE-Eigenschaft für zeitkontinuierliche ADP-Ansätze (vgl. (2.27)) analysiert. Hierbei werden eingangsaffine Nicht-Nullsummen-Differenzialspiele, die eine Generalisierung eingangsaffiner Optimalregelungsprobleme darstellen, betrachtet. Die Einhaltung der PE-Eigenschaft garantiert die Konvergenz der Critic-Gewichte bei Adaption mit einem Gradientenabstiegsverfahren. Wie Abschnitt 2.3 motiviert, spielen die in einem Signal vorhandenen Frequenzen eine wesentliche Rolle bei der Erfüllung der PE-Eigenschaft. Daher wird in der vorliegenden Arbeit untersucht, wie der Frequenzgehalt des Systemzustands durch die Verwendung polynomieller Basisfunktionen zur Critic-Funktionsapproximation beeinflusst wird. Der zentrale Beitrag ist schließlich durch neuartige hinreichende Frequenzbedingungen an den Systemzustand gegeben, welche gewährleisten, dass die für die Konvergenz benötigte PE-Eigenschaft erfüllt ist. Die präsentierten Bedingungen weisen zudem Freiheitsgrade bezüglich der für die Anregung verwendeten Frequenzen und Amplituden auf. Somit können anwendungsspezifische Anforderungen berücksichtigt werden. Ausgehend von den theoretischen Erkenntnissen werden Simulationsergebnisse präsentiert und diskutiert.

3 Zeitdiskrete ADP-basierte Solltrajektorienfolgeregelung

Dieses Kapitel beantwortet die Forschungsfrage, wie Solltrajektorienverläufe dargestellt werden können, um in zeitdiskrete, modellfreie ADP-Ansätze integrierbar zu sein (vgl. Forschungsfrage 1). Hierzu wird in Abschnitt 3.1 zunächst definiert, wodurch eine ADP-kompatible Solltrajektorienendarstellung charakterisiert ist. Anschließend werden zwei Ansätze präsentiert: In Abschnitt 3.2 wird eine zeitdiskrete, parametrisierte ADP-kompatible Solltrajektorienendarstellung vorgestellt, anschließend betrachtet Abschnitt 3.3 die Verwendung beliebiger Sollzustände auf einem endlichen Vorausschauhorizont. Ein abschließendes Resümee rundet das Kapitel ab.

3.1 Definition ADP-kompatibler zeitdiskreter Trajektorien

In diesem Abschnitt werden allgemeine Vorüberlegungen zur Funktionsweise von ADP-Methoden vorgenommen, um darauf aufbauend Anforderungen an den Entwurf ADP-basierter Trajektorienfolgeregler abzuleiten. Hierzu wird zunächst exemplarisch die zeitdiskrete Value Function

$$V^\mu(\mathbf{x}_k) = \sum_{\kappa=k}^{\infty} \gamma^{\kappa-k} r(\mathbf{x}_\kappa, \boldsymbol{\mu}(\mathbf{x}_\kappa)) \quad (3.1a)$$

$$= r(\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k)) + \gamma V^\mu(\mathbf{x}_{k+1}) \quad (3.1b)$$

$$= r(\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k)) + \gamma V^\mu(\mathbf{f}(\mathbf{x}_k) + \mathbf{g}(\mathbf{x}_k)\boldsymbol{\mu}(\mathbf{x}_k)) \quad (3.1c)$$

(vgl. (2.3a)) ohne Verwendung einer Solltrajektorie betrachtet. Jedoch lassen sich die hier gewonnenen Erkenntnisse auch auf Formulierungen unter Verwendung einer Q-Function (vgl. (2.7)) oder auf zeitkontinuierliche Darstellungen (vgl. (2.13) und (2.20b)) übertragen. Eine zentrale Eigenschaft der Value Function $V^\mu(\mathbf{x}_k)$ ist, dass eine funktionale Abhängigkeit vom aktuellen Zustand \mathbf{x}_k die (ggf. diskontierten) Gesamtkosten repräsentieren kann, die über einen unendlichen Zeithorizont anfallen. Die Betrachtung der Value Function in (3.1) offenbart, dass dies maßgeblich der Markov-Eigenschaft des Systemzustands \mathbf{x}_k zu verdanken ist, der aus der Dynamik

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k) + \mathbf{g}(\mathbf{x}_k)\boldsymbol{\mu}(\mathbf{x}_k) \quad (3.2)$$

des geschlossenen Regelkreises resultiert (vgl. [van12], [SLH⁺14]). Diese Repräsentation der Gesamtkosten in Form der Value Function $V^\mu(\mathbf{x}_k)$ ist mit dem betrachteten System (2.1),

dem mit der Value Function assoziierten Regelgesetz $\boldsymbol{\mu}(\mathbf{x}_k)$ und der Einschrittkostenfunktion $r(\mathbf{x}_k, \mathbf{u}_k)$ verknüpft. Ohne diese Markov-Eigenschaft würde der Systemzustand \mathbf{x}_k nicht als alleiniger, expliziter Parameter der Value Function $V^\mu(\mathbf{x}_k)$ genügen, um die Kosten über einen unendlichen Zeithorizont korrekt abbilden zu können, da aus dem aktuellen Zustand \mathbf{x}_k der weitere Zustandsverlauf, der implizit in $V^\mu(\mathbf{x}_k)$ enthalten ist, nicht bestimmt werden könnte³⁴.

Aufgrund dieser Markov-Eigenschaft ist es basierend auf (3.1b) zudem grundsätzlich möglich, mithilfe von Datentupeln

$$\mathcal{T}_k := \{\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k), r(\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k)), \mathbf{x}_{k+1}\} \quad (3.3)$$

(und ggf. γ) die Value Function $V^\mu(\mathbf{x}_k)$, bzw. bei Verwendung eines Funktionsapproximators (vgl. (2.21)) das Critic-Gewicht \mathbf{w} , zu bestimmen³⁵ (vgl. beispielsweise den Policy-Evaluation-Schritt der PI (2.22) oder der VI (2.25)).

Um im Folgenden Solltrajektorienverläufe $\mathbf{x}_{r,k}, \mathbf{x}_{r,k+1}, \dots$ in das Gütefunktional J und somit die Value Function V^μ zu integrieren, hängen die Einschrittkosten r neben dem Systemzustand \mathbf{x}_k und der Stellgröße \mathbf{u}_k im Allgemeinen zusätzlich vom Sollzustand $\mathbf{x}_{r,k}$ ab und es ergibt sich $r(\mathbf{x}_k, \mathbf{x}_{r,k}, \mathbf{u}_k)$. Würde eine Abhängigkeit von der Solltrajektorie $\mathbf{x}_{r,k}, \mathbf{x}_{r,k+1}, \dots$ nicht explizit in der Value Function verwendet und weiterhin eine Funktion der Form $V^\mu(\mathbf{x}_k)$ zugrunde gelegt, wie beispielsweise in der Arbeit von Yu et al. [YSH⁺17] (vgl. Abschnitt 2.2), so müsste diese Value Function $\forall \mathbf{x}_k$ die Gleichung

$$V^\mu(\mathbf{x}_k) = \sum_{\kappa=k}^{\infty} \gamma^{\kappa-k} r(\mathbf{x}_\kappa, \mathbf{x}_{r,\kappa}, \boldsymbol{\mu}(\mathbf{x}_\kappa)) = r(\mathbf{x}_k, \mathbf{x}_{r,k}, \boldsymbol{\mu}(\mathbf{x}_k)) + \gamma V^\mu(\mathbf{x}_{k+1}) \quad (3.4)$$

erfüllen (vgl. (3.1)). Dies wäre jedoch bei beliebigen Sollzuständen $\mathbf{x}_{r,k}$ nicht möglich, da eine Abhängigkeit der Sollzustände nicht von der Value Function $V^\mu(\mathbf{x}_k)$ erfasst werden könnte. Für den sehr eingeschränkten Spezialfall, dass $\mathbf{x}_{r,k}$ stets eindeutig durch eine zeitinvariante Abbildung $\mathbf{f}_{\mathbf{x}_r, \mathbf{x}} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ durch $\mathbf{x}_{r,k} = \mathbf{f}_{\mathbf{x}_r, \mathbf{x}}(\mathbf{x}_k)$ mit dem Systemzustand \mathbf{x}_k verknüpft ist, ließe sich eine Value Function $V^\mu(\mathbf{x}_k)$ definieren, die (3.4) erfüllt. Die damit gelernte Value Function wäre jedoch implizit nur für diesen direkt mit \mathbf{x}_k korrelierten Sollzustandsverlauf gültig und in der Anwendung kaum von Nutzen.

Den Sollzustand nur für den aktuellen und $n_h < \infty$ weitere Zeitschritte vorzugeben, wäre hingegen mit einem Optimierungsproblem mit endlichem Zeithorizont verknüpft, da der

³⁴ Ebenso ist es für klassische Methoden des RL, die das Lösen Markovscher Entscheidungsprozesse [Bel57b] anstreben, erforderlich, dass der Systemzustand \mathbf{x}_k die Markov-Eigenschaft erfüllt [SB18, S. 49], [KT03]. Bei Betrachtung stochastischer Markov-Entscheidungsprozesse ist somit die Zustandsübergangswahrscheinlichkeit nur vom aktuellen, nicht jedoch von vergangenen Zuständen abhängig.

³⁵ Da es sich bei (3.1b) um eine skalare Gleichung handelt, sind zur Schätzung eines h -dimensionalen Critic-Gewichts mindestens h unabhängige Datentupel erforderlich (vgl. beispielsweise [LVS12, S. 495]). Insbesondere kann der Critic unter Einbeziehung genügend vieler und entsprechend vielseitiger Datentupel \mathcal{T}_k (d. h. bei ausreichender Anregung des Systems, vgl. Abschnitt 2.3) aus den während des Trainingsprozesses betrachteten Zustandsübergängen $\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k), \mathbf{x}_{k+1}$ implizit die Systemdynamik und über die Einschrittkosten $r(\mathbf{x}_k, \boldsymbol{\mu}(\mathbf{x}_k))$ die zugrunde liegenden Gesamtkosten erlernen.

weitere Verlauf der Solltrajektorie $\mathbf{x}_{r,\kappa}$ für $\kappa > k + n_h$ unbekannt wäre³⁶. Die Verwendung eines unendlichen Zeithorizontes ist jedoch nach Bemerkung 2.1 dafür verantwortlich, dass in (3.1) auf beiden Seiten der Gleichung dieselbe Value Function $V^\mu(\cdot)$ auftaucht und diese somit durch denselben Funktionsapproximator beschrieben werden kann³⁷.

Würde eine Value Function eines Optimierungsproblems mit unendlichem Optimierungshorizont gesucht, die einen komplett beliebigen Solltrajektorienverlauf berücksichtigt, so müsste diese eine Abhängigkeit beliebiger Sollzustände $\mathbf{x}_{r,k}, \mathbf{x}_{r,k+1}, \dots$ bis ins Unendliche aufweisen, also durch

$$\begin{aligned} V^\mu(\mathbf{x}_k, \mathbf{x}_{r,k}, \mathbf{x}_{r,k+1}, \dots) &= \sum_{\kappa=k}^{\infty} \gamma^{\kappa-k} r(\mathbf{x}_\kappa, \mathbf{x}_{r,\kappa}, \boldsymbol{\mu}(\mathbf{x}_\kappa, \mathbf{x}_{r,\kappa}, \mathbf{x}_{r,\kappa+1}, \dots)) \\ &= r(\mathbf{x}_k, \mathbf{x}_{r,k}, \boldsymbol{\mu}(\mathbf{x}_k, \mathbf{x}_{r,k}, \mathbf{x}_{r,k+1}, \dots)) \\ &\quad + \gamma V^\mu(\mathbf{x}_{k+1}, \mathbf{x}_{r,k+1}, \mathbf{x}_{r,k+2}, \dots) \end{aligned} \quad (3.5)$$

definiert sein (vgl. [KPRH20]). Eine Value Function dieser Form wäre jedoch weder aus theoretischer noch aus praktischer Sicht handhabbar: Weder lassen sich unendlich viele Eingabeparameter verwenden, noch ist in den meisten praktischen Anwendungen der Sollverlauf (wie beispielsweise die Sollgeschwindigkeit oder Sollposition eines Fahrzeugs) bis ins Unendliche bekannt. Auch würde eine Value Function mit unendlich vielen beliebigen Referenzparametern mit einem unendlichdimensionalen Funktionsapproximator (vgl. Abschnitt 2.1.3) korrelieren.

Somit offenbart sich, dass der Verlauf der Solltrajektorie zwar als explizite Abhängigkeit in die Value Function aufgenommen werden muss, dies aber mithilfe einer endlichdimensionalen Repräsentation erfolgen sollte, um handhabbar zu sein. Sei diese endlichdimensionale Repräsentation im Zeitschritt k allgemein durch einen n_ζ -dimensionalen ($n_\zeta < \infty$) Parameter $\zeta_k \in \mathcal{Z} \subset \mathbb{R}^{n_\zeta}$ mit Definitionsbereich \mathcal{Z} und $\mathbf{x}_{r,k} = \mathbf{f}_{\mathbf{x}_r, \zeta}(\zeta_k)$ gegeben, wobei $\mathbf{f}_{\mathbf{x}_r, \zeta} : \mathcal{Z} \rightarrow \mathcal{X}$ eine zeitinvariante Abbildung darstellt. Damit nun ζ_k genügt, um den zeitlichen Verlauf von $\mathbf{x}_{r,\kappa}, \forall \kappa \geq k$, erfassen zu können, muss zudem für ζ_k die Markov-Eigenschaft gelten, also der Parameter $\zeta_{k+1} = \mathbf{f}_\zeta(\zeta_k)$ mit $\mathbf{f}_\zeta : \mathcal{Z} \rightarrow \mathcal{Z}$ zum Folgezeitpunkt lediglich vom aktuellen Parameter ζ_k abhängen³⁸. Aus praktischer Sicht lässt sich durch den Parameter ζ_k somit in jedem Zeitschritt k Information über den aktuellen Solltrajektorienverlauf im ADP-Mechanismus berücksichtigen. Beim Entwurf von $\mathbf{f}_{\mathbf{x}_r, \zeta}(\cdot)$ und $\mathbf{f}_\zeta(\cdot)$ kann, abhängig von der Anwendung und den Anforderungen, nicht zuletzt durch die Wahl der Dimension n_ζ des Parameters ζ ein Kompromiss zwischen der Komplexität und der Kompaktheit der

³⁶ Wird $\forall \kappa > k + n_h$ ein konstanter Sollzustand angenommen, so lässt sich hingegen wieder ein Optimierungsproblem mit unendlichem Optimierungshorizont formulieren. Dies wird in Abschnitt 3.3 vorgestellt.

³⁷ Die Verwendung desselben Funktionsapproximators auf beiden Seiten der Bellman-Gleichung würde bei einem Optimierungsproblem mit endlichem Optimierungshorizont zu einem systematischen Fehler führen. Bei sehr langem aber endlichem Optimierungshorizont oder bei einer starken Diskontierung ($\gamma \ll 1$) würde dieser systematische Fehler zwar reduziert, jedoch führt ersteres zu einem hochdimensionalen und zudem schwierig trainierbaren Funktionsapproximator und letzteres verfälscht einerseits ggf. das eigentliche Optimierungsziel und kann andererseits ein instabiles Gesamtsystem begünstigen (vgl. Abschnitt 4.2.2).

³⁸ Bzw. in zeitkontinuierlichen Fall die Änderungsrate $\dot{\zeta}(t)$ lediglich von $\zeta(t)$ abhängen (vgl. Abschnitt 4.1).

Solltrajektorienendarstellung eingegangen werden. Zusammenfassend lässt sich eine mit dem ADP-Formalismus kompatible Solltrajektorienendarstellung wie folgt definieren:

Definition 3.1 (Zeitdiskrete ADP-kompatible Solltrajektorienendarstellung)

Eine mit dem ADP-Formalismus kompatible Solltrajektorienendarstellung mit dem Sollzustand $\mathbf{x}_{\tau,k} \in \mathcal{X}$ ist durch zeitinvariante Funktionen $\mathbf{f}_{\mathbf{x}_{\tau},\zeta} : \mathcal{Z} \rightarrow \mathcal{X}$ und $\mathbf{f}_{\zeta} : \mathcal{Z} \rightarrow \mathcal{Z}$ mit

$$\mathbf{x}_{\tau,k} = \mathbf{f}_{\mathbf{x}_{\tau},\zeta}(\zeta_k) \quad (3.6)$$

und

$$\zeta_{k+1} = \mathbf{f}_{\zeta}(\zeta_k) \quad (3.7)$$

charakterisiert, wobei $\zeta \in \mathcal{Z} \subset \mathbb{R}^{n_{\zeta}}$, $n_{\zeta} < \infty$.

Basierend auf Definition 3.1 lässt sich die folgende formale Aussage formulieren³⁹.

Proposition 3.1

Die Solltrajektorienendarstellung sei nach Definition 3.1 ADP-kompatibel. Für ein durch die zeitdiskrete Systemdynamik (2.1) beschriebenes System gilt dann:

1. Sind die durch

$$\sum_{\kappa=k}^{\infty} \gamma^{\kappa-k} r(\mathbf{x}_{\kappa}, \mathbf{x}_{\tau,\kappa}, \boldsymbol{\mu}(\mathbf{x}_{\kappa}, \zeta_{\kappa})), \quad (3.8)$$

$0 < \gamma \leq 1$, gegebenen Gesamtkosten (vgl. (2.3a)), die von der Sollzustandstrajektorie $\mathbf{x}_{\tau,k}, \mathbf{x}_{\tau,k+1}, \dots$ abhängen, endlich⁴⁰, dann können sie durch eine Value Function der Form $V^{\mu}(\mathbf{x}_k, \zeta_k)$ beschrieben werden.

2. Sind die durch

$$r(\mathbf{x}_k, \mathbf{x}_{\tau,k}, \mathbf{u}_k) + \sum_{\kappa=k+1}^{\infty} \gamma^{\kappa-k} r(\mathbf{x}_{\kappa}, \mathbf{x}_{\tau,\kappa}, \boldsymbol{\mu}(\mathbf{x}_{\kappa}, \zeta_{\kappa})), \quad (3.9)$$

$0 < \gamma \leq 1$, gegebenen Gesamtkosten (vgl. (2.7)), welche von der Sollzustandstrajektorie $\mathbf{x}_{\tau,k}, \mathbf{x}_{\tau,k+1}, \dots$ abhängen, endlich, dann können sie durch eine Q-Function der Form $Q^{\mu}(\mathbf{x}_k, \zeta_k, \mathbf{u}_k)$ beschrieben werden.

³⁹ Ebenso gilt die Aussage von Proposition 3.1 auch für allgemeinere nichtlineare Systeme der Form $\mathbf{x}_{k+1} = \mathbf{F}(\mathbf{x}_k, \mathbf{u}_k)$, aus Gründen der Einheitlichkeit wurde hier jedoch der eingangsaflne Fall formuliert.

⁴⁰ Für sinnvoll gestellte Probleme schließt dies neben der Stabilisierbarkeit des Systems und einem zulässigen, d. h. insbesondere stabilisierenden, Regelgesetz $\boldsymbol{\mu}$ auch ein, dass durch die Wahl der Diskontierung γ bzw. des Solltrajektorienverlaufs die Gesamtkosten endlich bleiben.

Beweis:

Unter Verwendung von (2.1), (3.6) und (3.7) folgt

1. aus (3.8)

$$V^\mu(\mathbf{x}_k, \zeta_k) := \sum_{\kappa=k}^{\infty} \gamma^{\kappa-k} r(\mathbf{x}_\kappa, \mathbf{x}_{\tau,\kappa}, \boldsymbol{\mu}(\mathbf{x}_\kappa, \zeta_\kappa)) \quad (3.10a)$$

$$= r(\mathbf{x}_k, \mathbf{x}_{\tau,k}, \boldsymbol{\mu}(\mathbf{x}_k, \mathbf{x}_{\tau,k})) + \gamma V^\mu(\mathbf{x}_{k+1}, \zeta_{k+1}) \quad (3.10b)$$

$$= r(\mathbf{x}_k, \mathbf{f}_{\mathbf{x}_\tau, \zeta}(\zeta_k), \boldsymbol{\mu}(\mathbf{x}_k, \mathbf{f}_{\mathbf{x}_\tau, \zeta}(\zeta_k))) \\ + \gamma V^\mu(\mathbf{f}(\mathbf{x}_k) + \mathbf{g}(\mathbf{x}_k) \boldsymbol{\mu}(\mathbf{x}_k, \mathbf{f}_{\mathbf{x}_\tau, \zeta}(\zeta_k)), \mathbf{f}_\zeta(\zeta_k)) \quad (3.10c)$$

sowie

2. basierend auf (3.9) und analog zu (2.7) mit (3.10c)

$$Q^\mu(\mathbf{x}_k, \zeta_k, \mathbf{u}_k) := r(\mathbf{x}_k, \mathbf{x}_\tau, \mathbf{u}_k) + \gamma V^\mu(\mathbf{x}_{k+1}, \zeta_{k+1}) \quad (3.11a)$$

$$= r(\mathbf{x}_k, \mathbf{f}_{\mathbf{x}_\tau, \zeta}(\zeta_k), \mathbf{u}_k) \\ + \gamma V^\mu(\mathbf{f}(\mathbf{x}_k) + \mathbf{g}(\mathbf{x}_k) \boldsymbol{\mu}(\mathbf{x}_k, \mathbf{f}_{\mathbf{x}_\tau, \zeta}(\zeta_k)), \mathbf{f}_\zeta(\zeta_k)) \quad (3.11b)$$

$$= r(\mathbf{x}_k, \mathbf{f}_{\mathbf{x}_\tau, \zeta}(\zeta_k), \mathbf{u}_k) \\ + \gamma Q^\mu(\mathbf{f}(\mathbf{x}_k) + \mathbf{g}(\mathbf{x}_k) \boldsymbol{\mu}(\mathbf{x}_k, \mathbf{f}_{\mathbf{x}_\tau, \zeta}(\zeta_k)), \mathbf{f}_\zeta(\zeta_k), \boldsymbol{\mu}(\mathbf{x}_k, \mathbf{f}_{\mathbf{x}_\tau, \zeta}(\zeta_k))). \quad (3.11c)$$

□

Eine ADP-kompatible Solltrajektorienarstellung stellt somit sicher, dass die Repräsentation der Gesamtkosten in Form einer Value Function lediglich eine explizite Abhängigkeit von \mathbf{x}_k und ζ_k erfordert⁴¹. Zudem hängt eine Q-Function in diesem Fall nur explizit von \mathbf{x}_k , ζ_k und \mathbf{u}_k ab.

In der Arbeit von Yu et al. [YSH⁺17], bei der für jede neue Kombination aus Referenztrajektorie und Anfangszustand der Regler neu trainiert werden muss, ist die Kompatibilitätsforderung nach Definition 3.1 nicht erfüllt. Die Betrachtung der in den Abschnitten 2.2.2 und 2.2.3 diskutierten Methoden offenbart zwar, dass die dortigen zeitdiskreten Ansätze die Anforderung der ADP-Kompatibilität erfüllen⁴². Den existierenden Methoden gelingt es hierbei jedoch nicht, zu

⁴¹ An dieser Stelle sei noch zu bemerken, dass auch Referenztrajektorienparameter, deren Folgezustand aus einer endlichen Anzahl vergangener Parameter hervorgeht, nach Definition 3.1 kompatibel sind. Um dies zu verdeutlichen, werde exemplarisch $\bar{\zeta}_{k+1} = \bar{\mathbf{f}}(\bar{\zeta}_k, \bar{\zeta}_{k-1})$ mit $\bar{\zeta}_k \in \bar{\mathcal{Z}} \subset \mathbb{R}^{n_\zeta}$ betrachtet. Für den Referenzparameter $\zeta_k := [\bar{\zeta}_k^\top \quad \bar{\zeta}_{k-1}^\top]^\top$ gilt mit $\mathbf{f}(\zeta_k) := \bar{\mathbf{f}}(\bar{\zeta}_k, \bar{\zeta}_{k-1})$ aufgrund von $\zeta_{k+1} = [\bar{\mathbf{f}}^\top(\bar{\zeta}_k, \bar{\zeta}_{k-1}) \quad \bar{\zeta}_k^\top]^\top = [\mathbf{f}^\top(\zeta_k) \quad ([\mathbf{I}_{n_\zeta} \quad \mathbf{0}_{n_\zeta \times n_\zeta}] \zeta_k)^\top]^\top =: \mathbf{f}_\zeta(\zeta_k)$ die Markov-Eigenschaft.

⁴² Bei Ansätzen, bei denen ein stationärer Sollzustand vorgegeben wird, sind sowohl $\mathbf{f}_{\mathbf{x}_\tau, \zeta}(\cdot)$ als auch $\mathbf{f}_\zeta(\cdot)$ Identitätsabbildungen.

generalisieren, indem eine flexible, veränderliche, externe Vorgabe des Solltrajektorienverlaufs, wie in Abbildung 2.4 visualisiert, berücksichtigt werden kann.

Bemerkung 3.1

Zwar wird im Zeitschritt k angenommen, dass die durch ζ_k , $f_{x_r, \zeta}(\cdot)$ und $f_\zeta(\cdot)$ beschriebene Solltrajektorie bis ins Unendliche fortgesetzt wird, jedoch kann aufgrund der expliziten Abhängigkeit der Value Function von ζ_k dieser Parameter, der den Verlauf der Solltrajektorie definiert, prinzipiell jederzeit geändert werden. Wenngleich in jedem Zeitschritt k ein unendlicher Optimierungshorizont zugrunde gelegt wird, kann somit durch eine Anpassung von ζ_k eine lokale Beschreibung der Solltrajektorie, d. h. ein im Zeitschritt k auf den lokalen Verlauf der Referenztrajektorie angepasster Parameter ζ_k , verwendet werden. Jedoch muss bei einem datenbasierten Training eines solchen optimalen Trajektorienfolgereglers $\mu(x_k, \zeta_k)$ darauf geachtet werden, dass dem ADP-Algorithmus Datentupel

$$\mathcal{T}_k := \{x_k, \zeta_k, \mu(x_k, \zeta_k), r(x_k, \mu(x_k, \zeta_k)), x_{k+1}, \zeta_{k+1}\} \quad (3.12)$$

(vgl. (3.3)) präsentiert werden, welche nach Definition 3.1 zu kompatiblen Solltrajektorien gehören. Dies erlaubt, neben der impliziten Abhängigkeit der Value Function $V^\mu(x_k, \zeta_k)$ von der Diskontierung γ , der Systemdynamik $f(\cdot)$, $g(\cdot)$ und der Einschrittkostenfunktion $r(\cdot)$, auch die implizite Abhängigkeit vom Solltrajektorienverlauf, beschrieben durch $f_{x_r, \zeta}(\cdot)$ und $f_\zeta(\cdot)$, zu erlernen.

Die im vorliegenden Abschnitt erstmals formal definierte zeitdiskrete ADP-kompatible Solltrajektorienendarstellung liefert den ersten Beitrag zu der in Abschnitt 2.4.1 formulierten Forschungsfrage 1 und legt den Grundstein zur Entwicklung ADP-basierter Solltrajektorienfolgeregler.

Im weiteren Verlauf von Kapitel 3 werden Methoden vorgestellt, die extern vorgebbare Solltrajektorienverläufe durch nach Definition 3.1 kompatible Referenzverläufe entweder lokal approximieren oder auf einem endlichen Vorausschauhorizont nutzen und dabei, wie in Bemerkung 3.1 beschrieben, die Konsistenz der während des Adaptionvorgangs verwendeten Daten gewährleisten (vgl. Abbildung 3.1). Nachfolgend werden zwei zeitdiskrete ADP-basierte Solltrajektorienregler präsentiert. Die zentrale Idee beider Methoden ist dabei, eine Q-Funktion zu definieren, die explizit den zukünftigen Verlauf der Solltrajektorie berücksichtigt. Bei der in Abschnitt 3.2 vorgestellten Methode geschieht dies mithilfe von Parametern, welche den Verlauf repräsentieren, wohingegen der Ansatz in Abschnitt 3.3 die Sollzustände über einen gleitenden Vorausschauhorizont explizit in die Q-Funktion integrieren.

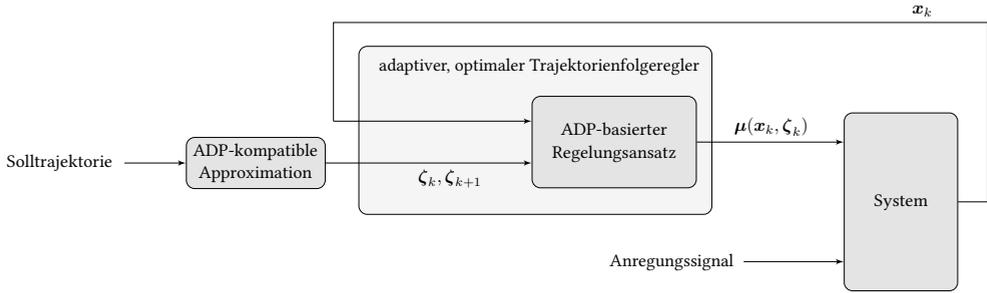


Abbildung 3.1: Aufgabe der ADP-kompatiblen Approximation einer Solltrajektorie ist es, aus einem im Zeitschritt k von außen vorgegebenen Referenztrajektorienverlauf eine nach Definition 3.1 kompatible Darstellung mit endlichdimensionalem Referenzparameter ζ_k zu erzeugen. Aus den Tupeln $\{\zeta_k, \zeta_{k+1}\}$, die einem ADP-basierten Regelungsansatz zur Verfügung gestellt werden, kann dieser implizit $f_{x_r, \zeta}(\cdot)$ und $f_\zeta(\cdot)$ und somit den lokal im Zeitschritt k approximierten Solltrajektorienverlauf berücksichtigen.

3.2 Zeitdiskrete ADP-kompatible parametrisierte Referenztrajektorie

Um im Zeitschritt k eine kompakte Darstellung des Einflusses des zukünftigen Solltrajektorienverlaufs auf die Gesamtkosten zu erhalten, wird im Folgenden eine Parametermatrix Z_k definiert, die Basisfunktionen $\rho(\cdot)$ gewichtet und den Referenztrajektorienverlauf x_r beschreibt. Diese Parametermatrix wird in die Q-Function $Q^\mu(x_k, Z_k, u_k)$ integriert, hierbei entspricht $\zeta_k = \text{vec}(Z_k)$ ⁴³. Dadurch repräsentiert die Q-Function explizit die Abhängigkeit der akkumulierten (ggf. diskontierten) Gesamtkosten von dem durch Z_k beschriebenen und jederzeit beeinflussbaren Solltrajektorienverlauf. Der dieser Q-Function zugehörige Regler hängt letztlich ebenfalls explizit von der durch Z_k parametrisierten Beschreibung der Referenztrajektorie ab, ein gelernter Regler muss bei einer Änderung des Solltrajektorienverlaufs also nicht erneut trainiert werden. Zudem wird im Optimierungsproblem nicht nur die aktuelle Abweichung vom Sollzustand, sondern der (parametrisierte) Sollzustandsverlauf, wie beispielsweise Straßen-, Geschwindigkeits-, Konzentrations- oder Temperaturprofile, berücksichtigt. Dies erlaubt eine flexible Vorausschau von ADP-Reglern und verspricht Vorteile gegenüber einer konstanten Sollwertvorgabe. Nachfolgend wird dieser Ansatz als *Parametrized Reference ADP* (PRADP) bezeichnet.⁴⁴

Nach Einführung dieser referenzabhängigen Q-Function können Funktionsapproximatoren verwendet werden, um basierend auf dem TD-Fehler sowohl die Q-Function als auch den optimalen Regler zu schätzen. Der im Folgenden präsentierte Ansatz stellt insbesondere aufgrund der speziellen Formulierung der solltrajektorienabhängigen Bellman-Gleichung unter Verwendung einer verschobenen Parametermatrix sicher, dass trotz beliebiger Vorgabe des

⁴³ Hierbei bezeichnet $\text{vec}(\cdot)$ die Vektorisierung einer Matrix, indem die Spalten vertikal konkateniert werden.

⁴⁴ Teile des vorliegenden Abschnitts wurden in [KRP⁺20] veröffentlicht.

Solltrajektorienparameters ADP-Kompatibilität innerhalb der Trainingsdaten gewährleistet ist. Für den wichtigen Fall linear-quadratischer optimaler Trajektorienfolgeregler wird die dem Trajektorienfolgeregelungsproblem zugrunde liegende Q-Function anschließend analysiert. Unter plausiblen Annahmen werden Existenz und Eindeutigkeit der optimalen Lösung des betrachteten Optimierungsproblems bewiesen. Zudem wird die Stabilität des geschlossenen Regelkreises für die optimale Lösung bei Verwendung einer diskontierten Gütefunktion untersucht und ein hinreichendes Stabilitätskriterium vorgestellt. Anschließend werden Simulationsergebnisse, welche die Flexibilität der neuartigen Methode offenbaren, präsentiert und diskutiert. Für die reale Anwendung der in diesem Kapitel entwickelten ADP-kompatiblen Solltrajektorienbeschreibung sei auf Kapitel 6 verwiesen.

3.2.1 Allgemeine Problemstellung

Betrachtet werde zunächst ein zeitdiskretes System mit unbekannter Systemdynamik⁴⁵

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k) + \mathbf{g}(\mathbf{x}_k)\mathbf{u}_k \quad (3.13)$$

mit dem diskreten Zeitschritt $k \in \mathbb{N}_{\geq 0}$, dem Systemzustand $\mathbf{x}_k \in \mathbb{R}^n$ und der Stellgröße $\mathbf{u}_k \in \mathbb{R}^p$. Das System (3.13) sei steuerbar auf der kompakten Menge $\mathcal{X} \subset \mathbb{R}^n$, die den Ursprung enthält (vgl. [KL15]). Weiterhin sei die durch die Parametermatrix

$$\mathbf{Z}_k = \begin{bmatrix} \mathbf{z}_{k,1}^\top \\ \mathbf{z}_{k,2}^\top \\ \vdots \\ \mathbf{z}_{k,n}^\top \end{bmatrix}, \quad (3.14)$$

$\mathbf{Z}_k \in \mathbb{R}^{n \times n_z}$, und gegebene Basisfunktionen $\boldsymbol{\rho}(\kappa) \in \mathbb{R}^{n_z}$ beschriebene Sollzustandstrajektorie $\mathbf{x}_r(\mathbf{Z}_k, \kappa) \in \mathbb{R}^n$ im Zeitschritt k durch

$$\mathbf{x}_r(\mathbf{Z}_k, \kappa) = \mathbf{Z}_k \boldsymbol{\rho}(\kappa) \quad (3.15)$$

definiert. Dabei beschreibt $\kappa \in \mathbb{N}_{\geq 0}$ den Zeitschritt auf der Solltrajektorie aus lokaler Perspektive des Zeitschritts k . Für $\kappa = 0$ ergibt sich somit der Sollzustand im Zeitschritt k , wohingegen für $\kappa > 0$ eine Vorausschau des Sollzustands für zukünftige Zeiten resultiert.

Das Ziel der optimalen Trajektorienfolgeregelung ist es, dass der Systemzustand $\mathbf{x}_{k+\kappa}$ dem Solltrajektorienverlauf $\mathbf{x}_r(\mathbf{Z}_k, \kappa)$, $\kappa = 0, 1, \dots$, optimal bezüglich eines Gütefunktionals J_k folgt. Konkret soll für ein System mit unbekannter Systemdynamik ein Regler $\boldsymbol{\mu}^*(\mathbf{x}_k, \mathbf{Z}_k)$ gefunden werden, der das Gütefunktional

$$J_k = \sum_{\kappa=0}^{\infty} \gamma^\kappa r(\mathbf{x}_{k+\kappa}, \mathbf{x}_r(\mathbf{Z}_k, \kappa), \mathbf{u}_{k+\kappa}) \quad (3.16)$$

⁴⁵ Grundsätzlich lässt sich der PRADP-Mechanismus auch auf allgemeine nichtlineare Systeme $\mathbf{F}(\mathbf{x}_k, \mathbf{u}_k)$ anwenden, jedoch erfolgt dann für die Optimierung im Policy-Improvement-Schritt im Allgemeinen kein direkter analytischer Ausdruck mehr.

minimiert. Dabei beschreibt $\gamma \in (0, 1]$ einen gegebenen Diskontierungsfaktor und $r(\cdot)$ nicht-negative Einschrittkosten, die beispielsweise Abweichungen des Systemzustands $\mathbf{x}_{k+\kappa}$ von $\mathbf{x}_r(\mathbf{Z}_k, \kappa)$ sowie die aufgebrauchte Stellenergie bestrafen können. Somit ergibt sich die nachfolgende Problemstellung.

Problem 3.1

Für eine durch \mathbf{Z}_k parametrisierte Sollzustandstrajektorie nach (3.15) werde die optimale Stellgrößensequenz, die das Kostenfunktional J_k gemäß (3.16) minimiert, mit $\mathbf{u}_k^*, \mathbf{u}_{k+1}^*, \dots$ und die damit verknüpften Kosten mit J_k^* bezeichnet. Zudem sei die Systemdynamik unbekannt. Gesucht ist die optimale Stellgröße $\mathbf{u}_k^* = \boldsymbol{\mu}^*(\mathbf{x}_k, \mathbf{Z}_k)$ in jedem Zeitschritt k in Abhängigkeit vom aktuellen Zustand \mathbf{x}_k und der aktuellen Solltrajektorienbeschreibung \mathbf{Z}_k .

3.2.2 Q-Function mit parametrierter Referenzdarstellung

Im Folgenden wird eine modifizierte Q-Function vorgestellt, deren minimierendes Regelgesetz eine Lösung $\boldsymbol{\mu}^*(\mathbf{x}_k, \mathbf{Z}_k)$ für Problem 3.1 darstellt. Diese Q-Function kann anschließend durch lineare Funktionsapproximatoren beschrieben und mithilfe einer ADP-Methode – beispielsweise dem LSPI-Ansatz [LP03] – ohne Kenntnis des Systemmodells aus Daten gelernt werden.

Um die durch J_k beschriebenen Kosten in (3.16) zu minimieren, muss die relative Zeit κ auf der aktuellen Solltrajektorie, die nach (3.15) durch \mathbf{Z}_k parametrisiert ist, berücksichtigt werden. Dies geschieht durch die Verwendung einer um κ Zeitschritte verschobenen Parametermatrix $\mathbf{Z}_k^{(\kappa)}$ nach der folgenden Definition.

Definition 3.2 (Verschobene Parametermatrix $\mathbf{Z}_k^{(\kappa)}$)

Die Matrix $\mathbf{Z}_k^{(\kappa)}$ sei derart definiert, dass

$$\mathbf{x}_r(\mathbf{Z}_k^{(\kappa)}, j) = \mathbf{x}_r(\mathbf{Z}_k, \kappa + j) \quad (3.17a)$$

$$\Leftrightarrow \mathbf{Z}_k^{(\kappa)} \boldsymbol{\rho}(j) = \mathbf{Z}_k \boldsymbol{\rho}(\kappa + j) \quad (3.17b)$$

gelte. Demnach ist

$$\mathbf{Z}_k^{(\kappa)} = \mathbf{Z}_k \mathbf{D}(\kappa) = \begin{bmatrix} \mathbf{z}_{k,1}^{(\kappa)\top} \\ \vdots \\ \mathbf{z}_{k,n}^{(\kappa)\top} \end{bmatrix} \quad (3.18)$$

eine modifizierte Version von $\mathbf{Z}_k = \mathbf{Z}_k^{(0)}$, sodass die zugehörige Referenztrajektorie um κ Zeitschritte verschoben ist, wobei $\mathbf{D}(\kappa)$ eine geeignete Matrix ist, sodass (3.17) gilt.

Anzumerken sei an dieser Stelle, dass $\mathbf{D}(\kappa)$ gemäß Definition 3.2 im Allgemeinen mehrdeutig ist, da für den Fall $n_z > 1$ das durch (3.17b) gegebene Gleichungssystem, das verwendet wird, um $\mathbf{Z}_k^{(\kappa)}$ zu ermitteln, unterbestimmt ist.

Weiterhin gelte die Kurzschreibweise $\mathbf{x}_r(\mathbf{Z}_k^{(\kappa)}) := \mathbf{x}_r(\mathbf{Z}_k^{(\kappa)}, 0) = \mathbf{x}_r(\mathbf{Z}_k, \kappa)$. Mit $\mathbf{Z}_k^{(\kappa)}$ nach Definition 3.2 sei die erweiterte Q-Function, die den Solltrajektorienverlauf, parametrisiert durch \mathbf{Z}_k , explizit einbezieht, wie folgt definiert.

Definition 3.3 (Q-Function mit parametrisiertem Solltrajektorienverlauf)

Sei

$$\begin{aligned} Q^\mu(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k) &:= r(\mathbf{x}_k, \mathbf{x}_r(\mathbf{Z}_k), \mathbf{u}_k) \\ &+ \sum_{\kappa=1}^{\infty} \gamma^\kappa r(\mathbf{x}_{k+\kappa}, \mathbf{x}_r(\mathbf{Z}_k^{(\kappa)}), \boldsymbol{\mu}(\mathbf{x}_{k+\kappa}, \mathbf{Z}_k^{(\kappa)})) \\ &= r(\mathbf{x}_k, \mathbf{x}_r(\mathbf{Z}_k), \mathbf{u}_k) \\ &+ \gamma Q^\mu(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)}, \boldsymbol{\mu}(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)})), \end{aligned} \quad (3.19)$$

wobei $\boldsymbol{\mu} : \mathbb{R}^n \times \mathbb{R}^{n \times n_z} \rightarrow \mathbb{R}^p$ das zu bewertende Regelgesetz bezeichne.

Mit dieser Definition repräsentiert $Q^\mu(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k)$ die akkumulierten (diskontierten) Kosten, wenn sich das System im Zustand \mathbf{x}_k befindet, die Stellgröße \mathbf{u}_k im Zeitschritt k angewandt wird, danach dem Regelgesetz $\boldsymbol{\mu}(\cdot)$ gefolgt wird, und die Sollzustandstrajektorie durch \mathbf{Z}_k parametrisiert ist. Basierend auf (3.19) ergibt sich die optimale Q-Function zu

$$\begin{aligned} Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k) &= r(\mathbf{x}_k, \mathbf{x}_r(\mathbf{Z}_k), \mathbf{u}_k) \\ &+ \min_{\boldsymbol{\mu}} \gamma Q^\mu(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)}, \boldsymbol{\mu}(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)})) \\ &= r(\mathbf{x}_k, \mathbf{x}_r(\mathbf{Z}_k), \mathbf{u}_k) \\ &+ \gamma Q^*(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)}, \boldsymbol{\mu}^*(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)})). \end{aligned} \quad (3.20)$$

Dabei bezeichnet $\boldsymbol{\mu}^*(\cdot)$ das optimale Regelgesetz, d. h. es gilt $\boldsymbol{\mu}^*(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)}) = \mathbf{u}_{k+1}^*$. Wie das nachfolgende Lemma, das die Zusammenhänge des klassischen Q-Learnings [WD92] auf den PRADP-Fall überträgt, zeigt, erweist sich diese Q-Function als hilfreich, um Problem 3.1 zu lösen.

Lemma 3.1

Die Stellgröße \mathbf{u}_k , die $Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k)$ minimiert, stellt eine Lösung für \mathbf{u}_k^* nach Problem 3.1 dar.

Beweis:

Mit (3.20) folgt

$$\begin{aligned}
\min_{\mathbf{u}_k} Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k) &= r(\mathbf{x}_k, \mathbf{x}_r(\mathbf{Z}_k), \mathbf{u}_k^*) + \gamma Q^*(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)}, \mathbf{u}_{k+1}^*) \\
&= \min_{\mathbf{u}_k, \mathbf{u}_{k+1}, \dots} \sum_{\kappa=0}^{\infty} \gamma^\kappa r(\mathbf{x}_\kappa, \mathbf{x}_r(\mathbf{Z}_k^{(\kappa)}), \mathbf{u}_\kappa) \\
&= J_k^*
\end{aligned} \tag{3.21}$$

und somit direkt die Aussage des Lemmas. \square

Wenn die optimale Q-Funktion $Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k)$ bekannt ist, ergibt sich nach Lemma 3.1 die gesuchte optimale Stellgröße \mathbf{u}_k^* durch

$$\mathbf{u}_k^* = \arg \min_{\mathbf{u}_k} Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k). \tag{3.22}$$

Im nächsten Abschnitt wird die Q-Funktion durch lineare Funktionsapproximatoren beschrieben und basierend auf dem TD-Fehler die unbekannte Q-Funktion geschätzt.

3.2.3 Funktionsapproximation und Policy Iteration der erweiterten Q-Funktion

Da klassisches, tabellarisches Q-Learning wertkontinuierliche Zustands- und Stellgrößenräume nicht angemessen handhaben kann [van12], wird die als stetig angenommene Q-Funktion durch einen linearen Funktionsapproximator

$$Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k) = \mathbf{w}^{*\top} \boldsymbol{\phi}(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k) + \epsilon(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k) \tag{3.23}$$

beschrieben. Hierbei bezeichnet $\mathbf{w}^* \in \mathbb{R}^h$ den unbekannt optimalen Gewichtsvektor, $\boldsymbol{\phi} \in \mathbb{R}^h$ einen geeigneten Basisfunktionsvektor und ϵ den Approximationsfehler (vgl. Abschnitt 2.1.3). Da \mathbf{w}^* a priori unbekannt ist, sei die geschätzte optimale Q-Funktion durch

$$\hat{Q}^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k) = \hat{\mathbf{w}}^\top \boldsymbol{\phi}(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k) \tag{3.24}$$

gegeben. Analog zu (3.22) folgt daraus das geschätzte optimale Regelgesetz

$$\hat{\boldsymbol{\mu}}^*(\mathbf{x}_k, \mathbf{Z}_k) = \arg \min_{\mathbf{u}_k} \hat{Q}^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k). \tag{3.25}$$

Basierend auf der parametrisierten Q-Funktion ergibt sich der zugehörige TD-Fehler [Sut88] wie nachfolgend definiert.

Definition 3.4 (TD-Fehler der erweiterten Q-Function)

Der TD-Fehler, der sich mit der geschätzten Q-Function $\hat{Q}^*(\cdot)$ (3.24) aus der Bellman-Gleichung (3.20) ergibt, sei durch

$$\begin{aligned} \delta_k &:= r(\mathbf{x}_k, \mathbf{x}_r(\mathbf{Z}_k), \mathbf{u}_k) + \gamma \hat{Q}^* \left(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)}, \hat{\boldsymbol{\mu}}^* \left(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)} \right) \right) \\ &\quad - \hat{Q}^* \left(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k \right) \\ &= r(\mathbf{x}_k, \mathbf{x}_r(\mathbf{Z}_k), \mathbf{u}_k) + \gamma \hat{\mathbf{w}}^\top \boldsymbol{\phi} \left(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)}, \hat{\boldsymbol{\mu}}^* \left(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)} \right) \right) \\ &\quad - \hat{\mathbf{w}}^\top \boldsymbol{\phi} \left(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k \right) \end{aligned} \quad (3.26)$$

definiert.

Da der TD-Fehler die Qualität der Approximation der Q-Function quantifiziert, wird das Gewicht $\hat{\mathbf{w}} \in \mathbb{R}^h$ gesucht, welches den quadratischen TD-Fehler δ_k^2 minimiert. Da (3.26) eine skalare Gleichung darstellt, werden hierzu $M \geq h$ Tupel

$$\mathcal{T}_k := \left\{ r_k, \hat{Q}_k^*, \hat{Q}_k^{*+} \right\}, \quad k = 1, \dots, M, \quad (3.27a)$$

mit

$$r_k := r(\mathbf{x}_k, \mathbf{x}_r(\mathbf{Z}_k), \mathbf{u}_k), \quad (3.27b)$$

$$\hat{Q}_k^* := \hat{\mathbf{w}}^\top \boldsymbol{\phi}_k := \hat{\mathbf{w}}^\top \boldsymbol{\phi}(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k), \quad (3.27c)$$

$$\hat{Q}_k^{*+} := \hat{\mathbf{w}}^\top \boldsymbol{\phi}_k^+ := \hat{\mathbf{w}}^\top \boldsymbol{\phi} \left(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)}, \hat{\boldsymbol{\mu}}^* \left(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)} \right) \right) \quad (3.27d)$$

verwendet. Diese M Tupel werden aus gemessenen oder simulierten Systemtrajektorien $(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1})$, den Solltrajektorienparametern \mathbf{Z}_k und verschobenen Solltrajektorienparametern $\mathbf{Z}_k^{(1)}$ gebildet. Anschließend kann $\hat{\mathbf{w}}$ mit einer geeigneten ADP-Methode, beispielsweise mit dem LSPI-Algorithmus [LP03], geschätzt werden. Mit den Tupeln $\mathcal{T}_k, k = 1, \dots, M$, folgt aus (3.26)⁴⁶

$$\underbrace{\begin{bmatrix} \delta_1 \\ \vdots \\ \delta_M \end{bmatrix}}_{=: \boldsymbol{\delta}} = \underbrace{\begin{bmatrix} r_1 \\ \vdots \\ r_M \end{bmatrix}}_{=: \mathbf{r}} + \underbrace{\left(\gamma \begin{bmatrix} \boldsymbol{\phi}_1^{+\top} \\ \vdots \\ \boldsymbol{\phi}_M^{+\top} \end{bmatrix} - \begin{bmatrix} \boldsymbol{\phi}_1^\top \\ \vdots \\ \boldsymbol{\phi}_M^\top \end{bmatrix} \right)}_{=: \boldsymbol{\Phi}} \hat{\mathbf{w}}. \quad (3.28)$$

Sofern die Anregungsbedingung

$$\text{Rang}(\boldsymbol{\Phi}^\top \boldsymbol{\Phi}) = h \quad (3.29)$$

erfüllt ist, existiert nach Åström und Wittenmark [ÅW95, Theorem 2.1] durch

$$\hat{\mathbf{w}} = (\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{r} \quad (3.30)$$

eine eindeutige Lösung $\hat{\mathbf{w}}$, welche $\boldsymbol{\delta}^\top \boldsymbol{\delta}$ minimiert.

⁴⁶ Dies stellt im Wesentlichen eine Erweiterung von [LP03, Abschnitt 5.1] auf ADP-Solltrajektorienfolgeregler mit parametrimtem Referenzverlauf dar.

Bemerkung 3.2

Die Verwendung von $\mathbf{Z}_k^{(1)} = \mathbf{Z}_k \mathbf{D}(1)$ (vgl. (3.18)) anstelle von beliebigen nachfolgenden Parametern \mathbf{Z}_{k+1} in den Trainingstupeln \mathcal{T}_k (3.27a) ist essenziell, da hierdurch (in Kombination mit der Systemdynamik (3.13)) gewährleistet wird, dass die für ADP-Methoden wichtige Markov-Eigenschaft erfüllt ist (vgl. Abschnitt 3.1 und [LV09]). Insbesondere wird hierdurch sichergestellt, dass die für den ADP-Formalismus verwendete Solltrajektorien-darstellung ADP-kompatibel nach Definition 3.1 ist, denn mit⁴⁷

$$\zeta_k := \text{vec}(\mathbf{Z}_k) \quad (3.31)$$

folgt in Analogie zu (3.6)

$$\mathbf{f}_{x_r, \zeta}(\zeta_k) := \mathbf{Z}_k \boldsymbol{\rho}(0) = \text{mat}(\zeta_k, n, n_z) \boldsymbol{\rho}(0) \quad (3.32)$$

(vgl. (3.15)) und analog zu (3.7)

$$\mathbf{f}_\zeta(\zeta_k) := \text{vec}(\mathbf{Z}_k^{(1)}) = \text{vec}(\mathbf{Z}_k \mathbf{D}(1)) = \text{vec}(\text{mat}(\zeta_k, n, n_z) \mathbf{D}(1)). \quad (3.33)$$

An dieser Stelle sei anzumerken, dass $\hat{\boldsymbol{\mu}}^*(\cdot)$, welches in \hat{Q}_k^{*+} nach (3.27d) verwendet wird, selbst eine Schätzung (für das optimale Regelgesetz) darstellt. Dieser Mechanismus, dass eine Schätzung basierend auf einer anderen Schätzung erfolgt, ist im Reinforcement Learning unter dem Begriff *bootstrapping* bekannt (vgl. [SB18]). Daher genügt eine einmalige Schätzung von $\hat{\boldsymbol{w}}$ basierend auf der Least-Squares-Lösung von (3.30) nicht, um die optimale Q-Funktion und den optimalen Regler zu ermitteln. Stattdessen wird eine Policy Iteration, beginnend bei einem zulässigen Initialgewicht $\hat{\boldsymbol{w}}^{[0]}$, durchgeführt. Dieses Vorgehen ist in Algorithmus 3.1 zusammengefasst, wobei der Schwellwert $e_{\hat{\boldsymbol{w}}}$ eine Abbruchbedingung darstellt.

Algorithmus 3.1 PRADP mithilfe des LSPI-Algorithmus

- 1: **Initialisiere** $l := 0, \hat{\boldsymbol{w}}^{[0]}$, sodass $\hat{\boldsymbol{\mu}}^{[0]}(\cdot)$ zulässig ist
- 2: **do**
- 3: Policy Evaluation: berechne $\hat{\boldsymbol{w}}^{[l+1]}$ basierend auf (3.30) mit $\hat{\boldsymbol{w}} = \hat{\boldsymbol{w}}^{[l+1]}$
- 4: Policy Improvement: berechne $\hat{\boldsymbol{\mu}}^{[l+1]}$ nach (3.25)
- 5: $l := l + 1$
- 6: **while** $\left\| \hat{\boldsymbol{w}}^{[l]} - \hat{\boldsymbol{w}}^{[l-1]} \right\|_2 > e_{\hat{\boldsymbol{w}}}$

⁴⁷ Hierbei bildet $\text{mat}(\cdot)$ aus einem Vektor eine Matrix, wobei $\text{mat}(\text{vec}(\mathbf{M}), \tilde{n}, p) = \mathbf{M} \in \mathbb{R}^{\tilde{n} \times p}$ gilt.

Bemerkung 3.3

Aufgrund der Verwendung einer Q -Function, die explizit von der grundsätzlich beliebigen Stellgröße \mathbf{u}_k abhängt, liegt ein Off-Policy-Verfahren vor (vgl. Abschnitt 2.1.4.4). Dies erweist sich insbesondere als vorteilhaft, da während des Aufzeichnens der M Datentupel \mathcal{T}_k , $k = 1, \dots, M$, üblicherweise eine Anregung des Systems stattfinden muss, um die durch (3.29) gegebene Anregungsbedingung zu erfüllen und einen erfolgreichen Policy-Evaluation-Schritt zu gewährleisten. Die Behavior Policy, d. h. die Stellgröße \mathbf{u}_k , die tatsächlich auf das System angewandt wird, kann daher zur Systemanregung verwendet werden (beispielsweise durch additive Überlagerung von Rauschen oder harmonischen Schwingungen, siehe Kapitel 5). Demgegenüber wird die Target Policy $\hat{\boldsymbol{\mu}}^*$, die dem geschätzten optimalen Regelgesetz nach (3.25) entspricht, im Ausdruck $\gamma \hat{Q}^*(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)}, \hat{\boldsymbol{\mu}}^*(\mathbf{x}_{k+1}, \mathbf{Z}_k^{(1)}))$ in (3.26) verwendet, weshalb die mit dem geschätzten optimalen Regelgesetz assoziierte Q -Function gelernt wird.

Mit $\hat{Q}^{[l]}(\cdot) = \hat{\mathbf{w}}^{[l]\top} \boldsymbol{\phi}(\cdot)$ und $Q^{\hat{\boldsymbol{\mu}}^{[l]}}$ nach (3.19), wobei $\boldsymbol{\mu} = \hat{\boldsymbol{\mu}}^{[l]}$ sei, lässt sich die in [LP03, Theorem 7.1] gegebene Konvergenzaussage des LSPI-Algorithmus direkt auf den auf parametrisierte Solltrajektorienverläufe erweiterten Fall übertragen.

Proposition 3.2 (Konvergenz der Q -Function nach [LP03, Theorem 7.1])

Sei für alle Iterationen l durch $\bar{\epsilon} \geq 0$ eine obere Schranke für den Approximationsfehler zwischen der zum Regelgesetz $\hat{\boldsymbol{\mu}}^{[l]}$ gehörenden geschätzten Q -Function $\hat{Q}^{[l]}$ und der wahren Q -Function $Q^{\hat{\boldsymbol{\mu}}^{[l]}}$ gegeben, d. h.

$$\left\| \hat{Q}^{[l]} - Q^{\hat{\boldsymbol{\mu}}^{[l]}} \right\|_{\infty} \leq \bar{\epsilon}, \quad \forall l = 1, 2, \dots \quad (3.34)$$

Dann gilt für Algorithmus 3.1:

$$\limsup_{l \rightarrow \infty} \left\| \hat{Q}^{[l]} - Q^* \right\|_{\infty} \leq \frac{2\gamma\bar{\epsilon}}{(1-\gamma)^2}. \quad (3.35)$$

Beweis:

Der Beweis ist in [LP03, Theorem 7.1] bzw. [BT96, Proposition 6.2] zu finden. \square

Lagoudakis und Parr [LP03] betonen, dass eine geeignete Wahl der Basisfunktionen $\boldsymbol{\phi}(\cdot)$ sowie die Erzeugung der Datentupel (d. h. eine geeignete Anregung) maßgeblich für die Fehlerschranke $\bar{\epsilon}$ verantwortlich sind. Nach Proposition 3.2 konvergiert Algorithmus 3.1 für adäquate Funktionsapproximatoren $\boldsymbol{\phi}(\cdot)$ sowie eine angemessene Anregung somit gegen eine Nachbarschaft der optimalen erweiterten Q -Function. Jedoch ist für eingangsaﬃne Systeme (3.13) bzw. generelle nichtlineare Systeme $\mathbf{F}(\mathbf{x}_k, \mathbf{u}_k)$ und allgemeine Kostenfunktionale (3.16) eine geeignete

Wahl an Basisfunktionen sowie deren Anzahl nach wie vor ein ungelöstes Problem⁴⁸ [WHL17]. Zudem erfordert eine geeignete Systemanregung im allgemeinen, nichtlinearen Fall, dass die Trainingsdaten alle relevanten Bereiche des Zustands- und Stellgrößenraums abdecken, um die vorhandenen Nichtlinearitäten angemessen zu berücksichtigen und optimale Gewichte \mathbf{w}^* lernen zu können. Die neue ADP-basierte Solltrajektorienfolgeregelungsmethode wird im Folgenden anhand linearer Systeme und quadratischer Gütefunktionale betrachtet. Dieser in zahlreichen regelungstechnischen Problemen relevante LQ-Fall ermöglicht analytische Einsichten in die Struktur der Q-Function $Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k)$ und somit eine geeignete Wahl der Basisfunktionen $\phi(\cdot)$ zur Funktionsapproximation.

3.2.4 Linear-quadratische optimale Trajektorienfolgeregelung mit parametrierter Referenz

Im Folgenden sei die Systemdynamik durch

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k \quad (3.36)$$

und das Kostenfunktional durch

$$\begin{aligned} J_k &= \sum_{\kappa=0}^{\infty} \gamma^{\kappa} \left[\left(\mathbf{x}_{k+\kappa} - \mathbf{x}_r \left(\mathbf{Z}_k^{(\kappa)} \right) \right)^{\top} \mathbf{Q} \left(\mathbf{x}_{k+\kappa} - \mathbf{x}_r \left(\mathbf{Z}_k^{(\kappa)} \right) \right) + \mathbf{u}_{k+\kappa}^{\top} \mathbf{R} \mathbf{u}_{k+\kappa} \right] \\ &=: \sum_{\kappa=0}^{\infty} \gamma^{\kappa} \left[\mathbf{e}_{k,\kappa}^{\top} \mathbf{Q} \mathbf{e}_{k,\kappa} + \mathbf{u}_{k+\kappa}^{\top} \mathbf{R} \mathbf{u}_{k+\kappa} \right] \end{aligned} \quad (3.37)$$

mit $\mathbf{e}_{k,\kappa} := \mathbf{x}_{k+\kappa} - \mathbf{x}_r \left(\mathbf{Z}_k^{(\kappa)} \right)$ beschrieben. Hierbei bestraft $\mathbf{Q} \in \mathbb{R}^{n \times n}$ die Abweichung des Systemzustands $\mathbf{x}_{k+\kappa}$ vom Sollzustand $\mathbf{x}_r \left(\mathbf{Z}_k^{(\kappa)} \right)$ und $\mathbf{R} \in \mathbb{R}^{p \times p}$ gewichtet den Stellaufwand. Zudem seien die folgenden Annahmen erfüllt.

Annahme 3.1

Sei $\mathbf{Q} = \mathbf{Q}^{\top} \succeq \mathbf{0}$, $\mathbf{R} = \mathbf{R}^{\top} \succ \mathbf{0}$, (\mathbf{A}, \mathbf{B}) steuerbar und (\mathbf{A}, \mathbf{C}) detektierbar, wobei \mathbf{C} derart definiert ist, dass $\mathbf{C}^{\top} \mathbf{C} = \mathbf{Q}$ gilt⁴⁹.

Annahme 3.2

Sei die Matrix $\mathbf{D}(\kappa)$, welche die verschobene Parametermatrix $\mathbf{Z}_k^{(\kappa)}$ nach (3.18) definiert (vgl. Definition 3.2), so, dass $|\lambda_j| < 1, \forall j = 1, \dots, n_z$, gilt, wobei λ_j die Eigenwerte von $\sqrt{\gamma} \mathbf{D}(1)$ sind.

⁴⁸ Wang et al. [WHL17] bezeichnen die geeignete Wahl des Funktionsapproximators als „more of an art than science“.

⁴⁹ Jede symmetrische, positiv semidefinite Matrix lässt sich als $\mathbf{Q} = \mathbf{C}^{\top} \mathbf{C}$ schreiben, beispielsweise durch Diagonalisierung (vgl. [Lib12, S. 195]).

Die Bedeutung dieser Annahmen wird in der folgenden Bemerkung eingeordnet.

Bemerkung 3.4

Annahme 3.1 ist üblich für linear-quadratische Optimalregelungsprobleme. Hierdurch wird die Existenz und Eindeutigkeit einer stabilisierenden Lösung der zeitdiskreten algebraischen Riccati-Gleichung des Regulationsproblems, das durch (3.36) und (3.37) für $\mathbf{Z}_k = \mathbf{0}$ (d. h. für einen Sollzustand $\mathbf{0}$ für alle k) gegeben ist, gewährleistet (vgl. [Kuč72, Theorem 8]).

Zudem ist es einleuchtend, dass für eine sinnvolle Problemstellung der Solltrajektorienverlauf $\mathbf{x}_r(\mathbf{Z}_k^{(\kappa)})$ derart definiert sein muss, dass ein Regelgesetz existiert, das mit einem endlichen Gütefunktional J_k verknüpft ist. Wie der nachfolgende Satz 3.1 zeigen wird, garantieren Annahme 3.1 und Annahme 3.2 die Existenz einer solchen Lösung.

Bevor das optimale Regelgesetz hergeleitet wird, wird der Trajektorienfehler $e_{k,\kappa}$ aus (3.37) durch

$$\begin{aligned}
 e_{k,\kappa} &= \mathbf{x}_{k+\kappa} - \mathbf{x}_r(\mathbf{Z}_k^{(\kappa)}) \\
 &= \mathbf{x}_{k+\kappa} - \mathbf{Z}_k^{(\kappa)} \boldsymbol{\rho}(0) \\
 &= \underbrace{\begin{bmatrix} \mathbf{I}_n & \begin{bmatrix} -\boldsymbol{\rho}(0) & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & -\boldsymbol{\rho}(0) \end{bmatrix} \end{bmatrix}^\top}_{=: M} \underbrace{\begin{bmatrix} \mathbf{x}_{k+\kappa} \\ \mathbf{z}_{k,1}^{(\kappa)} \\ \vdots \\ \mathbf{z}_{k,n}^{(\kappa)} \end{bmatrix}}_{=: \tilde{\mathbf{x}}_{k,\kappa}}, \tag{3.38}
 \end{aligned}$$

$\kappa = 0, 1, \dots$, ausgedrückt. \mathbf{I}_n bezeichnet die $n \times n$ -Einheitsmatrix und $\tilde{\mathbf{x}}_{k,\kappa}$ stellt den um den Referenzparameter $\mathbf{Z}_k^{(\kappa)}$ erweiterten Systemzustand $\mathbf{x}_{k+\kappa}$ dar. Der zugehörige Optimalregler ist durch den folgenden Satz gegeben.

Satz 3.1 (Optimales Regelgesetz des zeitdiskreten parametrischen Trajektorienfolgeregelungsproblems)

Seien eine Sollzustandsrepräsentation $\mathbf{x}_\tau(\mathbf{Z}_k^{(\kappa)})$ (vgl. (3.15)) und eine Verschiebungsmatrix $\mathbf{D}(\kappa)$ wie in Definition 3.2 gegeben.

1. Das optimale Regelgesetz, das (3.37) unter Berücksichtigung der Systemdynamik (3.36) minimiert, ist linear bezüglich $\tilde{\mathbf{x}}_{k,\kappa}$ in (3.38) und somit durch

$$\boldsymbol{\mu}^*(\mathbf{x}_{k+\kappa}, \mathbf{Z}_k^{(\kappa)}) = \mathbf{u}_{k+\kappa}^* = -\mathbf{K}^* \tilde{\mathbf{x}}_{k,\kappa}, \quad \kappa = 0, 1, \dots, \quad (3.39)$$

gegeben. Die optimale Verstärkungsmatrix \mathbf{K}^* ergibt sich zu

$$\mathbf{K}^* = (\gamma \tilde{\mathbf{B}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{B}} + \mathbf{R})^{-1} \gamma \tilde{\mathbf{B}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{A}}, \quad (3.40)$$

mit

$$\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{A} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{D}^\top(1) & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{D}^\top(1) \end{bmatrix}, \quad \tilde{\mathbf{B}} = \begin{bmatrix} \mathbf{B} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix}, \quad (3.41)$$

$\tilde{\mathbf{A}} \in \mathbb{R}^{n(n_z+1) \times n(n_z+1)}$, $\tilde{\mathbf{B}} \in \mathbb{R}^{n(n_z+1) \times p}$. In (3.40) bezeichnet $\tilde{\mathbf{P}}$ die Lösung der zeitdiskreten algebraischen Riccati-Gleichung

$$\tilde{\mathbf{P}} = \gamma \tilde{\mathbf{A}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{A}} - \gamma \tilde{\mathbf{A}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{B}} (\mathbf{R} + \tilde{\mathbf{B}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{B}})^{-1} \tilde{\mathbf{B}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{A}} + \tilde{\mathbf{Q}} \quad (3.42)$$

mit $\tilde{\mathbf{Q}} = \mathbf{M}^\top \mathbf{Q} \mathbf{M}$ und \mathbf{M} nach (3.38).

2. Des Weiteren sind unter Annahme 3.1 und Annahme 3.2 Existenz und Eindeutigkeit des optimalen Regelgesetzes $\boldsymbol{\mu}^*(\mathbf{x}_{k+\kappa}, \mathbf{Z}_k^{(\kappa)})$ gesichert.

Beweis:

1. Mit (3.38) kann das Gütefunktional (3.37) zu

$$J_k = \sum_{\kappa=0}^{\infty} \gamma^\kappa \left[\tilde{\mathbf{x}}_{k,\kappa}^\top \mathbf{M}^\top \mathbf{Q} \mathbf{M} \tilde{\mathbf{x}}_{k,\kappa} + \mathbf{u}_{k+\kappa}^\top \mathbf{R} \mathbf{u}_{k+\kappa} \right] \quad (3.43)$$

umgeformt werden. Mit (3.36) und (3.18) gilt zudem

$$\tilde{\mathbf{x}}_{k,\kappa+1} = \begin{bmatrix} \mathbf{A}\mathbf{x}_{k+\kappa} + \mathbf{B}\mathbf{u}_{k+\kappa} \\ D(1)^\top \mathbf{z}_{k,1}^{(\kappa)} \\ \vdots \\ D(1)^\top \mathbf{z}_{k,n}^{(\kappa)} \end{bmatrix} = \tilde{\mathbf{A}}\tilde{\mathbf{x}}_{k,\kappa} + \tilde{\mathbf{B}}\mathbf{u}_{k+\kappa}. \quad (3.44)$$

Mit γ , $\tilde{\mathbf{A}}$, $\tilde{\mathbf{B}}$, $\tilde{\mathbf{Q}}$ und \mathbf{R} ergibt sich somit ein *diskontiertes* Standard-LQ-Optimalregelungsproblem aus (3.43) und (3.44) für den erweiterten Zustand $\tilde{\mathbf{x}}_{k,\kappa}$. Dieses diskontierte Problem ist äquivalent zum undiskontierten LQ-Optimalregelungsproblem mit $\sqrt{\gamma}\tilde{\mathbf{A}}$, $\sqrt{\gamma}\tilde{\mathbf{B}}$, $\tilde{\mathbf{Q}}$ und \mathbf{R} (vgl. [GGH⁺18]). Für dieses undiskontierte Standard-LQ-Optimalregelungsproblem ist wohlbekannt, dass der optimale Regler linear bezüglich des Zustands (d. h. hier des erweiterten Zustands $\tilde{\mathbf{x}}_{k,\kappa}$) ist und die optimale Zustandsrückführungsmatrix durch (3.40) gegeben ist [LVS12, Abschnitt 2.4]. Daher gilt (3.39) und somit die erste Aussage von Satz 3.1.

2. Bezüglich der zweiten Aussage sei zunächst angemerkt, dass die Stabilisierbarkeit von $(\sqrt{\gamma}\tilde{\mathbf{A}}, \sqrt{\gamma}\tilde{\mathbf{B}})$ direkt aus Annahme 3.1 und Annahme 3.2 folgt, da (\mathbf{A}, \mathbf{B}) steuerbar ist und $|\lambda_j| < 1$, $\forall j = 1, \dots, n_z$, gilt. Zudem folgt aus $\mathbf{Q} \succeq \mathbf{0}$ direkt $\tilde{\mathbf{Q}} \succeq \mathbf{0}$. Weil (\mathbf{A}, \mathbf{C}) nach Annahme 3.1 detektierbar ist und da alle zusätzlichen Zustände in $\tilde{\mathbf{A}}$ im Vergleich zu \mathbf{A} aufgrund von Annahme 3.2 stabil sind, folgt, dass $(\sqrt{\gamma}\tilde{\mathbf{A}}, \tilde{\mathbf{C}})$ (mit $\tilde{\mathbf{C}}$ so, dass $\tilde{\mathbf{C}}^\top \tilde{\mathbf{C}} = \tilde{\mathbf{Q}}$) ebenfalls detektierbar ist. Schließlich folgt nach Kučera [Kuč72, Theorem 8], dass wegen $\tilde{\mathbf{Q}} \succeq \mathbf{0}$, $\mathbf{R} \succ \mathbf{0}$, $(\sqrt{\gamma}\tilde{\mathbf{A}}, \sqrt{\gamma}\tilde{\mathbf{B}})$ stabilisierbar und $(\sqrt{\gamma}\tilde{\mathbf{A}}, \tilde{\mathbf{C}})$ detektierbar eine eindeutige Lösung existiert. \square

Bemerkung 3.5

Satz 3.1 offenbart, dass im Fall bekannter Systemmatrizen \mathbf{A} und \mathbf{B} die optimale Lösung \mathbf{K}^* direkt durch das Lösen der zeitdiskreten algebraischen Riccati-Gleichung (vgl. [AL84]), die zu $\sqrt{\gamma}\tilde{\mathbf{A}}$, $\sqrt{\gamma}\tilde{\mathbf{B}}$, $\tilde{\mathbf{Q}}$ und \mathbf{R} gehört, resultiert.

Gleichung (3.44) zeigt letztlich auch, dass die, aus ADP-Sicht fundamentale, Markov-Eigenschaft gilt. Ausgehend von Satz 3.1 wird für den Fall unbekannter Systemmatrizen \mathbf{A} und \mathbf{B} die folgende LQ-PRADP-Problemstellung formuliert.

Problem 3.2 (LQ-PRADP)

Sei die durch die Matrizen \mathbf{A} und \mathbf{B} beschriebene Systemdynamik unbekannt. Gesucht ist die Reglermatrix \mathbf{K}^* (vgl. (3.39)), die das Kostenfunktional J_k (3.37) minimiert. Die optimale Stellgröße im Zeitschritt k ist dabei durch

$$\mathbf{u}_k^* = \boldsymbol{\mu}^* \left(\mathbf{x}_k, \mathbf{Z}_k^{(0)} \right) = -\mathbf{K}^* \tilde{\mathbf{x}}_{k,0} =: - \begin{bmatrix} \mathbf{K}_x^* & \mathbf{K}_z^* \end{bmatrix} \begin{bmatrix} \mathbf{x}_k \\ z_{k,1}^{(0)} \\ \vdots \\ z_{k,n}^{(0)} \end{bmatrix} \quad (3.45)$$

gegeben.

Bevor das optimale Regelgesetz \mathbf{K}^* ohne Kenntnis der Systemmatrizen \mathbf{A} und \mathbf{B} mit einem ADP-Ansatz gelernt wird, wird im folgenden Lemma zunächst die genaue Struktur der optimalen Q-Function $Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k)$, die zu Problem 3.2 gehört, analysiert.

Lemma 3.2 (Struktur der Q-Function mit parametrierter Solltrajektorien-darstellung)

Die zu Problem 3.2 gehörende Q-Function weist die quadratische Form

$$Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k) = \bar{\mathbf{y}}_k^\top \mathbf{H} \bar{\mathbf{y}}_k = \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \\ z_{k,1:n} \end{bmatrix}^\top \begin{bmatrix} h_{xx} & h_{xu} & h_{xz} \\ h_{ux} & h_{uu} & h_{uz} \\ h_{zx} & h_{zu} & h_{zz} \end{bmatrix} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \\ z_{k,1:n} \end{bmatrix} \quad (3.46)$$

mit $\bar{\mathbf{y}}_k := [\mathbf{x}_k^\top \quad \mathbf{u}_k^\top \quad z_{k,1:n}^\top]^\top := [\mathbf{x}_k^\top \quad \mathbf{u}_k^\top \quad z_{k,1}^\top \quad \dots \quad z_{k,n}^\top]^\top$ auf, wobei \mathbf{H} so gewählt sei, dass $\mathbf{H} = \mathbf{H}^\top$ gilt.

Beweis:

Mit (3.19) und (3.20) folgt

$$Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k) = r(\mathbf{x}_k, \mathbf{x}_r(\mathbf{Z}_k), \mathbf{u}_k) + \sum_{\kappa=1}^{\infty} \gamma^\kappa r \left(\mathbf{x}_{k+\kappa}, \mathbf{x}_r \left(\mathbf{Z}_k^{(\kappa)} \right), \boldsymbol{\mu}^* \left(\mathbf{x}_{k+\kappa}, \mathbf{Z}_k^{(\kappa)} \right) \right). \quad (3.47)$$

Wegen (3.36), (3.39) und (3.18) folgt, dass die Zustände $\mathbf{x}_{k+\kappa}$ und Stellgrößen $\boldsymbol{\mu}^* \left(\mathbf{x}_{k+\kappa}, \mathbf{Z}_k^{(\kappa)} \right)$ $\forall \kappa = 0, 1, \dots$ linear bezüglich $\bar{\mathbf{y}}_k$ sind. Aus dieser linearen Abhängigkeit und mit (3.38) folgt Linearität von $e_{k,\kappa}$ bezüglich $\bar{\mathbf{y}}_k$, $\forall \kappa \geq 0$. Durch die Linearität von $e_{k,\kappa}$ und $\boldsymbol{\mu}^*(\cdot)$ bezüglich $\bar{\mathbf{y}}_k$ und die quadratische Struktur von $r(\cdot)$ in (3.37) ist die Q-Function in (3.47) quadratisch bezüglich $\bar{\mathbf{y}}_k$ und folglich gilt (3.46). \square

Eine Konsequenz von Lemma 3.2 ist, dass die optimale Q-Function Q^* exakt durch einen Funktionsapproximator \hat{Q}^* beschrieben werden kann, wenn $\hat{w} = w^*$ den nicht-redundanten Elementen der Matrix $\mathbf{H} = \mathbf{H}^\top$ entspricht⁵⁰ und die Funktionsapproximatoren zu

$$\phi(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k) = \bar{\mathbf{y}}_k \otimes_r \bar{\mathbf{y}}_k \quad (3.48)$$

gewählt werden⁵¹. Basierend auf Lemma 3.2 kann das optimale Regelgesetz wie nachfolgend gezeigt durch die Matrix \mathbf{H} und ohne explizite Verwendung der Systemmatrizen \mathbf{A} und \mathbf{B} angegeben werden.

Satz 3.2 (Optimales Regelgesetz in Abhängigkeit von \mathbf{H})

Die eindeutige erweiterte optimale Zustandsrückführung, die J_k (3.37) minimiert, ist durch

$$\mathbf{u}_k^* = \boldsymbol{\mu}^*(\mathbf{x}_k, \mathbf{Z}_k) = -\mathbf{K}^* \tilde{\mathbf{x}}_{k,0} = -\underbrace{\mathbf{h}_{\text{uu}}^{-1} \begin{bmatrix} \mathbf{h}_{\text{ux}} & \mathbf{h}_{\text{uz}} \end{bmatrix}}_{\begin{bmatrix} \mathbf{K}_x^* & \mathbf{K}_z^* \end{bmatrix}} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{z}_{k,1:n} \end{bmatrix} \quad (3.49)$$

gegeben.

Beweis:

Nach Lemma 3.1 minimiert die Stellgröße \mathbf{u}_k^* , die $Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k)$ minimiert, auch J_k . Mit (3.46) und $\mathbf{H} = \mathbf{H}^\top$ liefert die notwendige Bedingung

$$\frac{\partial Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k)}{\partial \mathbf{u}_k} = 2(\mathbf{h}_{\text{ux}} \mathbf{x}_k + \mathbf{h}_{\text{uz}} \mathbf{z}_{k,1:n} + \mathbf{h}_{\text{uu}} \mathbf{u}_k) \stackrel{!}{=} \mathbf{0} \quad (3.50)$$

die in (3.49) gegebene Stellgröße \mathbf{u}_k^* . Des Weiteren zeigt

$$\frac{\partial^2 Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k)}{\partial \mathbf{u}_k^2} = 2\mathbf{h}_{\text{uu}}, \quad (3.51)$$

dass $\mathbf{h}_{\text{uu}} \succ \mathbf{0}$ gelten muss, damit die Stellgröße \mathbf{u}_k^* in (3.49) auch tatsächlich J_k in (3.37) minimiert. Daher wird im Folgenden $\mathbf{h}_{\text{uu}} \succ \mathbf{0}$ gezeigt. Sei $Q_{\text{reg}}^*(\mathbf{x}_k, \mathbf{u}_k)$ die optimale Q-Function, die den Regulationsfall mit $\mathbf{x}_r(\mathbf{Z}_k^{(\kappa)}) = \mathbf{x}_r(\mathbf{0}, \kappa) = \mathbf{0}$ beschreibt. Dann gilt

$$Q^*(\mathbf{x}_k, \mathbf{0}, \mathbf{u}_k) = Q_{\text{reg}}^*(\mathbf{x}_k, \mathbf{u}_k), \quad \forall \mathbf{x}_k \in \mathbb{R}^n, \mathbf{u}_k \in \mathbb{R}^p. \quad (3.52)$$

Zudem gilt nach Bradtke et al. [BYB94]

$$\begin{aligned} Q_{\text{reg}}^*(\mathbf{x}_k, \mathbf{u}_k) &= \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix}^\top \begin{bmatrix} \mathbf{h}_{\text{reg,xx}} & \mathbf{h}_{\text{reg,xu}} \\ \mathbf{h}_{\text{reg,ux}} & \mathbf{h}_{\text{reg,uu}} \end{bmatrix} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix}^\top \begin{bmatrix} \gamma \mathbf{A}^\top \mathbf{P} \mathbf{A} + \mathbf{Q} & \gamma \mathbf{A}^\top \mathbf{P} \mathbf{B} \\ \gamma \mathbf{B}^\top \mathbf{P} \mathbf{A} & \gamma \mathbf{B}^\top \mathbf{P} \mathbf{B} + \mathbf{R} \end{bmatrix} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix} \end{aligned} \quad (3.53)$$

⁵⁰ Nicht-redundant meint hierbei alle Elemente h_{ij} der symmetrischen Matrix \mathbf{H} mit $i \leq j$. Elemente von \hat{w} , die zu Nebendiagonalelementen von \mathbf{H} gehören, werden zudem mit dem Faktor 2 multipliziert.

⁵¹ Der Operator \otimes_r berechnet für identische Operanden das reduzierte Kronecker-Produkt mit nicht-redundanten Elementen.

für den Regulationsfall. Hierbei ist P die Lösung der zeitdiskreten algebraischen Riccati-Gleichung

$$P = \gamma A^T P A - \gamma A^T P B (R + B^T P B)^{-1} B^T P A + Q. \quad (3.54)$$

Unter Annahme 3.1 existiert eine eindeutige Lösung $P = P^T \succeq 0$ (vgl. [Kuč72, Theorem 8]). Daher folgt aus (3.52) und (3.53)

$$h_{uu} = h_{\text{reg,uu}} = \gamma B^T P B + R \succ 0 \quad (3.55)$$

und somit die gesuchte Aussage. \square

Da in (3.16) potenziell ein diskontiertes Gütemaß verwendet wird und Satz 3.1 aufgrund der Analogie der diskontierten Problemstellung zum nicht-diskontierten Problem mit der Systemmatrix $\sqrt{\gamma} \tilde{A}$ und der Eingangsmatrix $\sqrt{\gamma} \tilde{B}$ zunächst nur Stabilität des fiktiven Systems

$$\tilde{x}_{k+1} = \sqrt{\gamma} \tilde{A} \tilde{x}_k - \sqrt{\gamma} \tilde{B} \tilde{K}^* \tilde{x}_k \quad (3.56)$$

gewährleistet, muss nachfolgend noch überprüft werden, ob für die gewählte Diskontierung $0 < \gamma \leq 1$ das geregelte System

$$x_{k+1} = A x_k - B K_x^* x_k - B K_\zeta^* \zeta_k \quad (3.57)$$

für beschränkte Referenzparametervorgaben $\|\zeta_k\|_2 < \infty$ mit $\zeta_k = \text{vec}(Z_k)$ stabil ist, d. h. ob sämtliche Eigenwerte von $A - B K_x^*$ innerhalb des Einheitskreises liegen. In Anlehnung an die Eingangs-Zustands-Stabilität [Kha02, Definition 4.7] sei dieser Zusammenhang wie nachfolgend definiert als zeitdiskrete Referenz-Zustands-Stabilität bezeichnet.

Definition 3.5 (Zeitdiskrete Referenz-Zustands-Stabilität)

Ein geregeltes System ist Referenz-Zustands-stabil, wenn für beschränkte Referenzparametervorgaben

$$\|\zeta_k\|_2 < \infty$$

mit $\zeta_k = \text{vec}(Z_k)$ beschränkte Systemzustände

$$\|x_k\|_2 < \infty, \forall k \geq 0,$$

resultieren.

Bei Kenntnis von A und B können die Eigenwerte von $A - B K_x^*$ mithilfe des berechneten oder gelernten Reglers $K^* = [K_x^* \quad K_z^*]$ direkt überprüft werden. Da die durch A und B beschriebene Systemdynamik jedoch insbesondere im Fall adaptiver Optimalregler nicht immer bekannt ist, wird nachfolgend ein hinreichendes Kriterium gegeben, mit dem ausschließlich anhand von Q , R und der ggf. mittels eines ADP-Ansatzes gelernten Q-Funktion $Q^*(x_k, Z_k, u_k)$ überprüft werden kann, ob die Reglermatrix K^* das System (3.57) stabilisiert.

Satz 3.3 (Stabilitätskriterium des Solltrajektorienfolgeerregungsproblems)

Sei H wie in (3.46) definiert. Aus $Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k)$ folgt $P^* := \mathbf{h}_{xx} - \mathbf{h}_{xu} \mathbf{h}_{uu}^{-1} \mathbf{h}_{ux}$ und $\mathbf{K}_x^* = \mathbf{h}_{uu}^{-1} \mathbf{h}_{ux}$ (vgl. (3.49)). Wenn

$$b_\gamma := \frac{1}{\lambda_{\min} \left(\mathbf{P}^* (\mathbf{P}^* - \mathbf{Q} - \mathbf{K}_x^{*\top} \mathbf{R} \mathbf{K}_x^*)^{-1} \right)} < \gamma \leq 1 \quad (3.58)$$

gilt, dann liegen sämtliche Eigenwerte von $\mathbf{A} - \mathbf{B} \mathbf{K}_x^*$ innerhalb des Einheitskreises und das mit \mathbf{K}^* geregelte System ist Referenz-Zustands-stabil nach Definition 3.5.

Beweis:

Nach Satz 3.2 und wegen (3.52) folgt, dass mit Kenntnis von $Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k)$ auch die optimale Q-Funktion $Q_{\text{reg}}^*(\mathbf{x}_k, \mathbf{u}_k)$ des Regulierungsproblems ohne Vorgabe einer Solltrajektorie bekannt ist. Wegen

$$\begin{aligned} V^*(\mathbf{x}_k) &= Q_{\text{reg}}^*(\mathbf{x}_k, \boldsymbol{\mu}^*(\mathbf{x}_k)) \\ &= \mathbf{x}_k^\top \mathbf{h}_{xx} \mathbf{x}_k + 2 \mathbf{x}_k^\top \mathbf{h}_{xu} \boldsymbol{\mu}^*(\mathbf{x}_k) + \boldsymbol{\mu}^{*\top}(\mathbf{x}_k) \mathbf{h}_{uu} \boldsymbol{\mu}^*(\mathbf{x}_k) \\ &= \mathbf{x}_k^\top \mathbf{h}_{xx} \mathbf{x}_k - 2 \mathbf{x}_k^\top \mathbf{h}_{xu} \mathbf{h}_{uu}^{-1} \mathbf{h}_{ux} \mathbf{x}_k + \mathbf{x}_k^\top \mathbf{h}_{xu} \mathbf{h}_{uu}^{-1} \mathbf{h}_{uu} \mathbf{h}_{uu}^{-1} \mathbf{h}_{ux} \mathbf{x}_k \\ &= \mathbf{x}_k^\top \underbrace{(\mathbf{h}_{xx} - \mathbf{h}_{xu} \mathbf{h}_{uu}^{-1} \mathbf{h}_{ux})}_{=: \mathbf{P}^*} \mathbf{x}_k \end{aligned} \quad (3.59)$$

stellt \mathbf{P}^* die positiv definite Lösung der diskreten algebraischen Riccati-Gleichung (3.54) dar und das durch

$$\mathbf{x}_{k+1} = \underbrace{\sqrt{\gamma} \mathbf{A}}_{=: \mathbf{A}_\gamma} \mathbf{x}_k - \underbrace{\sqrt{\gamma} \mathbf{B} \mathbf{K}_x^*}_{=: \mathbf{B}_\gamma} \mathbf{x}_k \quad (3.60)$$

gegebene System wird durch \mathbf{K}_x^* stabilisiert. Somit liegen alle Eigenwerte von $\mathbf{A}_\gamma - \mathbf{B}_\gamma \mathbf{K}_x^*$ innerhalb des Einheitskreises. Im Folgenden wird überprüft, ob für ein gegebenes γ auch Stabilität des durch $\frac{1}{\sqrt{\gamma}} (\mathbf{A}_\gamma - \mathbf{B}_\gamma \mathbf{K}_x^*) = \mathbf{A} - \mathbf{B} \mathbf{K}_x^*$ beschriebenen geschlossenen Regelkreises gilt.

Hierzu werde zunächst untersucht, für welche Konstanten $c \geq 1$ das System

$$\mathbf{x}_{k+1} = c (\mathbf{A}_\gamma - \mathbf{B}_\gamma \mathbf{K}_x^*) \mathbf{x}_k \quad (3.61)$$

stabil ist. Betrachtet werde dazu der Lyapunov-Kandidat $V^*(\mathbf{x}_k) = \mathbf{x}_k^\top \mathbf{P}^* \mathbf{x}_k$. Somit gilt

$$V^*(\mathbf{x}_{k+1}) - V^*(\mathbf{x}_k) = \mathbf{x}_k^\top \underbrace{\left((c \mathbf{A}_\gamma - c \mathbf{B}_\gamma \mathbf{K}_x^*)^\top \mathbf{P}^* (c \mathbf{A}_\gamma - c \mathbf{B}_\gamma \mathbf{K}_x^*) - \mathbf{P}^* \right)}_{=: \tilde{\mathbf{M}}} \mathbf{x}_k \quad (3.62)$$

und (3.61) ist garantiert stabil, wenn \bar{M} negativ definit ist. Durch Umformung ergibt sich

$$\begin{aligned}
\bar{M} &= c^2 A_\gamma^\top P^* A_\gamma - P^* - c^2 K_x^{*\top} B_\gamma^\top P^* A_\gamma - c^2 A_\gamma^\top P^* B_\gamma K_x^* + c^2 K_x^{*\top} B_\gamma^\top P^* B_\gamma K_x^* \\
&\stackrel{(3.54)}{=} -c^2 Q + c^2 P^* - P^* + c^2 A_\gamma^\top P^* B_\gamma (B_\gamma^\top P^* B_\gamma + R)^{-1} B_\gamma^\top P^* A_\gamma \\
&\quad - c^2 K_x^{*\top} B_\gamma^\top P^* A_\gamma - c^2 A_\gamma^\top P^* B_\gamma K_x^* + c^2 K_x^{*\top} B_\gamma^\top P^* B_\gamma K_x^* \\
&= -c^2 Q - (1 - c^2) P^* - c^2 K_x^{*\top} B_\gamma^\top P^* A_\gamma + c^2 K_x^{*\top} B_\gamma^\top P^* B_\gamma K_x^* \\
&= -c^2 Q - (1 - c^2) P^* - c^2 K_x^{*\top} (B_\gamma^\top P^* B_\gamma + R) K_x^* + c^2 K_x^{*\top} B_\gamma^\top P^* B_\gamma K_x^* \\
&= -c^2 Q - (1 - c^2) P^* - c^2 K_x^{*\top} R K_x^* \\
&= c^2 (P^* - Q - K_x^{*\top} R K_x^*) - P^* \stackrel{!}{<} 0,
\end{aligned} \tag{3.63}$$

wobei die aus $K_x^* = (B_\gamma^\top P^* B_\gamma + R)^{-1} B_\gamma^\top P^* A_\gamma$ resultierende Beziehung $B_\gamma^\top P^* A_\gamma = (B_\gamma^\top P^* B_\gamma + R) K_x^*$ genutzt wird. Somit ist \bar{M} stets negativ definit, wenn

$$1 \leq c^2 < \lambda_{\min} \left(P^* (P^* - Q - K_x^{*\top} R K_x^*)^{-1} \right) \tag{3.64}$$

gilt. Mit $c^2 = \frac{1}{\gamma}$ folgt schließlich die Aussage des Satzes. \square

Das nachfolgende Beispiel verdeutlicht, dass Satz 3.3 ein hinreichendes Stabilitätskriterium für den geschlossenen Regelkreis $x_{k+1} = (A - BK_x^*) x_k$ liefert.

Beispiel 3.1 Für das durch

$$\begin{aligned}
A &= \begin{bmatrix} 3 & 0 \\ 0 & 0,4 \end{bmatrix}, & B &= \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}, \\
Q &= \begin{bmatrix} 1 & 2 \\ 2 & 5 \end{bmatrix}, & R &= \begin{bmatrix} 10 & 0 \\ 0 & 1 \end{bmatrix}
\end{aligned} \tag{3.65}$$

und unterschiedliche Diskontierungen γ gegebene Optimierungsproblem folgen die Zusammenhänge in Tabelle 3.1. Hierbei ist $\gamma > b_\gamma$ nach Satz 3.3 hinreichend für die Stabilität von $A - BK_x^*$, ohne A und B explizit zu kennen.

γ	$ \text{eig}(\mathbf{A}_\gamma - \mathbf{B}_\gamma \mathbf{K}_x^*) $	$ \text{eig}(\mathbf{A} - \mathbf{B} \mathbf{K}_x^*) $	b_γ	$\mathbf{A} - \mathbf{B} \mathbf{K}_x^*$ stabil	$\mathbf{A} - \mathbf{B} \mathbf{K}_x^*$ stabil nach Satz 3.3
1	0,4883; 0,0286	0,4883; 0,0286	0,4337	ja	ja
0,6	0,6055; 0,0354	0,7818; 0,0457	0,5498	ja	ja
0,5	0,6484; 0,0380	0,9170; 0,0537	0,5836	ja	k. A. ⁵²
0,4	0,6976; 0,0412	1,1031; 0,0651	0,6142	nein	k. A.

Tabelle 3.1: Eigenwerte des fiktiven Systems $\mathbf{A}_\gamma - \mathbf{B}_\gamma \mathbf{K}_x^*$ sowie des realen Systems $\mathbf{A} - \mathbf{B} \mathbf{K}_x^*$ bei Regelung mit \mathbf{K}_x^* , welches aus der optimalen Lösung des mit γ diskontierten linear-quadratischen Optimierungsproblems mit den Parametern aus (3.65) resultiert. Die hinreichende untere Schranke b_γ folgt aus Satz 3.3.

Gemäß Satz 3.2 können, falls \mathbf{H} (oder alternativ \mathbf{w}^*) bekannt ist, sowohl Q^* als auch $\boldsymbol{\mu}^*$ direkt berechnet werden. Deshalb ist das Ziel des PRADP-Ansatzes, \mathbf{w}^* mithilfe von Datentupeln \mathcal{T}_k , $k = 1, \dots, M$, (vgl. (3.27a)) zu bestimmen. Da nach Lemma 3.2 die optimale Q-Function $Q^*(\mathbf{x}_k, \mathbf{Z}_k, \mathbf{u}_k)$ zu Problem 3.2 quadratisch bezüglich $\bar{\mathbf{y}}_k$ ist, können hierfür quadratische Funktionsapproximatoren gewählt werden. Zudem liefert (3.49) in Satz 3.2 einen analytischen Zusammenhang zwischen \mathbf{H} und dem optimalen Regelgesetz $\boldsymbol{\mu}^*(\mathbf{x}_k, \mathbf{Z}_k)$. Daher können \mathbf{w}^* und somit auch der optimale Regler $\boldsymbol{\mu}^*(\mathbf{x}_k, \mathbf{Z}_k)$ datenbasiert und ohne explizite Kenntnis der Systemmatrizen \mathbf{A} und \mathbf{B} gelernt werden – beispielsweise durch Verwendung der LSPI in Algorithmus 3.1. Hierbei ergibt sich der Policy-Improvement-Schritt direkt zu

$$\hat{\boldsymbol{\mu}}^{[l+1]}(\mathbf{x}_k, \mathbf{Z}_k) = - \left(\mathbf{h}_{\text{uu}}^{[l+1]} \right)^{-1} \begin{bmatrix} \mathbf{h}_{\text{ux}}^{[l+1]} & \mathbf{h}_{\text{uz}}^{[l+1]} \end{bmatrix} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{z}_{k,1:n} \end{bmatrix}. \quad (3.66)$$

Eine konkrete ADP-Implementierung wird im Folgenden präsentiert. Zudem werden Simulationsergebnisse gezeigt und diskutiert.

3.2.5 Implementierung und simulativer Vergleich des PRADP

In diesem Abschnitt werden die Vorteile des vorgestellten PRADP-Verfahrens gegenüber einem weitverbreiteten Ansatz aus der Literatur, der annimmt, der Solltrajektorienverlauf folge global einer unbekanntenen Exosystemdynamik, simulativ gezeigt. Dazu wird zunächst beschrieben, wie ein beliebiger Solltrajektorienverlauf durch die beispielhafte und ADP-kompatible Wahl kubischer Polynome aus lokaler Sicht im Zeitschritt k approximiert werden kann. Anschließend wird ein Beispielproblem definiert. Für dieses Beispielproblem wird einerseits der in dieser Arbeit vorgestellte PRADP-Ansatz und andererseits die besagte Vergleichsmethode aus der Literatur trainiert. Für beide resultierenden Regler werden Simulationsergebnisse zu drei unterschiedlichen Szenarien betrachtet. Eine vergleichende Diskussion schließt diesen Abschnitt ab.

⁵² Keine Aussage möglich, Lyapunov-Kandidat liefert ein hinreichendes Kriterium.

3.2.5.1 Lokale Solltrajektorienapproximation durch kubische Polynome

Der im vorliegenden Kapitel präsentierte Ansatz löst optimale Solltrajektorienfolgeregelungsprobleme, deren Gütefunktional ADP-kompatible Solltrajektorien berücksichtigt (vgl. (3.37)). Um den anwendungsnahen Fall beliebiger Sollvorgaben $\mathbf{x}_r, \text{soll}, k$ handhaben zu können, wird im Folgenden ein Verfahren vorgestellt, das aus einem beliebigen Sollzustandsverlauf, der auf einem endlichen, gleitenden Horizont der Länge n_h gegeben ist, in jedem Zeitschritt k eine ADP-kompatible Approximation erzeugt. Exemplarisch⁵³ werden hierzu kubische Polynome gewählt, die einen Kompromiss zwischen der Approximationsfähigkeit in der lokalen Nachbarschaft des aktuellen Zeitschritts und einer moderaten Zunahme der Anzahl zu lernender Gewichte darstellen.

Damit $\mathbf{x}_r(\mathbf{Z}_k^{(\kappa)})$ kubischen Polynomen entspricht, werde

$$\boldsymbol{\rho}(\kappa) = [(\kappa\Delta t)^3 \quad (\kappa\Delta t)^2 \quad \kappa\Delta t \quad 1]^\top \quad (3.67)$$

gewählt, wobei Δt die Abtastzeit darstellt⁵⁴. Eine Transformationsmatrix $\mathbf{D}(\kappa)$, die zu einer verschobenen Parametermatrix $\mathbf{Z}_k^{(\kappa)}$ nach Definition 3.2 führt, folgt aus dem durch

$$\begin{aligned} \mathbf{x}_r(\mathbf{Z}_k^{(\kappa+j)}) &= \mathbf{x}_r(\mathbf{Z}_k, \kappa + j) \\ &= \mathbf{Z}_k \boldsymbol{\rho}(\kappa + j) \\ &= \mathbf{Z}_k \begin{bmatrix} ((\kappa + j)\Delta t)^3 \\ ((\kappa + j)\Delta t)^2 \\ (\kappa + j)\Delta t \\ 1 \end{bmatrix} \\ &= \mathbf{Z}_k \underbrace{\begin{bmatrix} 1 & 3\kappa\Delta t & 3(\kappa\Delta t)^2 & (\kappa\Delta t)^3 \\ 0 & 1 & 2\kappa\Delta t & (\kappa\Delta t)^2 \\ 0 & 0 & 1 & \kappa\Delta t \\ 0 & 0 & 0 & 1 \end{bmatrix}}_{=: \mathbf{D}(\kappa)} \boldsymbol{\rho}(j) \\ &= \mathbf{Z}_k^{(\kappa)} \boldsymbol{\rho}(j) \end{aligned} \quad (3.68)$$

⁵³ Die Approximationstheorie stellt ein eigenes Forschungsgebiet dar und soll in der vorliegenden Arbeit nicht vertieft werden. Für Untersuchungen hinsichtlich der Approximationsfähigkeit von Polynomen auf einem endlichen Intervall sei z. B. auf [DS10] verwiesen. Während nachfolgend beispielhaft kubische Polynome verwendet werden, werden im Anwendungskapitel (Kapitel 6) lineare (Abschnitt 6.1) und quadratische (Abschnitt 6.2) Polynome mit konstanten Sollzustandsvorgaben verglichen. Eine geeignete Wahl der ADP-kompatiblen Solltrajektorienapproximation, wie der zugehörige Polynomgrad oder der gleitende Vorausschauhorizont n_h , hängt insbesondere von der konkreten Anwendung ab. Somit können für jede konkrete Problemstellung beispielsweise der gewünschte Solltrajektorienverlauf oder auch mögliches Vorwissen über die Zeitkonstante des zu regelnden Systems bei der Wahl eines ADP-kompatiblen Solltrajektorienapproximators berücksichtigt werden.

⁵⁴ Andere Approximationen können durch eine andere Wahl von $\boldsymbol{\rho}(\kappa)$ erzielt werden. Beispielsweise führt $\boldsymbol{\rho}(\kappa) = [\kappa\Delta t \quad 1]^\top$ zu linearer Interpolation und $\boldsymbol{\rho}(\kappa) = 1$ zu einer konstanten Sollzustandsvorgabe.

gegebenen Zusammenhang⁵⁵. Für einen beliebigen Solltrajektorienverlauf $\mathbf{x}_{r,\text{soll},k}$ wird ein Parameter \mathbf{Z}_k in jedem Zeitschritt k benötigt, sodass

$$\mathbf{x}_r(\mathbf{Z}_k^{(\kappa)}) = \begin{bmatrix} x_{r,1}(\mathbf{Z}_k^{(\kappa)}) \\ \vdots \\ x_{r,n}(\mathbf{Z}_k^{(\kappa)}) \end{bmatrix}, \quad \kappa = 0, 1, \dots, \quad (3.69)$$

eine Approximation von

$$\mathbf{x}_{r,\text{soll},k+\kappa} = \begin{bmatrix} x_{r,\text{soll},k+\kappa,1} \\ \vdots \\ x_{r,\text{soll},k+\kappa,n} \end{bmatrix} \quad (3.70)$$

darstellt. Um diesen Parameter \mathbf{Z}_k zu bestimmen, der den lokalen Solltrajektorienverlauf approximiert, sei die beliebige Solltrajektorie zur Laufzeit über einen gleitenden Vorausschauhorizont von n_h Zeitschritten bekannt, d. h. $\mathbf{x}_{r,\text{soll},k:k+n_h-1}$ sei gegeben. Dann kann \mathbf{Z}_k beispielsweise mithilfe einer gewichteten Least-Squares-Regression bestimmt werden. Sei hierzu

$$\mathbf{x}_{r,\text{soll},k:k+n_h-1} := [\mathbf{x}_{r,\text{soll},k} \quad \dots \quad \mathbf{x}_{r,\text{soll},k+n_h-1}],$$

$$\boldsymbol{\rho}_{0:n_h-1} := \begin{bmatrix} \boldsymbol{\rho}^\top(0) \\ \boldsymbol{\rho}^\top(1) \\ \dots \\ \boldsymbol{\rho}^\top(n_h-1) \end{bmatrix} \quad \text{und} \quad \mathbf{W}_p := \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & \beta & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \beta^{n_h-1} \end{bmatrix}, \quad (3.71)$$

wobei \mathbf{W}_p die Gewichtungsmatrix mit den Least-Squares-Gewichtungen β darstellt. Nachfolgend werde beispielhaft $\beta = \gamma$ gewählt, sodass analog zur Diskontierung der Kosten (vgl. J_k in (3.16) und (3.37)) zukünftige Zeitschritte weniger wichtig bei der Least-Squares-Regression sind. Folglich ergibt sich die lokale Solltrajektorienapproximation durch den Parameter

$$\mathbf{Z}_k = \mathbf{x}_{r,\text{soll},k:k+n_h-1} \underbrace{\mathbf{W}_p \boldsymbol{\rho}_{0:n_h-1} (\boldsymbol{\rho}_{0:n_h-1}^\top \mathbf{W}_p \boldsymbol{\rho}_{0:n_h-1})^{-1}}_{=: \mathbf{P}_{\text{LS}}}. \quad (3.72)$$

Aus Anwendungssicht besonders günstig ist hierbei die Tatsache, dass sich bei gegebenem n_h und gegebenem $\boldsymbol{\rho}(\cdot)$ die Matrix \mathbf{P}_{LS} im Voraus berechnen lässt, sodass sich die Berechnung von \mathbf{Z}_k nach (3.72) auf eine einzelne Matrixmultiplikation beschränkt. Für den exemplarisch gewählten gleitenden Vorausschauhorizont $n_h = 10$, der auch im folgenden Simulationsbeispiel verwendet wird, ist ein beispielhafter Ausschnitt eines Solltrajektorienverlaufs $x_{r,\text{soll},k,1}$ des ersten Zustands x_1 sowie dessen lokale Approximation für die Beispielzeitschritte $k = 105$ und $k = 106$ in Abbildung 3.2 veranschaulicht.

⁵⁵ Mittels Koeffizientenvergleich ist zudem ersichtlich, dass für Polynome vom Grad d allgemein

$$\mathbf{D}(\kappa) = \begin{bmatrix} \binom{d}{0}(\kappa\Delta t)^0 & \binom{d}{1}(\kappa\Delta t)^1 & \dots & \binom{d}{d}(\kappa\Delta t)^d \\ 0 & \binom{d-1}{0}(\kappa\Delta t)^0 & \dots & \binom{d-1}{d-1}(\kappa\Delta t)^{d-1} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & \binom{0}{0}(\kappa\Delta t)^0 \end{bmatrix} \text{ gewählt werden kann.}$$

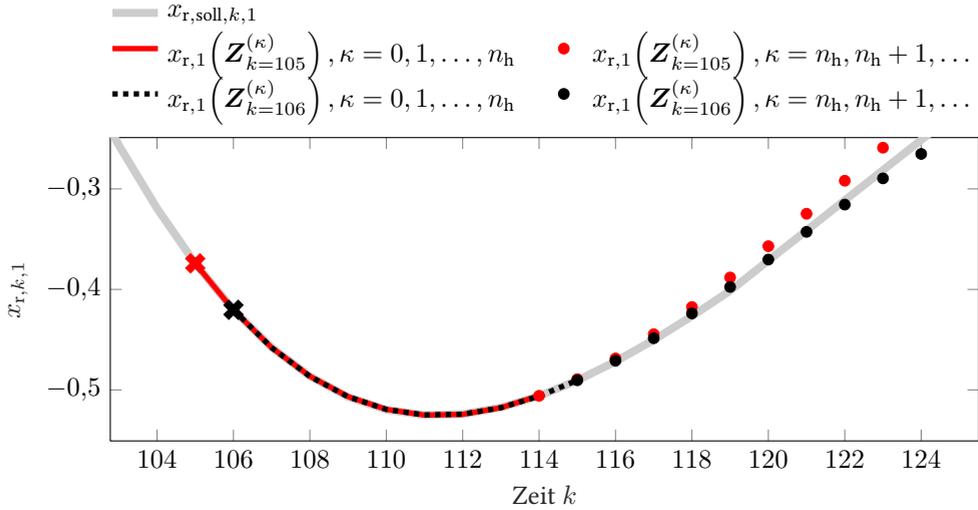


Abbildung 3.2: Beispielhafter Ausschnitt des Solltrajektorienverlaufs $x_{r,soll,k,1}$ des ersten Zustands x_1 (grau) sowie die lokale Approximation durch kubische Polynome zu den Zeitschritten $k = 105$ (rot) und $k = 106$ (schwarz), die für $Z_{k=105}$ bzw. $Z_{k=106}$ resultieren. Der Vorausschauhorizont, der für die lokale Solltrajektorienapproximation, d. h. die Anpassung des Solltrajektorienparameters Z_k in jedem Zeitschritt k , genutzt wird, ist in diesem Beispiel $n_h = 10$ (rote durchgezogene und schwarze gestrichelte Linie).

3.2.5.2 Beispielsystem und Gütefunktional

Im Folgenden werde ein Feder-Masse-Dämpfer-System mit der Systemdynamik

$$\dot{\mathbf{x}}(t) = \begin{bmatrix} 0 & 1 \\ -\frac{c_{\text{sys}}}{m_{\text{sys}}} & -\frac{d_{\text{sys}}}{m_{\text{sys}}} \end{bmatrix} \mathbf{x}(t) + \begin{bmatrix} 0 \\ \frac{1}{m_{\text{sys}}} \end{bmatrix} u(t) \quad (3.73)$$

mit $m_{\text{sys}} = 0,5 \text{ kg}$, $c_{\text{sys}} = 0,1 \text{ N m}^{-1}$ und $d_{\text{sys}} = 0,1 \text{ kg s}^{-1}$ betrachtet. Die Diskretisierung dieses Systems mittels Tustin-Approximation mit einer Abtastzeit von $\Delta t = 0,1 \text{ s}$ liefert

$$\mathbf{x}_{k+1} = \begin{bmatrix} 0,9990 & 0,0990 \\ -0,0198 & 0,9792 \end{bmatrix} \mathbf{x}_k + \begin{bmatrix} 0,0099 \\ 0,1979 \end{bmatrix} u_k. \quad (3.74)$$

Durch x_1 wird die Position und durch x_2 die Geschwindigkeit der Masse m_{sys} beschrieben, während die Stellgröße u_k einer auf die Masse wirkenden Kraft entspricht. Die durch (3.73) bzw. (3.74) beschriebene Systemdynamik sei dem Regler im Folgenden nicht bekannt.

Mit dem Ziel, die Position der Masse (d. h. x_1) einem Solltrajektorienverlauf folgen zu lassen, werde

$$\mathbf{Q} = \begin{bmatrix} 100 & 0 \\ 0 & 0 \end{bmatrix} \text{ und } \mathbf{R} = 1 \quad (3.75)$$

gesetzt, um Abweichungen des ersten Systemzustands vom parametrisierten Sollzustandsverlauf nach (3.37) zu bestrafen. Der Diskontierungsfaktor werde zu $\gamma = 0,9$ gewählt. Für dieses Beispiel sind Annahme 3.1 und Annahme 3.2 erfüllt, d. h. nach Satz 3.1 existiert der optimale Regler $\mu^* \left(\mathbf{x}_{k+\kappa}, \mathbf{Z}_k^{(\kappa)} \right)$ und ist zudem eindeutig. Außerdem folgt aus Satz 3.3 $b_\gamma = 0,6263$. Somit führt die optimale Lösung des mit $\gamma = 0,9 > b_\gamma$ diskontierten Optimierungsproblems zu einem nach Definition 3.5 Referenz-Zustands-stabilen System.

3.2.5.3 Trainingsvorgang

Der in diesem Kapitel vorgestellte ADP-basierte Solltrajektorienfolgeregler, der mit einer ADP-kompatiblen, parametrischen Darstellung des Sollzustandsverlaufs arbeitet, wird im Folgenden trainiert, um anschließend validiert zu werden. Zu Vergleichszwecken wird zudem ein ADP-Trajektorienfolgeregelungsansatz verwendet, dem, genau wie in den Arbeiten von Luo et al. [LLHW16] und Kiumarsi et al. [KLM⁺14], die Annahme zugrunde liegt, der Solltrajektorienverlauf werde direkt und global durch ein Exosystem $\mathbf{x}_{r,k+1} = \mathbf{f}_{\mathbf{x}_r}(\mathbf{x}_{r,k})$ erzeugt (vgl. auch Abschnitt 2.2.2). Somit hängt die Q-Function $Q(\mathbf{x}_k, \mathbf{x}_{r,k}, \mathbf{u}_k)$ (vgl. [KLM⁺14]) nur vom aktuellen Sollzustand $\mathbf{x}_{r,k}$ im Zeitschritt k ab, wohingegen der Parameter \mathbf{Z}_k der vorgestellten PRADP-Methode potenziell mehr Information über den Sollzustandsverlauf beinhalten kann. Für die Vergleichsmethode ist nach [KLM⁺14, Abschnitt 5.1] die Q-Function $Q(\mathbf{x}_k, \mathbf{x}_{r,k}, \mathbf{u}_k)$ quadratisch bezüglich \mathbf{x}_k , \mathbf{u}_k und $\mathbf{x}_{r,k}$, sofern die Exosystemdynamik $\mathbf{f}_{\mathbf{x}_r}(\mathbf{x}_{r,k})$ linear ist, und die Policy Iteration nach [KLM⁺14, Algorithm 3] kann durchgeführt werden.

Die beiden untersuchten modellfreien ADP-Trajektorienfolgeregelungsmethoden werden jeweils mit $M = 500$ Datentupeln \mathcal{T}_k trainiert. Während des Datenaufzeichnens wird zur Systemanregung Rauschen, das aus einer Normalverteilung mit Mittelwert 0 und Standardabweichung 1 erzeugt wird, zur Stellgröße u_k addiert. Der Solltrajektorienverlauf während des Trainings werde durch

$$\begin{bmatrix} x_{r,\text{soll},k+1,1} \\ x_{r,\text{soll},k+1,2} \end{bmatrix} = \mathbf{x}_{r,\text{soll},k+1} = \mathbf{f}_{\mathbf{x}_r}(\mathbf{x}_{r,\text{soll},k}) = \underbrace{\begin{bmatrix} 0,9988 & 0,0500 \\ -0,0500 & 0,9988 \end{bmatrix}}_{\mathbf{F}_{\text{ref}}} \mathbf{x}_{r,\text{soll},k} \quad (3.76)$$

mit dem Initialzustand $\mathbf{x}_{r,\text{soll},0} = [0 \ 1]^\top$ beschrieben⁵⁶. Die Vergleichsmethode gemäß [KLM⁺14] bzw. [LLHW16] verwendet $\mathbf{x}_{r,k} = \mathbf{x}_{r,\text{soll},k}$, wohingegen für die PRADP-Methode der nach Abschnitt 3.2.5.1 bestimmte Parameter \mathbf{Z}_k , der ein kubisches Polynom parametrisiert und den Sollzustandsverlauf ab dem jeweiligen Zeitschritt k beschreibt, genutzt wird.

Der PRADP-basierte Regler wird dann entsprechend Algorithmus 3.1 trainiert, wobei die Abbruchbedingung $e_{\hat{\mathbf{w}}} = 10^{-6}$ gewählt wird. Der Policy-Improvement-Schritt erfolgt nach (3.66). Aufgrund von Lemma 3.2 werden Basisfunktionen $\phi(\cdot)$ gewählt, die quadratisch bezüglich \mathbf{x}_k , \mathbf{u}_k und $\mathbf{z}_{k,1:n}$ sind. Weiterhin wird mangels besseren Wissens $\hat{\mathbf{w}}^{[0]}$ so initialisiert, dass

⁵⁶ Diese lineare Exosystemdynamik wird insbesondere verwendet, um die Voraussetzungen der Vergleichsmethode zu erfüllen. Für die PRADP-Methode wäre dies, wie auch im Anwendungskapitel 6 ersichtlich, grundsätzlich nicht nötig, da die Verwendung von \mathbf{Z}_k und $\mathbf{Z}_k^{(1)}$ in (3.26) ADP-Kompatibilität gewährleistet.

$\hat{\boldsymbol{\mu}}^{[0]} = \mathbf{0}$ gilt⁵⁷. Der gesamte Ablauf ist in Abbildung 3.3 veranschaulicht. Die Vergleichsmethode ist wie in [KLM⁺14, Algorithm 3] beschrieben implementiert, wobei ebenfalls die Abbruchbedingung $e_{\hat{\boldsymbol{w}}} = 10^{-6}$ sowie $\hat{\boldsymbol{\mu}}^{[0]} = \mathbf{0}$ genutzt wird.

3.2.5.4 Untersuchung der gelernten Regler

Zunächst werden die aus Daten und ohne Kenntnis der Systemmatrizen \mathbf{A} und \mathbf{B} gelernten ADP-Regler mit den analytisch berechneten Grundwahrheiten verglichen. Anschließend werden verschiedene Simulationsszenarien betrachtet, um Unterschiede zwischen dem PRADP-Regler und dem Vergleichsverfahren hinsichtlich ihrer Flexibilität und Performanz herauszuarbeiten.

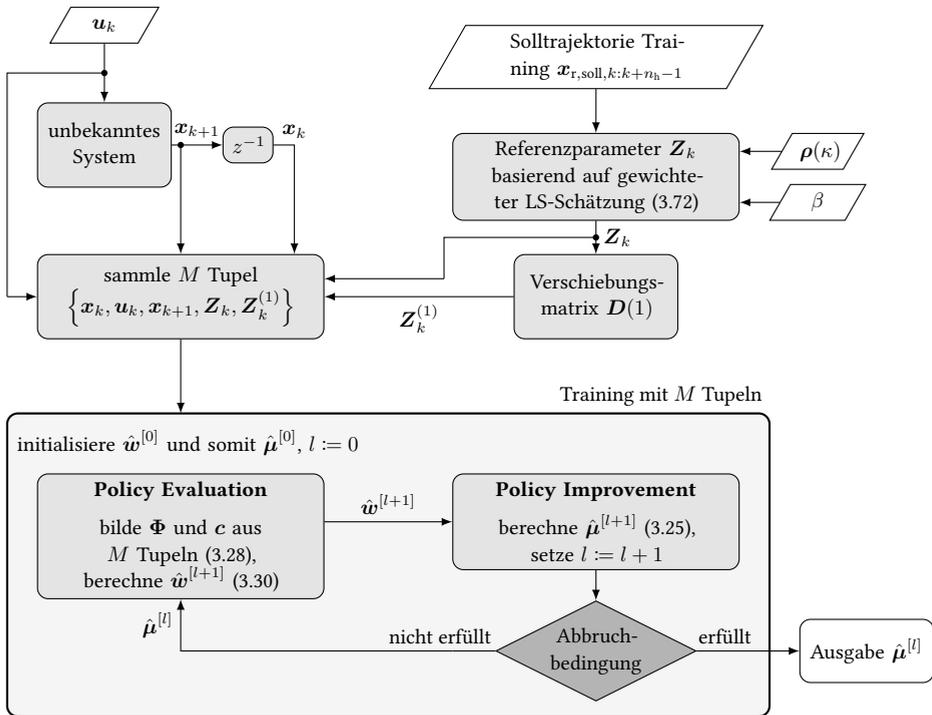


Abbildung 3.3: Ablaufschema des Datenaufzeichnungs- und Trainingsprozesses des PRADP-Algorithmus. Eingabe-größen sind hierbei die Stellgröße u_k und Solltrajektorie $x_{r,soll,k:k+n_h-1}$ während des Trainings, die Basisfunktionen $\rho(\kappa)$ zur Referenzapproximation und der Parameter β für die gewichtete Least-Squares-Approximation. Sobald die Abbruchbedingung (siehe Algorithmus 3.1) erfüllt ist, wird der geschätzte optimale Regler $\hat{\boldsymbol{\mu}}^{[l]}$ ausgegeben.

⁵⁷ Dies kann erreicht werden, indem die Gewichte, die zu h_{uu} gehören, so gesetzt werden, dass $h_{uu} > \mathbf{0}$ gilt, wohingegen alle anderen Gewichte zu null gesetzt werden (vgl. Lemma 3.2).

Für die Verwendung der parametrisierten Solltrajektorie ergibt sich die optimale Reglermatrix \mathbf{K}^* , die unter Verwendung des vollständigen Systemwissens nach Satz 3.1 und Bemerkung 3.5 berechnet werden kann, zu

$$\boldsymbol{\mu}^*(\mathbf{x}_k, \mathbf{Z}_k) = - \underbrace{\begin{bmatrix} 6,30 & 2,26 & -0,31 & -0,97 & -2,37 & -6,40 & 0 & 0 & 0 & 0 \end{bmatrix}}_{\mathbf{K}^*} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{z}_{k,1:n} \end{bmatrix}. \quad (3.77)$$

Der Vergleich des gelernten PRADP-Reglers $\mathbf{K}_{\text{PRADP}}$ mit dieser Grundwahrheit \mathbf{K}^* liefert $\|\mathbf{K}_{\text{PRADP}} - \mathbf{K}^*\|_2 = 6,51 \cdot 10^{-14}$. Somit stimmt der ohne Kenntnis der Systemdynamik gelernte PRADP-Regler bis auf numerische Ungenauigkeiten mit der optimalen Lösung überein.

Für die Vergleichsmethode ergibt sich der optimale Regler, nachfolgend durch $\boldsymbol{\mu}_{\text{exo}}^*(\cdot)$ und die Reglermatrix $\mathbf{K}_{\text{exo}}^*$ bezeichnet, für das gegebene Beispielproblem nach [KLM⁺14, (58)] zu

$$\boldsymbol{\mu}_{\text{exo}}^*(\mathbf{x}_k, \mathbf{x}_{r,k}) = - \underbrace{\begin{bmatrix} 6,30 & 2,26 & -6,28 & -1,18 \end{bmatrix}}_{\mathbf{K}_{\text{exo}}^*} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{x}_{r,k} \end{bmatrix}. \quad (3.78)$$

Die Reglermatrix $\mathbf{K}_{\text{exo,ADP}}$, die wie in Abschnitt 3.2.5.3 beschrieben aus Daten gelernt wird, stimmt ebenfalls mit der Grundwahrheit überein, da $\|\mathbf{K}_{\text{exo,ADP}} - \mathbf{K}_{\text{exo}}^*\|_2 = 2,60 \cdot 10^{-13}$ gilt. Somit liegt auch hier ein erfolgreich gelernter Regler vor.

Um die Performanz des in dieser Arbeit vorgestellten PRADP-Reglers und der Vergleichsmethode zu untersuchen, werden drei unterschiedliche Szenarien betrachtet. In allen Szenarien werden die zuvor gelernten Regler ohne weitere Anpassungen verwendet, zudem wird jeweils der Anfangszustand zu $\mathbf{x}_0 = [1 \ 0]^\top$ gesetzt. Die unterschiedlichen Szenarien werden im Folgenden beschrieben.

1. Der Sollzustandsverlauf wird durch dieselbe Exosystemdynamik erzeugt, die während des Trainings verwendet wurde, d. h. \mathbf{F}_{ref} wie in (3.76). Der resultierende Verlauf des Zustands x_1 ist in Abbildung 3.4a gezeigt. In Abbildung 3.4b sind die Parameter $\mathbf{z}_{k,1} = [z_{k,1,1} \ \dots \ z_{k,1,n_z}]^\top$ mit $n_z = 4$ gegeben. Diese parametrieren das kubische Polynom, das $x_{r,1}(\mathbf{Z}_k, 0)$ beschreibt.
2. Der Sollzustandsverlauf wird ebenfalls durch eine zeitinvariante Exosystemdynamik erzeugt, jedoch durch eine andere Dynamik als die während des Trainings verwendete, d. h. es gilt $\mathbf{F}_{\text{ref, val}} \neq \mathbf{F}_{\text{ref}}$. Der Verlauf des Sollzustands wird durch

$$\mathbf{x}_{r,k+1} = \underbrace{\begin{bmatrix} 0,9987 & 0,0030 \\ -0,1998 & 0,9987 \end{bmatrix}}_{=:\mathbf{F}_{\text{ref, val}}} \mathbf{x}_{r,k} \quad (3.79)$$

mit $\mathbf{x}_{r,0} = [10 \ 1]^\top$ beschrieben. Der resultierende Verlauf des Zustands x_1 sowie die Polynomkoeffizienten, die den Sollverlauf approximieren, sind in Abbildung 3.5 gegeben. Um darüber hinaus die Qualität der Trajektorienfolgeeregler im Sinne der resultierenden Kosten vergleichen zu können, sind die Einschrittkosten $r(\mathbf{x}_k, \mathbf{x}_r(\mathbf{Z}_k), \mathbf{u}_k)$ sowie die akkumulierten Kosten $\sum_{i=0}^k r(\mathbf{x}_i, \mathbf{x}_r(\mathbf{Z}_i), \mathbf{u}_i)$ gezeigt.

3. Im dritten Szenario wird der Sollzustandsverlauf nicht durch ein zeitinvariantes Exosystem generiert, sondern eine benutzerdefinierte Vorgabe für $x_{r,\text{soll},k,1}$ verwendet und anschließend, wie in Abschnitt 3.2.5.1 beschrieben, durch $x_{r,k,1}$ approximiert. Der sich ergebende Sollverlauf ist in Abbildung 3.6a gezeigt. Zudem wird in diesem Beispiel $x_{r,k,2} = 0, \forall k$, gesetzt. Dies ist dadurch motiviert, dass der zweite Systemzustand aufgrund der Wahl von Q in (3.75) ohnehin nicht bestraft wird. Die Ergebnisse sind in Abbildung 3.6 veranschaulicht.

3.2.5.5 Diskussion der Simulationsergebnisse

Wie den Abbildungen 3.4a, 3.5a und 3.6a zu entnehmen ist, kann der PRADP-Regler in allen drei Szenarien dem Sollzustandsverlauf folgen. Da der Regler durch den Parameter Z_k nicht nur Informationen über den im Zeitschritt k aktuellen Sollzustand, sondern auch über dessen (lokal approximierten) zukünftigen Verlauf besitzt, weist der PRADP-Regler vorausschauendes Verhalten auf, anstatt rein reaktiv zu handeln. Dies ist in den Abbildungen 3.4a, 3.5a und 3.6a zu sehen, da der Systemzustand dem Sollverlauf direkt folgt und es zu keiner zeitlichen Verzögerung kommt.

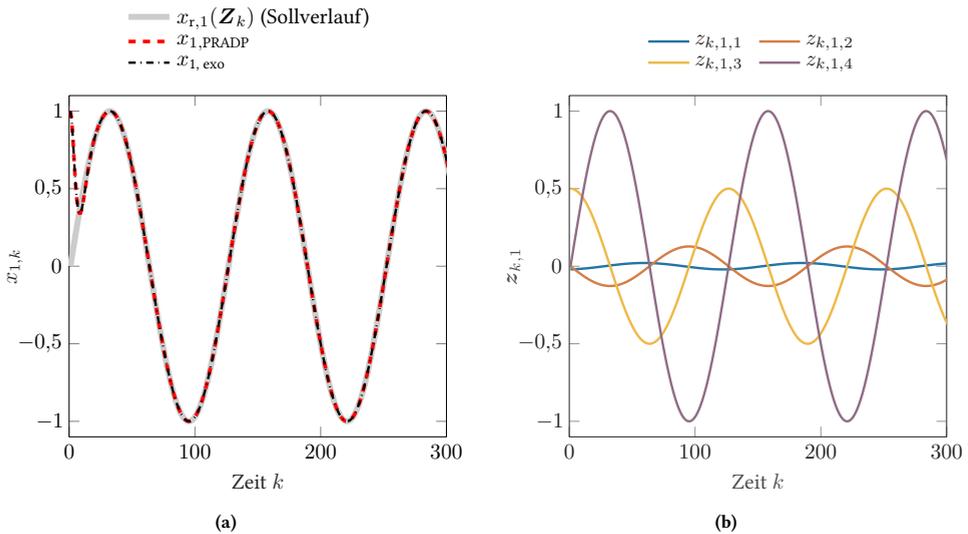


Abbildung 3.4: Ergebnis der vorgestellten PRADP-Methode im Vergleich zu einem Regler nach [LLHW16], [KLM⁺14] für Szenario 1, bei dem der Solltrajektorienverlauf durch F_{ref} generiert wird.

(a): Der approximierte Solltrajektorienverlauf $x_{r,1}(Z_k, 0)$ ist in Grau dargestellt, der resultierende Zustand bei Verwendung des vorgestellten PRADP-Algorithmus in Rot und das Ergebnis der Vergleichsmethode in Schwarz.

(b): Parametervektor $z_{k,1}$ der kubischen Polynome, welche den Solltrajektorienverlauf $x_{r,1}(Z_k, 0)$ beschreiben.

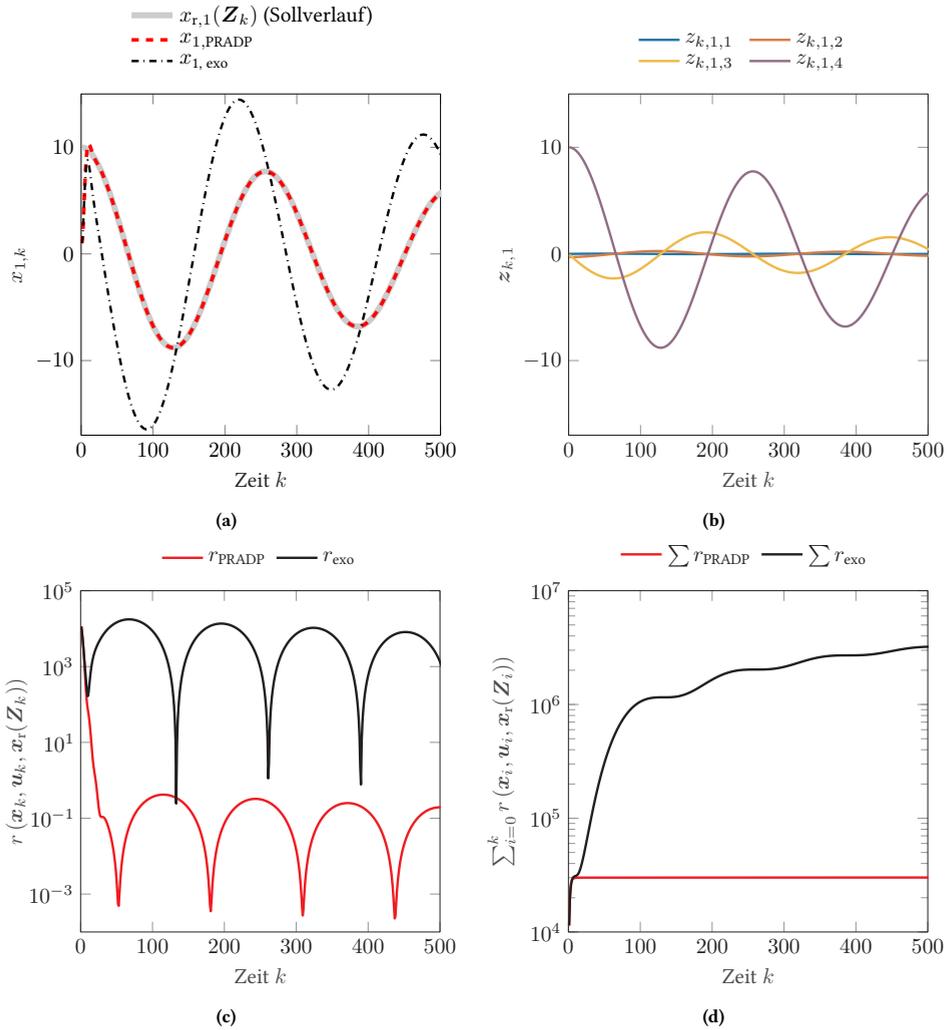


Abbildung 3.5: Ergebnis der vorgestellten PRADP-Methode im Vergleich zu einem Regler nach [LLHW16], [KLM⁺14] für Szenario 2, bei dem der Solltrajektorienverlauf durch $\mathbf{F}_{ref, val} \neq \mathbf{F}_{ref}$ generiert wird.

(a): Der approximierte Solltrajektorienverlauf $x_{r,1}(\mathbf{Z}_k, 0)$ ist in Grau dargestellt, der resultierende Zustand bei Verwendung des vorgestellten PRADP-Algorithmus in Rot und das Ergebnis der Vergleichsmethode in Schwarz.

(b): Parametervektor $z_{k,1}$ der kubischen Polynome, welche den Solltrajektorienverlauf $x_{r,1}(\mathbf{Z}_k, 0)$ beschreiben.

(c): Die Einschrittkosten $r(x_k, x_r(\mathbf{Z}_k), u_k)$, welche bei Verwendung der PRADP-Methode entstehen, sind in Rot gezeigt, wohingegen die mit der Vergleichsmethode verbundenen Kosten in Schwarz visualisiert sind (logarithmische Ordinate).

(d): Die akkumulierten Kosten $\sum_{i=0}^k r(x_i, x_r(\mathbf{Z}_i), u_i)$ der PRADP-Methode sind in Rot, die der Vergleichsmethode in Schwarz gegeben (logarithmische Ordinate).

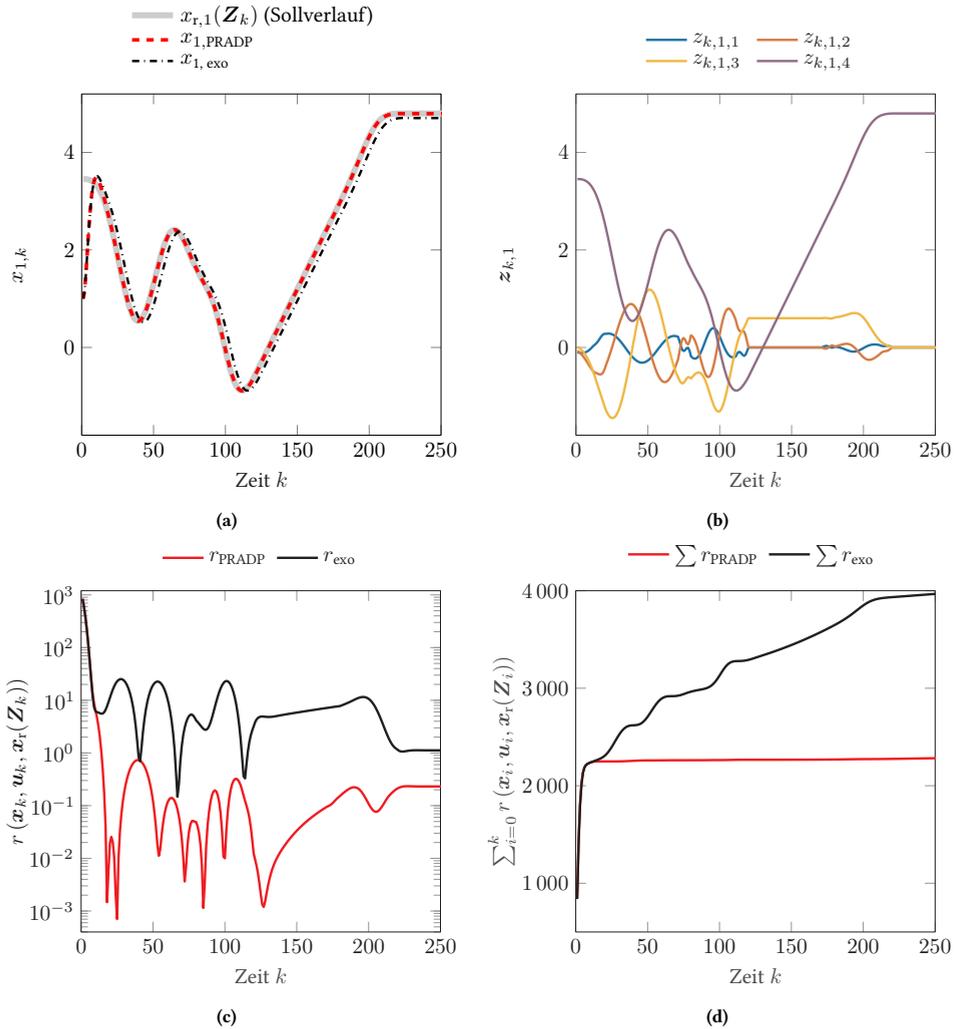


Abbildung 3.6: Ergebnis der vorgestellten PRADP-Methode im Vergleich zu einem Regler nach [LLHW16], [KLM⁺14] für Szenario 3, bei dem der Solltrajektorienverlauf keiner Dynamik folgt, sondern einer beliebigen benutzerdefinierten Vorgabe entspricht.

(a): Der approximierte Solltrajektorienverlauf $x_{r,1}(\mathbf{Z}_k, 0)$ ist in Grau dargestellt, der resultierende Zustand bei Verwendung des vorgestellten PRADP-Algorithmus in Rot und das Ergebnis der Vergleichsmethode in Schwarz.

(b): Parametervektor $z_{k,1}$ der kubischen Polynome, welche den Solltrajektorienverlauf $x_{r,1}(\mathbf{Z}_k, 0)$ beschreiben.

(c): Die Einschrittkosten $r(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_r(\mathbf{Z}_k))$, welche bei Verwendung der PRADP-Methode entstehen, sind in Rot gezeigt, wohingegen die mit der Vergleichsmethode verbundenen Kosten in Schwarz visualisiert sind (logarithmische Ordinate).

(d): Die akkumulierten Kosten $\sum_{i=0}^k r(\mathbf{x}_i, \mathbf{u}_i, \mathbf{x}_r(\mathbf{Z}_i))$ der PRADP-Methode sind in Rot, die der Vergleichsmethode in Schwarz gegeben.

Die Vergleichsmethode nach [LLHW16] bzw. [KLM⁺14] führt zu ähnlichen Ergebnissen, solange der Solltrajektorienverlauf mit der Exosystemdynamik F_{ref} , die während des Trainings verwendet wurde, übereinstimmt, wie dies in Szenario 1 der Fall ist (vgl. Abbildung 3.4). Dieser Vergleichsregler weist jedoch teilweise signifikante Abweichungen vom Solltrajektorienverlauf auf, sobald dieser nicht mehr durch die Dynamik F_{ref} beschrieben wird, d. h. sobald (3.76) nicht mehr gilt. Dies offenbart sich in Szenario 2 (Abbildung 3.5) und Szenario 3 (Abbildung 3.6). Obwohl der Vergleichsregler dem Sollzustand in Szenario 3 näherungsweise folgt⁵⁸, hat der Regler keinerlei Kenntnis über den weiteren Verlauf des Sollzustands. Durch diese mangelnde Vorausschaufähigkeit kann der Regler nur reagieren und es kommt zu einem zeitlichen Versatz zwischen der Trajektorie des Zustands x_1 und dem gewünschten Verlauf.

Zudem offenbaren die Einschrittkosten in den Abbildungen 3.5c und 3.6c sowie die akkumulierten Kosten in den Abbildungen 3.5d und 3.6d, dass die Kosten der Vergleichsmethode signifikant über den Kosten des PRADP-Ansatzes liegen, sobald der Solltrajektorienverlauf nicht mehr durch die während des Trainings verwendete Exosystemdynamik F_{ref} beschrieben wird. Dies ist insbesondere auf größere Abweichungen von der Solltrajektorie zurückzuführen.

Abschließend lässt sich zusammenfassen, dass die beiden betrachteten ADP-Methoden zwar in Szenario 1 vergleichbare Ergebnisse liefern, die in dieser Arbeit vorgestellte PRADP-Methode jedoch in Szenario 2 und Szenario 3 überlegen ist. Dies liegt insbesondere daran, dass der PRADP-Regler zur Laufzeit stets lokal den aktuellen Solltrajektorienverlauf mithilfe des Parameters Z_k approximiert und daher im Gegensatz zum Vergleichsregler nicht annehmen muss, dass der Sollzustandsverlauf derselben Exosystemdynamik wie während des Trainingsvorgangs folgt.

3.3 Zeitdiskrete ADP-kompatible Referenztrajektorie auf einem endlichem Vorausschauhorizont

Nachdem im letzten Abschnitt eine ADP-kompatible parametrisierte Darstellung des Sollzustandsverlaufs vorgestellt und in eine Q-Function integriert wurde, wird im vorliegenden Abschnitt die Verwendung eines beliebigen Sollzustandsverlaufs auf einem endlichen Vorausschauhorizont untersucht⁵⁹. Hierzu wird zunächst eine ADP-kompatible Problemformulierung gegeben. Konkret wird ein linear-quadratisches⁶⁰ Optimierungsproblem mit unendlichem Zeithorizont angesetzt (vgl. Bemerkung 2.1). In dieses Optimierungsproblem geht in jedem Zeitschritt k der beliebige, auf einem gleitenden Vorausschauhorizont der Länge n_h gegebene, Sollzustandsverlauf ein. Um eine ADP-kompatible Darstellung nach Definition 3.1 zu erhalten und das Kostenfunktional auf einem unendlichen Zeithorizont definieren zu können, wird

⁵⁸ Dies ist dem Zustandsrückführungsterm $K_{x,\text{exo}}^*$ des Reglers K_{exo}^* , der identisch zum Zustandsrückführungsterm des PRADP-Reglers ist, sowie dem Regleranteil, welcher der aktuellen Sollzustandsposition zugeordnet ist und der bei beiden Methoden ähnliche Werte aufweist (d. h. $-6,40$ in (3.77) bzw. $-6,28$ in (3.78)), zu verdanken.

⁵⁹ Teile des vorliegenden Abschnitts wurden in einer wissenschaftlichen Fachzeitschrift veröffentlicht [KWFH20].

⁶⁰ Linear-quadratische Optimierungsprobleme sind für zahlreiche regelungstechnische Probleme relevant.

der Sollzustandsverlauf jenseits des Vorausschauhorizonts als konstant angenommen. Für diese Problemstellung lässt sich eine Q-Function definieren, die den beliebigen Sollverlauf auf dem gleitenden Horizont explizit berücksichtigt. Eine genaue Analyse der Form dieser erweiterten Q-Function für linear-quadratische Problemstellungen führt einerseits zu Existenz- und Eindeutigkeitsaussagen bezüglich der optimalen Lösung und erlaubt schließlich eine kompakte Approximationsstruktur. Diese Struktur wird anschließend genutzt, um mithilfe einer modellfreien ADP-Methode das optimale Regelgesetz zu erlernen. Konvergenzaussagen der iterativen Off-policy-Methode werden auf den hier vorgestellten Solltrajektorienfolgeregelungsfall übertragen. Simulationsergebnisse, die den prädiktiven Charakter der vorgestellten Methode offenbaren, und eine anschließende Diskussion runden diesen Abschnitt ab.

3.3.1 Problemdefinition

Betrachtet werde ein lineares, zeitdiskretes System

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k \quad (3.80)$$

mit den diskreten Zeitschritten $k \in \mathbb{N}_{\geq 0}$, dem Zustandsvektor $\mathbf{x}_k \in \mathbb{R}^n$ und dem Eingangsvektor $\mathbf{u}_k \in \mathbb{R}^p$. Es wird angenommen, dass das System (\mathbf{A}, \mathbf{B}) steuerbar ist, jedoch seien $\mathbf{A} \in \mathbb{R}^{n \times n}$ und $\mathbf{B} \in \mathbb{R}^{n \times p}$ unbekannt.

In jedem Zeitschritt k sei der beliebige Sollzustandsverlauf $\mathbf{x}_{r,\text{soll},i}^{\{k\}} \in \mathbb{R}^n$ auf einem Vorausschauhorizont mit einer Länge von n_h Zeitschritten gegeben, wobei i den Zeitindex auf diesem Vorausschauhorizont bezeichnet. Jenseits des Horizonts n_h werde der Sollzustandsverlauf zu $\mathbf{0}$ angenommen. Somit ergibt sich aus Sicht des Zeitschritts k insgesamt der Sollzustandsverlauf

$$\mathbf{x}_{r,i}^{\{k\}} := \begin{cases} \mathbf{x}_{r,\text{soll},i}^{\{k\}}, & \text{für } i = k, \dots, k + n_h \\ \mathbf{0}, & \text{für } i > k + n_h. \end{cases} \quad (3.81)$$

Das Ziel ist im Folgenden, ohne Kenntnis von \mathbf{A} und \mathbf{B} ein Regelgesetz zu bestimmen, sodass der Systemzustand \mathbf{x}_i dem Solltrajektorienverlauf $\mathbf{x}_{r,i}^{\{k\}}$ (3.81) optimal bezüglich des Gütefunktional

$$\begin{aligned} J_k &= \sum_{i=k}^{\infty} \gamma^{i-k} \frac{1}{2} ((\mathbf{x}_i - \mathbf{x}_{r,i})^\top \mathbf{Q} (\mathbf{x}_i - \mathbf{x}_{r,i}) + \mathbf{u}_i^\top \mathbf{R} \mathbf{u}_i) \\ &=: \sum_{i=k}^{\infty} \gamma^{i-k} \underbrace{\frac{1}{2} (\mathbf{e}_i^\top \mathbf{Q} \mathbf{e}_i + \mathbf{u}_i^\top \mathbf{R} \mathbf{u}_i)}_{=: r(\mathbf{x}_i, \mathbf{u}_i, \mathbf{x}_{r,i}) =: r_i}, \end{aligned} \quad (3.82)$$

welches es zu minimieren gilt, folgt. Hierbei wird die Kurzschreibweise $\mathbf{x}_{r,i} := \mathbf{x}_{r,i}^{\{k\}}$ verwendet. Zudem bezeichnet $\mathbf{e}_i := \mathbf{x}_i - \mathbf{x}_{r,i}$ im Zeitschritt i die Abweichung des Systemzustands \mathbf{x}_i vom Sollzustand $\mathbf{x}_{r,i}$. Weiterhin sei $\mathbf{Q} \in \mathbb{R}^{n \times n}$, $\mathbf{Q} = \mathbf{Q}^\top \succeq \mathbf{0}$, eine symmetrische, positiv semidefinite Matrix, welche Abweichungen des Zustands \mathbf{x}_i vom Sollzustand $\mathbf{x}_{r,i}$

bestraft. Zudem sei (A, C) detektierbar, wobei C so definiert ist, dass $C^\top C = Q$ gilt. Der Stellaufwand wird durch die symmetrische, positiv definite Matrix $R \in \mathbb{R}^{m \times m}$, $R = R^\top \succ 0$, bestraft. Durch $r_i \in \mathbb{R}$ werden die Einschrittkosten bezeichnet und $\gamma \in (0, 1]$ stellt einen Diskontierungsfaktor dar⁶¹.

Die Wahl des Kostenfunktional (3.82) ist dabei wie folgt motiviert: Einerseits ist der Sollzustandsverlauf üblicherweise nur begrenzt in die Zukunft bekannt (beispielsweise der Straßenverlauf beim autonomen Fahren), daher wird ein endlicher Vorausschauhorizont n_h verwendet. Andererseits basiert die in Abschnitt 3.3.3 verwendete ADP-Methode auf einem Kostenfunktional mit unendlichem Zeithorizont, wodurch eine effiziente Darstellung der Bellman-Gleichung ermöglicht wird (vgl. Bemerkung 2.1). Daher wird das Kostenfunktional (3.82) auf einem unendlichen Optimierungshorizont definiert, berücksichtigt dabei jedoch einen auf einem endlichen Vorausschauhorizont gegebenen Sollzustandsverlauf. Die nachfolgende Proposition zeigt, dass für die gegebene Problemstellung ein eindeutiger optimaler linearer Regler existiert.

Annahme 3.3

Sei (A, B) steuerbar, $Q = Q^\top \succeq 0$ und $R = R^\top \succ 0$, sowie (A, C) detektierbar mit $C^\top C = Q$.

Proposition 3.3

Unter Annahme 3.3 existiert ein eindeutiger optimaler linearer Regler $\mu^*(\tilde{x}_k)$ mit $\tilde{x}_k := \tilde{x}_{k,k} := \left[\mathbf{x}_k^\top \quad \mathbf{x}_{r,k}^{\{k\}\top} \quad \dots \quad \mathbf{x}_{r,k+n_h}^{\{k\}\top} \right]^\top$, der das durch (3.80)–(3.82) beschriebene optimale Solltrajektorienfolge-
regelungsproblem löst.

Beweis:

Die zeitliche Propagation des Sollzustands $\mathbf{x}_{r,i}^{\{k\}}$ in (3.81) kann durch

$$\begin{bmatrix} \mathbf{x}_{r,i+1}^{\{k\}} \\ \vdots \\ \mathbf{x}_{r,i+n_h+1}^{\{k\}} \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{0}_{n n_h \times n} & \mathbf{I}_{n n_h} \\ \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n n_h} \end{bmatrix}}_{=: D} \underbrace{\begin{bmatrix} \mathbf{x}_{r,i}^{\{k\}} \\ \vdots \\ \mathbf{x}_{r,i+n_h}^{\{k\}} \end{bmatrix}}_{=: \mathbf{x}_{r,i+i+n_h}^{\{k\}}}, \quad i = k, k+1, \dots, \quad (3.83)$$

wobei $\mathbf{0}_{n_1 \times n_2}$ eine $n_1 \times n_2$ Nullmatrix und $\mathbf{I}_{n n_h}$ eine $n n_h \times n n_h$ Einheitsmatrix bezeichnet, ausgedrückt werden. Der erweiterte Systemzustand

$$\tilde{\mathbf{x}}_{k,i} := \left[\mathbf{x}_i^\top \quad \mathbf{x}_{r,i}^{\{k\}\top} \quad \dots \quad \mathbf{x}_{r,i+n_h}^{\{k\}\top} \right]^\top, \quad i = k, k+1, \dots, \quad (3.84)$$

⁶¹ Falls beschränkte Einschrittkosten r_i vorliegen und diese für $i \rightarrow \infty$ zudem verschwinden, ist das Kostenfunktional J_k selbst für $\gamma = 1$ endlich. Im Fall nicht-verschwindender aber beschränkter Einschrittkosten r_i sorgt $\gamma < 1$ dafür, dass J_k endlich ist (vgl. [SB18, Abschnitt 3.3]). Nicht-verschwindende Einschrittkosten r_i treten beispielsweise auf, wenn Stellenergie benötigt wird, um das System im Zustand $\mathbf{0}$ zu halten.

folgt somit der durch

$$\tilde{\mathbf{x}}_{k,i+1} = \begin{bmatrix} \mathbf{A}\mathbf{x}_i + \mathbf{B}\mathbf{u}_i \\ \mathbf{D}\mathbf{x}_{r,i:i+n_h}^{\{k\}} \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{A} & \mathbf{0}_{n \times n(n_h+1)} \\ \mathbf{0}_{n(n_h+1) \times n} & \mathbf{D} \end{bmatrix}}_{=:\tilde{\mathbf{A}}} \underbrace{\begin{bmatrix} \mathbf{x}_i \\ \mathbf{x}_{r,i:i+n_h}^{\{k\}} \end{bmatrix}}_{=:\tilde{\mathbf{x}}_{k,i}} + \underbrace{\begin{bmatrix} \mathbf{B} \\ \mathbf{0}_{n(n_h+1) \times p} \end{bmatrix}}_{=:\tilde{\mathbf{B}}} \mathbf{u}_i \quad (3.85)$$

gegebenen Dynamik. Zudem kann $\mathbf{e}_i^\top \mathbf{Q} \mathbf{e}_i$ in (3.82) durch

$$\mathbf{e}_i^\top \mathbf{Q} \mathbf{e}_i = \begin{bmatrix} \mathbf{x}_i \\ \mathbf{x}_{r,i}^{\{k\}} \\ \vdots \\ \mathbf{x}_{r,i+n_h}^{\{k\}} \end{bmatrix}^\top \underbrace{\begin{bmatrix} \mathbf{Q} & -\mathbf{Q} & \mathbf{0}_{n \times nn_h} \\ -\mathbf{Q} & \mathbf{Q} & \mathbf{0}_{n \times nn_h} \\ \mathbf{0}_{nn_h \times n} & \mathbf{0}_{nn_h \times n} & \mathbf{0}_{nn_h \times nn_h} \end{bmatrix}}_{=:\tilde{\mathbf{Q}}} \begin{bmatrix} \mathbf{x}_i \\ \mathbf{x}_{r,i}^{\{k\}} \\ \vdots \\ \mathbf{x}_{r,i+n_h}^{\{k\}} \end{bmatrix} = \tilde{\mathbf{x}}_{k,i}^\top \tilde{\mathbf{Q}} \tilde{\mathbf{x}}_{k,i} \quad (3.86)$$

ausgedrückt werden. Folglich lässt sich das durch (3.80)–(3.82) beschriebene Optimierungsproblem als diskontiertes Standard-LQ-Optimierungsproblem mit γ , $\tilde{\mathbf{A}}$, $\tilde{\mathbf{B}}$, $\tilde{\mathbf{Q}}$ und \mathbf{R} ausdrücken. Die gesuchte Aussage folgt analog zum Beweis von Satz 3.1. Hierbei wird genutzt, dass alle Eigenwerte von $\sqrt{\gamma} \tilde{\mathbf{D}}$ im Koordinatenursprung und somit innerhalb des Einheitskreises liegen. \square

Bemerkung 3.6

Eine mögliche Alternative zur Annahme $\mathbf{x}_{r,i}^{\{k\}} = \mathbf{0}, \forall i > k + n_h$, in (3.81) ist, anzunehmen, dass der Sollzustand jenseits des Vorausschauhorizonts der Länge n_h auf einen konstanten Wert $\neq \mathbf{0}$ gesetzt wird. Falls beispielsweise $\mathbf{x}_{r,i}^{\{k\}} = \mathbf{x}_{r,k+n_h}^{\{k\}}, \forall i > k + n_h$, gewählt wird, folgt für \mathbf{D} in (3.83)

$$\mathbf{D} := \begin{bmatrix} \mathbf{0}_{n(n_h-1) \times n} & \mathbf{I}_{n(n_h-1)} & \mathbf{0}_{n(n_h-1) \times n} \\ \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n(n_h-1)} & \mathbf{I}_n \\ \mathbf{0}_{n \times n} & \mathbf{0}_{n \times n(n_h-1)} & \mathbf{I}_n \end{bmatrix}. \quad (3.87)$$

Für $\gamma < 1$ bleibt die Aussage von Proposition 3.3 dann weiterhin bestehen. Die Lösungen dieser beiden Formulierungen nähern sich für kleiner werdende Diskontierungen γ und große Vorausschauhorizonte n_h einander an.

Bemerkung 3.7

Der Beweis von Proposition 3.3 offenbart, dass die Solltrajektorien­darstellung ADP-kompatibel gemäß Definition 3.1 ist. Hierbei ist

$$\zeta_k := \begin{bmatrix} \mathbf{x}_{r,k}^{\{k\}\top} & \dots & \mathbf{x}_{r,k+n_h}^{\{k\}\top} \end{bmatrix}^\top, \quad (3.88)$$

außerdem gilt

$$\mathbf{f}_{\mathbf{x}_r, \zeta}(\zeta_k) := \begin{bmatrix} \mathbf{I}_n & \mathbf{0}_{n \times n_h} \end{bmatrix} \zeta_k \quad (3.89)$$

(vgl. (3.6)) und analog zu (3.7)

$$\mathbf{f}_\zeta(\zeta_k) := \mathbf{D}\zeta_k. \quad (3.90)$$

Der optimale Regler, der nach Proposition 3.3 existiert, soll, wie in der folgenden Problemdefinition zusammengefasst, ohne Kenntnis der Systemmatrizen \mathbf{A} und \mathbf{B} ermittelt werden.

Problem 3.3

Sei $\mathbf{u}_k^*, \mathbf{u}_{k+1}^*, \dots$ die Stellgrößen­sequenz, die das Gütefunktional J_k in (3.82) für ein durch (3.80) beschriebenes System minimiert. Zudem seien die Systemmatrizen \mathbf{A} und \mathbf{B} unbekannt. Finde $\mathbf{u}_k^* = \boldsymbol{\mu}^*(\tilde{\mathbf{x}}_k)$, d. h. den optimalen Regler in Abhängigkeit von \mathbf{x}_k und $\mathbf{x}_{r,\text{soll},k}, \mathbf{x}_{r,\text{soll},k+1}, \dots, \mathbf{x}_{r,\text{soll},k+n_h}$.

3.3.2 Solltrajektorienabhängige Q-Function

Die zugrunde liegende Idee zur Lösung von Problem 3.3 ist die Definition einer solltrajektorienabhängigen Q-Function, d. h. einer Zustands-Aktions-Solltrajektorien-Nutzenfunktion im Gegensatz zur üblichen Zustands-Aktions-Nutzenfunktion. Diese Q-Function wird derart konstruiert, dass die sie minimierende Stellgröße eine Lösung für Problem 3.3 darstellt.

Im vorliegenden Abschnitt wird diese solltrajektorienabhängige Q-Function zunächst definiert, sowie deren analytische Lösung hergeleitet. Die analytische Lösung liefert wichtige Einblicke und beantwortet insbesondere die Frage einer geeigneten und möglichst effizienten Funktionsapproximation der Q-Function. Diese Funktionsapproximation wird in Abschnitt 3.3.3 für einen modellfreien ADP-Ansatz benötigt, der den optimalen Regler ohne Kenntnis der Systemmatrizen \mathbf{A} und \mathbf{B} lernt.

Zunächst werde der Fall eines endlichen Optimierungshorizonts der Länge K betrachtet, später erfolgt $K \rightarrow \infty$ ⁶². Die solltrajektorienabhängige Q-Function für einen Optimierungshorizont der Länge K sei wie nachfolgend gegeben definiert.

⁶² Bei der nachfolgenden Notation ist eine Unterscheidung des Optimierungshorizonts K sowie des Vorausschauhorizonts n_h des Solltrajektorienverlaufs wichtig. Zudem muss zwischen dem aktuellen Zeitschritt k , in dem sich das System befindet und von dem ausgehend das Optimierungsproblem gelöst werden soll, und der Zeit κ auf dem Optimierungshorizont der Länge K unterschieden werden.

Definition 3.6 (Solltrajektorienabhängige Q-Funktion)

Die solltrajektorienabhängige Q-Funktion für einen Optimierungshorizont der Länge K mit der Kurzschreibweise ${}^K Q_\kappa := Q_{K-\kappa}(\mathbf{x}_{k+\kappa}, \mathbf{u}_{k+\kappa}, \mathbf{x}_{r,k+\kappa}, \dots, \mathbf{x}_{r,k+K})$ sei durch die Bellman-Gleichung

$$\begin{aligned} {}^K Q_\kappa &:= r_{k+\kappa} + \gamma \sum_{i=k+\kappa+1}^{k+K} \gamma^{i-(k+\kappa+1)} r_i |_{\mathbf{u}_i^*} \\ &= r_{k+\kappa} + \gamma {}^K Q_{\kappa+1} |_{\mathbf{u}_{k+\kappa+1}^*} \end{aligned} \quad (3.91)$$

mit

$${}^K Q_K := r(\mathbf{x}_{k+K}, \mathbf{u}_{k+K}, \mathbf{x}_{r,k+K}) \quad (3.92)$$

und $r_i := r(\mathbf{x}_i, \mathbf{u}_i, \mathbf{x}_{r,i})$ wie in (3.82) definiert. Dabei bezeichnet $\kappa \in \mathbb{N}_{\geq 0}$, $\kappa < K$, die Zeit auf dem aktuellen Optimierungshorizont der Länge K , beginnend bei k . Zudem gelte $\mathbf{x}_{r,i} := \mathbf{x}_{r,i}^{\{k\}} = \mathbf{0}$, $\forall i > k + n_h$ (siehe (3.81)).⁶³

Folglich beschreibt ${}^K Q_\kappa$ die akkumulierten diskontierten Kosten vom Zeitschritt $k + \kappa$ bis zum Zeitschritt $k + K$, wenn die Stellgröße $\mathbf{u}_{k+\kappa}$ im Zeitschritt $k + \kappa$ und anschließend die optimalen Stellgrößen $\mathbf{u}_{k+\kappa+1}^*, \dots, \mathbf{u}_{k+K}^*$, welche die Kosten auf dem verbliebenen Horizont minimieren, verwendet werden. Dabei ist der Sollzustandsverlauf wie in (3.81) auf einem gleitenden Vorausschauhorizont der Länge n_h gegeben. Für das durch

$${}^K J_k = \sum_{i=k}^{k+K} \gamma^{i-k} \underbrace{\frac{1}{2} (\mathbf{e}_i^\top \mathbf{Q} \mathbf{e}_i + \mathbf{u}_i^\top \mathbf{R} \mathbf{u}_i)}_{=r_i} \quad (3.93)$$

definierte Kostenfunktional auf einem Optimierungshorizont der Länge K gilt das folgende Lemma.

Lemma 3.3

Die Stellgröße \mathbf{u}_k , welche die solltrajektorienabhängige Q-Funktion ${}^K Q_0$ minimiert, stellt eine optimale Lösung \mathbf{u}_k^* dar, die ${}^K J_k$ minimiert.

⁶³ Die Notation $r_i |_{\mathbf{u}_i^*}$ meint, dass die optimale Stellgröße \mathbf{u}_i^* in die Funktion $r(\cdot)$ nach (3.82) eingesetzt wird. Analog folgt ${}^K Q_{\kappa+1} |_{\mathbf{u}_{k+\kappa+1}^*}$ durch die Verwendung von $\mathbf{u}_{k+\kappa+1} = \mathbf{u}_{k+\kappa+1}^*$ in ${}^K Q_{\kappa+1}$.

Beweis:

Aus (3.93) und

$$\begin{aligned} \min_{\mathbf{u}_k} {}^K Q_0 &= r_k |_{\mathbf{u}_k^*} + \gamma {}^K Q_1 |_{\mathbf{u}_{k+1}^*} \\ &= {}^K Q_0 |_{\mathbf{u}_k^*} = \sum_{i=k}^{k+K} \gamma^{i-k} r_i |_{\mathbf{u}_i^*} \end{aligned} \quad (3.94)$$

folgt

$$\min_{\mathbf{u}_k, \dots, \mathbf{u}_{k+K}} {}^K J_k = \sum_{i=k}^{k+K} \gamma^{i-k} r_i |_{\mathbf{u}_i^*} = {}^K Q_0 |_{\mathbf{u}_k^*}. \quad (3.95)$$

□

Für den Grenzübergang

$$\lim_{K \rightarrow \infty} {}^K J_k = J_k \quad (3.96)$$

folgt J_k nach (3.82). Mit der Q-Funktion ${}^K Q_0$ nach Definition 3.6 und als Resultat von Lemma 3.3 lässt sich Problem 3.3 daher in die folgende äquivalente Problemstellung umformulieren.

Problem 3.4

Seien die Systemmatrizen \mathbf{A} und \mathbf{B} unbekannt. Im Zeitschritt k sei das System im Zustand \mathbf{x}_k , zudem sei der Sollzustand auf einem Horizont der Länge n_h durch $\mathbf{x}_{\tau,k}, \mathbf{x}_{\tau,k+1}, \dots, \mathbf{x}_{\tau,k+n_h}$ gegeben. Finde $\mathbf{u}_k^* = \boldsymbol{\mu}^*(\tilde{\mathbf{x}}_k)$, sodass die Q-Funktion

$$Q_0 := \lim_{K \rightarrow \infty} {}^K Q_0 \quad (3.97)$$

minimiert wird.

Zunächst werde angenommen, dass die Systemmatrizen \mathbf{A} und \mathbf{B} bekannt sind, um die analytische Lösung von ${}^K Q_0$, wie im nachfolgenden Satz 3.4 gezeigt, zu untersuchen. Im nächsten Schritt wird dann eine iterative ADP-Methode verwendet, um ohne Kenntnis von \mathbf{A} und \mathbf{B} , basierend auf dem TD-Fehler, die optimale Q-Funktion und das optimale Regelgesetz zu ermitteln.

Satz 3.4 (Analytische Lösung von ${}^K Q_0$)

Für $K \geq n_h$ ist die zum Gütefunktional ${}^K J_k$ in (3.93) gehörende Q-Function ${}^K Q_0$ (vgl. Definition 3.6) durch

$${}^K Q_0 = \frac{1}{2} \begin{bmatrix} \mathbf{x}_k^\top & \mathbf{u}_k^\top & \mathbf{x}_{r,k}^\top & \cdots & \mathbf{x}_{r,k+n_h}^\top & \mathbf{0}^\top \end{bmatrix} \mathbf{H}_K \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \\ \mathbf{x}_{r,k} \\ \vdots \\ \mathbf{x}_{r,k+n_h} \\ \mathbf{0} \end{bmatrix} \quad (3.98)$$

gegeben. Hierbei sei $\mathbf{H}_K = \mathbf{H}_K^\top \in \mathbb{R}^{((K+2)n+p) \times ((K+2)n+p)}$ mit

$$\mathbf{H}_K = \begin{bmatrix} h_{xx} & h_{xu} & h_{xxr,0} & h_{xxr,1} & h_{xxr,2} & \cdots & h_{xxr,K} \\ h_{ux} & h_{uu} & \mathbf{0} & h_{uxr,1} & h_{uxr,2} & \cdots & h_{uxr,K} \\ h_{x_r,0x} & \mathbf{0} & h_{x_r,0x_r,0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ h_{x_r,1x} & h_{x_r,1u} & \mathbf{0} & h_{x_r,1x_r,1} & \mathbf{0} & \cdots & \mathbf{0} \\ h_{x_r,2x} & h_{x_r,2u} & \mathbf{0} & \mathbf{0} & h_{x_r,2x_r,2} & \cdots & h_{x_r,2x_r,K} \\ h_{x_r,3x} & h_{x_r,3u} & \mathbf{0} & \mathbf{0} & h_{x_r,3x_r,2} & \cdots & h_{x_r,3x_r,K} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ h_{x_r,Kx} & h_{x_r,Ku} & \mathbf{0} & \mathbf{0} & h_{x_r,Kx_r,2} & \cdots & h_{x_r,Kx_r,K} \end{bmatrix}. \quad (3.99)$$

Die genauen Werte von \mathbf{H}_K folgen aus den Berechnungen im nachfolgenden Beweis.

Beweis:

Der Beweis ist in Anhang A.1 skizziert. □

Für $K \rightarrow \infty$ sei \mathbf{H} die nordwestliche $((n_h + 2)n + p) \times ((n_h + 2)n + p)$ -Teilmatrix von \mathbf{H}_K . Dann folgt wegen $\mathbf{x}_{r,i} = \mathbf{0}, \forall i > k + n_h$,

$$Q_0 = \lim_{K \rightarrow \infty} {}^K Q_0 = \frac{1}{2} \tilde{\mathbf{y}}_k^\top \mathbf{H} \tilde{\mathbf{y}}_k, \quad (3.100)$$

wobei $\tilde{\mathbf{y}}_k := [\mathbf{x}_k^\top \quad \mathbf{u}_k^\top \quad \mathbf{x}_{r,k}^\top \quad \cdots \quad \mathbf{x}_{r,k+n_h}^\top]^\top$.

Daher ist die Q-Function des Solltrajektorienfolgeregelungsproblems (vgl. Problem 3.3 bzw. Problem 3.4) quadratisch bezüglich des Systemzustands \mathbf{x}_k , der Stellgröße \mathbf{u}_k und des Sollzustandsverlaufs $\mathbf{x}_{r,k}, \dots, \mathbf{x}_{r,k+n_h}$ und zudem vollständig durch die Matrix \mathbf{H} parametrisiert. Somit stellt die hier präsentierte Q-Function Q_0 eine Verallgemeinerung der aus dem Regulationsfall [BYB94], [Lan97] bekannten Struktur der Q-Function dar. Weiterhin ist Q_0 nicht nur quadratisch, vielmehr weist \mathbf{H} eine spezielle Struktur auf. Die Relevanz der aus Satz 3.4 folgenden Struktur ist in der nachfolgenden Bemerkung zusammengefasst.

Bemerkung 3.8

Satz 3.4 liefert Aussagen über die genaue analytische Lösung der vorgestellten sollzustands-abhängigen Q -Function. Die damit einhergehende Erkenntnis über die exakte Struktur der Q -Function ermöglicht eine effiziente Wahl eines geeigneten Funktionsapproximators und erlaubt somit eine signifikante Reduktion der benötigten zu lernenden Gewichte (vgl. Lemma 3.4 und Lemma 3.5 sowie Abschnitt 3.3.5) bei der späteren ADP-Umsetzung.

3.3.3 Modellfreies Erlernen der Q -Function

Im Folgenden wird die Q -Function, die explizit vom Sollzustandsverlauf auf einem gleitenden Vorausschauhorizont der Länge n_h abhängt, ohne Kenntnis der Systemmatrizen \mathbf{A} und \mathbf{B} gelernt. Konkret wird dabei \mathbf{H} (vgl. (3.100)) aus Zustandsübergängen und damit verknüpften Einschrittkosten ermittelt.

Die optimale Stellgröße $\mathbf{u}_k^* = \boldsymbol{\mu}^*(\tilde{\mathbf{x}}_k)$, die für $\kappa = 0$ und $K \rightarrow \infty$ äquivalent zu (A.17) ist, kann dann, wie im nachfolgenden Korollar gegeben, direkt aus \mathbf{H} bestimmt werden. Dabei stellt (A.18) sicher, dass \mathbf{u}_k^* tatsächlich die Q -Function Q_0 minimiert.

Korollar 3.1 (Optimaler Regler in Abhängigkeit von \mathbf{H})

Aus Lemma 3.3 und Satz 3.4 folgt direkt, dass die nach Problem 3.4 optimale Stellgröße im Zeitschritt k durch

$$\mathbf{u}_k^* = \boldsymbol{\mu}^*(\tilde{\mathbf{x}}_k) = -\mathbf{h}_{uu}^{-1} \begin{bmatrix} \mathbf{h}_{ux} & \mathbf{h}_{ux,r,1} & \cdots & \mathbf{h}_{ux,r,n_h} \end{bmatrix} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{x}_{r,k+1} \\ \vdots \\ \mathbf{x}_{r,k+n_h} \end{bmatrix} \quad (3.101)$$

$$=: - \underbrace{\begin{bmatrix} \mathbf{K}_x^* & \mathbf{K}_{\text{ref}}^* \end{bmatrix}}_{=: \mathbf{K}^*} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{x}_{r,k+1} \\ \vdots \\ \mathbf{x}_{r,k+n_h} \end{bmatrix}$$

gegeben ist.

Im Gegensatz zu bestehenden ADP-Reglern (vgl. Abschnitt 2.2) bezieht das optimale Regelgesetz (3.101) den exakten Sollzustandsverlauf $\mathbf{x}_{r,k+1}, \dots, \mathbf{x}_{r,k+n_h}$ auf einem Vorausschauhorizont der Länge n_h explizit ein und ist daher für flexible Solltrajektorienverläufe geeignet.

Bemerkung 3.9

Nach [Kuč72, Theorem 8] folgt, dass das optimale Regelgesetz $\mu^*(\tilde{\mathbf{x}}_k)$ aus Proposition 3.3 das fiktive System

$$\tilde{\mathbf{x}}_{k+1} = \sqrt{\gamma} \tilde{\mathbf{A}} \tilde{\mathbf{x}}_k + \sqrt{\gamma} \tilde{\mathbf{B}} \mu^*(\tilde{\mathbf{x}}_k) \quad (3.102)$$

stabilisiert. Im Fall $\gamma = 1$ folgt daraus direkt, dass das geregelte System

$$\mathbf{x}_{k+1} = (\mathbf{A} - \mathbf{B}\mathbf{K}_x^*) \mathbf{x}_k - \mathbf{B}\mathbf{K}_{\text{ref}}^* \begin{bmatrix} \mathbf{x}_{r,k+1} \\ \vdots \\ \mathbf{x}_{r,k+n_h} \end{bmatrix} \quad (3.103)$$

für beschränkte Sollzustandsvorgaben $\mathbf{x}_{r,k+1}, \dots, \mathbf{x}_{r,k+n_h}$ Referenz-Zustands-stabil nach Definition 3.5 ist. Für $\gamma < 1$ muss, analog zu Abschnitt 3.2.4, überprüft werden, ob sämtliche Eigenwerte von $\mathbf{A} - \mathbf{B}\mathbf{K}_x^*$ im Inneren des Einheitskreises liegen. Falls \mathbf{A} und \mathbf{B} nicht explizit bekannt sind, kann hierzu anhand der Matrix \mathbf{H} aus (3.99) das durch Satz 3.3 gegebene hinreichende Kriterium verwendet werden.

3.3.3.1 Funktionsapproximation der erweiterten Q-Function

Um die Q-Function datenbasiert und ohne explizite Kenntnis der Systemmatrizen \mathbf{A} und \mathbf{B} zu lernen, wird Q_0 durch einen Funktionsapproximator (vgl. Abschnitt 2.1.3) parametrisiert. Anschließend wird der quadratische TD-Fehler genutzt, um eine Value Iteration (vgl. Abschnitt 2.1.4.2) durchzuführen, deren Resultat die geschätzten optimalen Q-Function-Gewichte sowie das zugehörige optimale Regelgesetz darstellt. Hierzu sei die geschätzte Q-Function \hat{Q}_0 durch eine Summe gewichteter Basisfunktionen gegeben:

$$\hat{Q}_0 = \hat{\mathbf{w}}^\top \phi(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{r,k}, \dots, \mathbf{x}_{r,k+n_h}) = \hat{\mathbf{w}}^\top \phi(\tilde{\mathbf{y}}_k). \quad (3.104)$$

Dabei stellt $\hat{\mathbf{w}} \in \mathbb{R}^h$ den zu schätzenden Gewichtsvektor dar und $\phi : \mathbb{R}^{(n_h+2)n+p} \rightarrow \mathbb{R}^h$ ist ein Basisfunktionsvektor. Im Gegensatz zu klassischen Q-Functions bezieht \hat{Q}_0 in (3.104) somit explizit den durch $\mathbf{x}_{r,k}, \mathbf{x}_{r,k+1}, \dots, \mathbf{x}_{r,k+n_h}$ beschriebenen Solltrajektorienverlauf ein. Das nachfolgende Lemma besagt, dass die Q-Function Q_0 mithilfe quadratischer Basisfunktionen exakt durch \hat{Q}_0 repräsentiert werden kann.

Lemma 3.4

Mit

$$h = \frac{1}{2}((n_h + 2)n + p)((n_h + 2)n + p + 1) - (n^2(2n_h - 1) + pn) \quad (3.105)$$

quadratischen Basisfunktionen $\phi = [\phi_1 \ \dots \ \phi_h]^\top$ existiert ein optimaler Gewichtsvektor $\hat{\mathbf{w}} = \mathbf{w}^*$, sodass $\hat{Q}_0 = Q_0$ gilt.

Beweis:

Aufgrund der Symmetrie von \mathbf{H}_K in (3.99) und der darin auftauchenden Nullen weist \mathbf{H} maximal h nicht-redundante Elemente auf, wobei h durch (3.105) gegeben ist. Für $\bar{h} = 1, \dots, h$ werden quadratische Basisfunktionen, deren Elemente die Form

$$\phi_{\bar{h}} = \begin{cases} \tilde{y}_{k,i} \tilde{y}_{k,j}, & \text{für } i \neq j \\ \frac{1}{2} \tilde{y}_{k,i}^2, & \text{für } i = j, \end{cases} \quad (3.106)$$

aufweisen, definiert. Hierbei indiziere i, j die entsprechenden nicht-redundanten Elemente von \mathbf{H} , \tilde{y}_k sei wie in (3.100) definiert und $\tilde{y}_{k,i}$ bezeichne das i -te Element des Vektors \tilde{y}_k . Folglich ist Q_0 in (3.100) gerade dann äquivalent zu \hat{Q}_0 in (3.104), wenn die Gewichte \hat{w}_i , $i = 1, \dots, h$, identisch zu den entsprechenden nicht-redundanten Elementen von \mathbf{H} sind, d. h. wenn $\hat{\mathbf{w}} = \mathbf{w}^*$ gilt. \square

Somit gibt h in Lemma 3.4 die maximal benötigte Anzahl an Basisfunktionen und zugehörigen zu lernenden Gewichten an, um Q_0 exakt parametrieren zu können. In der Anwendung liegen jedoch häufig dünnbesetzte Gewichtsmatrizen \mathbf{Q} im Gütefunktional (3.82) vor. Daher kann die Anzahl der benötigten Basisfunktionen und somit die Dimension von $\hat{\mathbf{w}}$ weiter reduziert werden. Hierfür wird die aus Satz 3.4 und dem zugehörigen Beweis resultierende Struktur von \mathbf{H} , die auch ohne explizite Kenntnis der Zahlenwerte in \mathbf{A} und \mathbf{B} exakt bekannt ist, genutzt. Dies wird im folgenden Lemma näher ausgeführt.

Lemma 3.5

Falls die o -te Zeile und o -te Spalte von \mathbf{Q} null ist, dann sind die

- o -te Spalte von $\mathbf{h}_{\mathbf{x}\mathbf{x}_r,i} \forall i \in \{0, \dots, n_h\}$,
- o -te Spalte von $\mathbf{h}_{\mathbf{u}\mathbf{x}_r,i} \forall i \in \{0, \dots, n_h\}$ und
- o -te Zeile und o -te Spalte von $\mathbf{h}_{\mathbf{x}_r,i \mathbf{x}_r,j} \forall i, j \in \{0, \dots, n_h\}$

ebenfalls null. Folglich reduziert sich die Anzahl h der nicht-redundanten Gewichte in \mathbf{H} , und somit die Dimension des zu schätzenden Gewichtsvektors $\hat{\mathbf{w}}$, zu

$$h = (n - q)(n - q + 1) \left(\frac{n_h}{2} + 1 \right) + \frac{1}{2} n(n + 1) \quad (3.107)$$

$$+ (p + n_h(n - q))(p + n) + \frac{1}{2} (n_h - 2)(n_h - 1)(n - q)^2.$$

Hierbei wurde $\mathbf{h}_{\mathbf{x}\mathbf{x}_r,0} = \mathbf{Q}$ und $\mathbf{h}_{\mathbf{x}_r,0 \mathbf{x}_r,0} = -\mathbf{Q}$ berücksichtigt, weiterhin bezeichnet q die Anzahl der Zeilen bzw. Spalten von \mathbf{Q} , die null sind.

Beweis:

Die Aussage folgt direkt aus Betrachtung von (A.15) unter Berücksichtigung von (A.11)–(A.14). \square

Obwohl Lemma 3.5 sehr technischer Natur ist, ist die Kenntnis über die Dünnbesetztheit der Matrix \mathbf{H} essenziell für die Implementierung effizienter ADP-Regler mit einer überschaubaren Anzahl h zu schätzender Elemente von $\hat{\mathbf{w}}$. Dank der genauen Kenntnis der spezifischen Struktur der analytischen Lösung von Q_0 , wie sie in Satz 3.4 und dessen Beweis hergeleitet wurde, genügen somit nach Lemma 3.5 h quadratische Basisfunktionen mit h nach (3.107), um Q_0 exakt zu parametrieren (vgl. Beispiel 3.1).

Beispiel 3.1 (Reduzierung der Dimension h des Critic-Gewichts $\hat{\mathbf{w}}$)

Tabelle 3.2 illustriert die Reduzierung der Dimension h des Critic-Gewichts $\hat{\mathbf{w}}$ durch Lemma 3.4 bzw. Lemma 3.5 im Vergleich dazu, nur die Symmetrie der Matrix \mathbf{H} zu berücksichtigen, für ein System zweiter Ordnung und ein System sechster Ordnung.

		System 1 ($n = 2, p = 1$)		System 2 ($n = 6, p = 1$)	
	q	1	1	5	5
	n_h	10	100	10	100
\mathbf{H} symmetrisch	h	325	21115	2701	188191
Lemma 3.4	h	247	20317	2011	181021
Lemma 3.5	h	84	5259	146	5681

Tabelle 3.2: Die Berücksichtigung der Struktur der Matrix \mathbf{H} erlaubt bei den betrachteten Beispielsystemen eine signifikante Verringerung der Dimension h des zu schätzenden Gewichtsvektors $\hat{\mathbf{w}}$.

Bemerkung 3.10

Die zu lernenden Gewichte $\hat{\mathbf{w}}$ hängen nicht von der Solltrajektorie ab. Somit verdeutlichen (3.104) und Lemma 3.5 auch, dass die vorgestellte Q -Function, die explizit vom Solltrajektorienverlauf abhängt, über beliebige Solltrajektorien, die auf einem Vorausschauhorizont der Länge n_h gegeben sind, generalisiert.

3.3.3.2 Datenbasierter Lernalgorithmus der erweiterten Q -Function

In diesem Abschnitt wird der auf einer Value Iteration basierende ADP-Algorithmus vorgestellt, der die Gewichte $\hat{\mathbf{w}}$ basierend auf der Minimierung des quadrierten TD-Fehlers und ohne Kenntnis der Systemmatrizen \mathbf{A} und \mathbf{B} adaptiert. Hierzu sei zunächst

$$\begin{aligned}\hat{Q}_1 &:= \hat{\mathbf{w}}^\top \phi(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}, \mathbf{x}_{r,k+1}, \dots, \mathbf{x}_{r,k+n_h}, \mathbf{0}) \\ &= \hat{\mathbf{w}}^\top \phi(\tilde{\mathbf{y}}_{k+1})\end{aligned}\quad (3.108)$$

mit $\tilde{\mathbf{y}}_{k+1} = [\mathbf{x}_{k+1}^\top \quad \mathbf{u}_{k+1}^\top \quad \mathbf{x}_{r,k+1}^\top \quad \cdots \quad \mathbf{x}_{r,k+n_h}^\top \quad \mathbf{0}^\top]^\top$ definiert⁶⁴. Die geschätzte optimale Stellgröße $\hat{\mathbf{u}}_{k+\kappa}^*$, die auf der geschätzten Q-Function $\hat{Q}_\kappa \big|_{\mathbf{u}_{k+\kappa}}$ basiert, werde zudem durch

$$\begin{aligned} \hat{\mathbf{u}}_{k+\kappa}^* &:= \arg \min_{\mathbf{u}_{k+\kappa}} \hat{Q}_\kappa \big|_{\mathbf{u}_{k+\kappa}} \\ &= -\mathbf{K} [\mathbf{x}_{k+\kappa}^\top \quad \mathbf{x}_{r,k+\kappa+1}^\top \quad \cdots \quad \mathbf{x}_{r,k+\kappa+n_h}^\top]^\top \end{aligned} \quad (3.109)$$

mit $\mathbf{K} := \hat{\mathbf{h}}_{\text{uu}}^{-1} [\hat{\mathbf{h}}_{\text{ux}} \quad \hat{\mathbf{h}}_{\text{ux},1} \quad \cdots \quad \hat{\mathbf{h}}_{\text{ux},n_h}]$ (vgl. (3.101)) definiert. Dabei sind durch $\hat{\mathbf{h}}_{(\cdot)}$ analog zu (3.99) die Teilmatrizen von $\hat{\mathbf{H}}$ gegeben, wobei $\hat{\mathbf{H}}$ die auf $\hat{\mathbf{w}}$ basierende Schätzung der Matrix \mathbf{H} bezeichnet. Somit resultiert $\hat{\mathbf{u}}_{k+\kappa}^* = \mathbf{u}_{k+\kappa}^*$, falls $\hat{\mathbf{w}} = \mathbf{w}^*$ gilt und die optimale Reglermatrix $\mathbf{K} = \mathbf{K}^*$ ergibt sich, wenn $\hat{\mathbf{H}} = \mathbf{H}$ gilt. Basierend auf der Bellman-Gleichung (3.91) wird nachfolgend der TD-Fehler δ_k (vgl. [Sut88]) definiert, der für $\hat{\mathbf{w}} = \mathbf{w}^*$ verschwindet.

Definition 3.7 (TD-Fehler der solltrajektorienabhängigen Q-Function)

Der TD-Fehler δ_k , d. h. der Approximationsfehler der Bellman-Gleichung (3.91), der aufgrund einer Abweichung des geschätzten Gewichts $\hat{\mathbf{w}}$ von \mathbf{w}^* resultiert, sei durch

$$\begin{aligned} \delta_k &:= r_k + \gamma \hat{Q}_1 \big|_{\hat{\mathbf{u}}_{k+1}^*} - \hat{Q}_0 \\ &= r_k + \gamma \hat{\mathbf{w}}^\top \phi(\tilde{\mathbf{y}}_{k+1}^*) - \hat{\mathbf{w}}^\top \phi(\tilde{\mathbf{y}}_k) \end{aligned} \quad (3.110)$$

definiert. Dabei sei $\tilde{\mathbf{y}}_{k+1}^* := \tilde{\mathbf{y}}_{k+1} \big|_{\mathbf{u}_{k+1} = \hat{\mathbf{u}}_{k+1}^*}$.

Um die Schätzung der sollzustandsabhängigen Q-Function \hat{Q}_0 sowie die daraus resultierende geschätzte optimale Stellgröße $\hat{\mathbf{u}}_k^*$ zu verbessern, wird im Folgenden eine Value Iteration (vgl. Abschnitt 2.1.4.2) verwendet. Im Policy-Evaluation-Schritt wird die Gewichtsschätzung $\hat{\mathbf{w}}^{[l]}$, die in der Iteration l die geschätzte Q-Function repräsentiert, angepasst. Im nachfolgenden Policy-Improvement-Schritt wird, ausgehend von der neuen Schätzung $\hat{\mathbf{w}}^{[l]}$ des Q-Function-Gewichts und der zugehörigen Matrix $\mathbf{H}^{[l]}$, das Regelgesetz $\mathbf{K}^{[l]}$, wie in (3.109) gezeigt, adaptiert. Zur Berechnung von δ_k in (3.110) wird $\hat{\mathbf{u}}_{k+1}^*$ benötigt, jedoch ist das optimale Gewicht \mathbf{w}^* a priori unbekannt. Daher wird zunächst $\hat{\mathbf{w}}^{[0]} = \mathbf{0}$ und $\hat{\mathbf{u}}_{k+1}^{*[0]} = \mathbf{0}$ initialisiert. Letzteres geschieht durch die Wahl von $\mathbf{K}^{[0]} = \mathbf{0}$ ⁶⁵. Im Policy-Evaluation-Schritt ist das Ziel,

⁶⁴ Für die in Bemerkung 3.6 beschriebene alternative Problemformulierung, welche $\mathbf{x}_{r,i}^{\{k\}} = \mathbf{x}_{r,k+n_h}^{\{k\}}$, $\forall i > k + n_h$, annimmt, müsste an dieser Stelle $\phi(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}, \mathbf{x}_{r,k+1}, \dots, \mathbf{x}_{r,k+n_h}, \mathbf{x}_{r,k+n_h})$ verwendet und $\tilde{\mathbf{y}}_{k+1} = [\mathbf{x}_{k+1}^\top \quad \mathbf{u}_{k+1}^\top \quad \mathbf{x}_{r,k+1}^\top \quad \cdots \quad \mathbf{x}_{r,k+n_h}^\top \quad \mathbf{x}_{r,k+n_h}^\top]^\top$ gesetzt werden.

⁶⁵ Bei vorhandenem Vorwissen kann die Initialisierung zugunsten einer schnelleren Konvergenz entsprechend angepasst werden.

eine neue Schätzung $\hat{\boldsymbol{w}}^{[l+1]}$ zu erlangen, sodass $\left(\delta_k^{[l]}\right)^2$ minimiert wird, wobei in Analogie zu (3.110)

$$\delta_k^{[l]} := r_k + \gamma \hat{\boldsymbol{w}}^{[l]\top} \boldsymbol{\phi} \left(\tilde{\boldsymbol{y}}_{k+1}^{*[l]} \right) - \hat{\boldsymbol{w}}^{[l+1]\top} \boldsymbol{\phi} \left(\tilde{\boldsymbol{y}}_k \right) \quad (3.111)$$

gesetzt wird. In (3.111) sei $\tilde{\boldsymbol{y}}_{k+1}^{*[l]} := \tilde{\boldsymbol{y}}_{k+1} \big|_{\boldsymbol{u}_{k+1} = \hat{\boldsymbol{u}}_{k+1}^{*[l]}}$, wobei $\hat{\boldsymbol{u}}_{k+1}^{*[l]}$ aus (3.109) mit $\boldsymbol{K} = \boldsymbol{K}^{[l]}$ folgt. Mit $\hat{\boldsymbol{w}}^{[l]} \in \mathbb{R}^h$ wie in Lemma 3.5 ergibt sich, dass $\delta_k^{[l]}$ mindestens zu $M \geq h$ unterschiedlichen Zeitschritten betrachtet werden muss, um eine Least-Squares-Schätzung unter Nutzung von M Datentupeln durchführen zu können. Damit resultiert

$$\hat{\boldsymbol{w}}^{[l+1]} = \arg \min_{\hat{\boldsymbol{w}}^{[l+1]}} \sum_{j=k-M+1}^k \left(\delta_j^{[l]} \right)^2. \quad (3.112)$$

Es seien weiterhin

$$\boldsymbol{\Phi} := \left[\boldsymbol{\phi} \left(\tilde{\boldsymbol{y}}_{k-M+1} \right) \quad \dots \quad \boldsymbol{\phi} \left(\tilde{\boldsymbol{y}}_k \right) \right]^\top \quad (3.113)$$

und

$$\boldsymbol{r} := \begin{bmatrix} r_{k-M+1} + \gamma \hat{\boldsymbol{w}}^{[l]\top} \boldsymbol{\phi} \left(\tilde{\boldsymbol{y}}_{k-M+2}^{*[l]} \right) \\ \vdots \\ r_k + \gamma \hat{\boldsymbol{w}}^{[l]\top} \boldsymbol{\phi} \left(\tilde{\boldsymbol{y}}_{k+1}^{*[l]} \right) \end{bmatrix} \quad (3.114)$$

definiert. Falls die Anregungsbedingung

$$\text{Rang} \left(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} \right) = h \quad (3.115)$$

erfüllt ist, existiert eine eindeutige Lösung $\hat{\boldsymbol{w}}^{[l+1]}$, die (3.112) minimiert. Diese ist durch

$$\hat{\boldsymbol{w}}^{[l+1]} = \left(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} \right)^{-1} \boldsymbol{\Phi}^\top \boldsymbol{r} \quad (3.116)$$

gegeben (vgl. [ÄW95, Theorem 2.1]).

Der Policy-Improvement-Schritt ergibt sich dann aus der neuen Gewichtsschätzung $\hat{\boldsymbol{w}}^{[l+1]}$ und der zugehörigen Matrix $\boldsymbol{H}^{[l+1]}$ zu

$$\boldsymbol{K}^{[l+1]} = \boldsymbol{h}_{\text{uu}}^{[l+1]-1} \begin{bmatrix} \boldsymbol{h}_{\text{ux}}^{[l+1]} & \boldsymbol{h}_{\text{ux}_r,1}^{[l+1]} & \dots & \boldsymbol{h}_{\text{ux}_r,n_h}^{[l+1]} \end{bmatrix} \quad (3.117)$$

(vgl. Korollar 3.1). Diese Iteration wird, bei festem Zeitschritt k , so oft wiederholt, bis sich $\hat{\boldsymbol{w}}^{[l]}$ nicht mehr signifikant ändert, d. h. $\left\| \hat{\boldsymbol{w}}^{[l]} - \hat{\boldsymbol{w}}^{[l-1]} \right\|_2$ einen Schwellwert $e_{\hat{\boldsymbol{w}}}$ unterschreitet.

Wenngleich die Q-Function zu der durch $\boldsymbol{K}^{[l]}$ gegebenen Target Policy $\hat{\boldsymbol{\mu}}^{*[l]}(\cdot)$, die in Form von $\hat{\boldsymbol{u}}_{k+1}^{*[l]}$ in $\tilde{\boldsymbol{y}}_{k+1}^{*[l]}$ nach (3.111) berücksichtigt wird, ausgewertet werden soll, wurde bislang noch nicht diskutiert, welche Behavior Policy verwendet wird (vgl. Abschnitt 2.1.4.4). Die Behavior Policy stellt dabei die Stellgröße \boldsymbol{u}_k dar, die auf das System angewandt wird und welche

in den Ausdrücken r_k und $\tilde{\mathbf{y}}_k$ in (3.111) auftaucht. Während die Target Policy die aktuelle Schätzung des optimalen Regelgesetzes darstellt, kann die Behavior Policy genutzt werden, um die in (3.115) gegebene Anregungsbedingung zu erfüllen. Während der Adaptionsprozess aktiv ist, sei

$$\mathbf{u}_k = \tilde{\mathbf{u}}_k^{*[l]} := \hat{\mathbf{u}}_k^{*[l]} + \boldsymbol{\xi}_k. \quad (3.118)$$

Dabei dient $\boldsymbol{\xi}_k \sim \mathcal{N}(\mathbf{0}; \mathbf{V})$ als Anregungsrauschen, welches den Systemzustand \mathbf{x}_k anregt. Zur Anregung des Sollzustandsverlaufs kann dieser ebenfalls durch additives Gaußsches Rauschen ξ_{ref} überlagert werden, um bei einem während des Trainings möglicherweise abwechslungsarmen Sollzustandsverlauf dennoch ausreichende Anregung sicherzustellen.

Der gesamte ADP-Algorithmus, der ohne Kenntnis der Systemmatrizen \mathbf{A} und \mathbf{B} lernt, einer auf einem gleitenden Vorausschauhorizont der Länge n_h gegebenen, beliebigen Solltrajektorie optimal im Sinne des Gütefunktional (3.82) zu folgen, ist in Algorithmus 3.2 gegeben.

Algorithmus 3.2 ADP-Algorithmus mit Sollverlauf auf endlichem Vorausschauhorizont

```

1: Initialisiere  $M$ ,  $\hat{\mathbf{w}} := \hat{\mathbf{w}}^{[0]} := \mathbf{0}$ ,  $\mathbf{K}^{[0]} := \mathbf{0}$ 
2: for  $k = 0, 1, 2, \dots$  do
3:   wende die Stellgröße  $\tilde{\mathbf{u}}_k^{*[l]}$  nach (3.118) auf das System (3.80) an
4:   if  $k + 1 \bmod M = 0$  then
5:      $l := 0$ ,  $\hat{\mathbf{w}}^{[l]} := \hat{\mathbf{w}}$ 
6:     do
7:       Policy-Evaluation-Schritt: berechne  $\hat{\mathbf{w}}^{[l+1]}$  mit (3.112)
8:       Policy-Improvement-Schritt: berechne  $\mathbf{K}^{[l+1]}$  mit (3.117)
9:        $l := l + 1$ 
10:    while  $\left\| \hat{\mathbf{w}}^{[l]} - \hat{\mathbf{w}}^{[l-1]} \right\|_2 > e_{\hat{\mathbf{w}}}$ 
11:       $\hat{\mathbf{w}} := \hat{\mathbf{w}}^{[l]}$ 
12:    end if
13: end for

```

Bemerkung 3.11

Die in Algorithmus 3.2 verwendete iterative Methode aus Policy-Evaluation- und Policy-Improvement-Schritt stellt eine Value Iteration dar. Dies manifestiert sich in der Definition von δ_k in (3.111), wobei \hat{Q}_1 auf $\hat{\mathbf{w}}^{[l]}$ basiert und \hat{Q}_0 auf $\hat{\mathbf{w}}^{[l+1]}$ (vgl. (2.25)).

Weiterhin gehört der Algorithmus zur Klasse der Off-Policy-Methoden (vgl. Abschnitt 2.1.4.4), da die Behavior Policy $\tilde{\mathbf{u}}_k^{*[l]} = \hat{\mathbf{u}}_k^{*[l]} + \boldsymbol{\xi}_k$ angewandt wird, während die zur Target Policy $\hat{\mathbf{u}}_k^{*[l]}$ gehörende Q -Function gelernt wird.

3.3.3.3 Konvergenzanalyse des Lernalgorithmus

In diesem Abschnitt wird die Konvergenz der geschätzten Solltrajektorienabhängigen Q-Function \hat{Q}_0 gegen die optimale Q-Function Q_0 untersucht. Konkret wird gezeigt, dass für $l \rightarrow \infty$ die Konvergenz des Gewichts $\hat{w}^{[l]} \rightarrow w^*$ und damit auch der Matrix $\mathbf{H}^{[l]} \rightarrow \mathbf{H}$ folgt. Somit konvergiert die Value Iteration für die erweiterte Sollzustandsabhängige Q-Function gegen das optimale Regelgesetz, d. h. $\mathbf{K}^{[l]} \rightarrow \mathbf{K}^*$ folgt.

Die Konvergenzanalyse ist dabei wie folgt strukturiert: Zuerst wird gezeigt, dass die Value Iteration, d. h. die Iteration zwischen (3.112) und (3.117), äquivalent zu einer durch $\mathbf{H}^{[l]}$ gegebenen Matrixfolge ist. Im zweiten Schritt wird bewiesen, dass diese Matrixfolge im Sinne von $\mathbf{0} \preceq \mathbf{H}^{[l]} \preceq \mathbf{Y}$ nach oben beschränkt ist und zudem $\mathbf{0} \preceq \mathbf{H}^{[l]} \preceq \mathbf{H}^{[l+1]}$ gilt. Daraus folgt die Konvergenz der Matrixfolge. Im letzten Schritt wird gezeigt, dass die konvergierte Folge die Bellman-Gleichung erfüllt und das zugehörige Regelgesetz optimal ist.

Das nachfolgende Lemma stellt zunächst eine Erweiterung von [ATLAK07, Lemma 1] auf den Solltrajektorienregelungsfall dar und zeigt, dass die verwendete Value Iteration äquivalent zu einer Matrixfolge $\mathbf{H}^{[l]}$ ist.

Lemma 3.6

Sei $\mathbf{H}^{[0]} = \mathbf{0}$, $\mathbf{R}^\top = \mathbf{R} \succ \mathbf{0}$, $\mathbf{Q}^\top = \mathbf{Q} \succeq \mathbf{0}$ und (\mathbf{A}, \mathbf{B}) steuerbar. Die durch (3.112) und (3.117) beschriebene Value Iteration ist äquivalent zur Iteration

$$\mathbf{H}^{[l+1]} = \mathbf{G} + \gamma \mathbf{M} \left(-\mathbf{K}^{[l]} \right)^\top \mathbf{H}^{[l]} \mathbf{M} \left(-\mathbf{K}^{[l]} \right) \quad (3.119)$$

mit

$$\mathbf{G} := \begin{bmatrix} \mathbf{Q} & \mathbf{0} & -\mathbf{Q} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} & \mathbf{0} & \mathbf{0} \\ -\mathbf{Q} & \mathbf{0} & \mathbf{Q} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (3.120)$$

und

$$\mathbf{M} \left(-\mathbf{K}^{[l]} \right) := \begin{bmatrix} \mathbf{A} & \mathbf{B} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \\ -\mathbf{K}_x^{[l]} \mathbf{A} & -\mathbf{K}_x^{[l]} \mathbf{B} & \mathbf{0} & \mathbf{0} & -\mathbf{K}_1^{[l]} & \dots & -\mathbf{K}_{n_h-1}^{[l]} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I}_n & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I}_n & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{I}_n \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \end{bmatrix}, \quad (3.121)$$

wobei $\mathbf{K}^{[l]} = \begin{bmatrix} \mathbf{K}_x^{[l]} & \mathbf{K}_1^{[l]} & \dots & \mathbf{K}_{n_h}^{[l]} \end{bmatrix} =: \mathbf{K} \left(\mathbf{H}^{[l]} \right)$ analog zu (3.117) aus $\mathbf{H}^{[l]}$ folgt.

Beweis:

Mit quadratischen Basisfunktionen $\phi(\cdot)$ (vgl. Lemma 3.4 und Lemma 3.5) ergibt sich

$$\hat{\mathbf{w}}^{[l]\top} \phi(\tilde{\mathbf{y}}_{k+1}^{*[l]}) = \frac{1}{2} \tilde{\mathbf{y}}_{k+1}^{*[l]\top} \mathbf{H}^{[l]} \tilde{\mathbf{y}}_{k+1}^{*[l]}. \quad (3.122)$$

Daraus folgt aufgrund von $\tilde{\mathbf{y}}_{k+1}^{*[l]} = \mathbf{M}(-\mathbf{K}^{[l]}) \tilde{\mathbf{y}}_k$

$$r_k + \gamma \hat{\mathbf{w}}^{[l]\top} \phi(\tilde{\mathbf{y}}_{k+1}^{*[l]}) = \frac{1}{2} \tilde{\mathbf{y}}_k^\top \underbrace{\left(\mathbf{G} + \gamma \mathbf{M}(-\mathbf{K}^{[l]})^\top \mathbf{H}^{[l]} \mathbf{M}(-\mathbf{K}^{[l]}) \right)}_{=\mathbf{H}^{[l+1]}} \tilde{\mathbf{y}}_k. \quad (3.123)$$

Aus (3.123), (3.113) und (3.114) folgt, dass (3.116) äquivalent zu

$$\begin{aligned} \hat{\mathbf{w}}^{[l+1]} &= (\Phi^\top \Phi)^{-1} \Phi^\top \begin{bmatrix} \frac{1}{2} \tilde{\mathbf{y}}_{k-M+1}^\top \mathbf{H}^{[l+1]} \tilde{\mathbf{y}}_{k-M+1} \\ \vdots \\ \frac{1}{2} \tilde{\mathbf{y}}_k^\top \mathbf{H}^{[l+1]} \tilde{\mathbf{y}}_k \end{bmatrix} \\ &= \underbrace{(\Phi^\top \Phi)^{-1} (\Phi^\top \Phi)}_{=I} \text{vecr}(\mathbf{H}^{[l+1]}) \end{aligned} \quad (3.124)$$

ist⁶⁶, wobei im letzten Schritt der Zusammenhang

$$\frac{1}{2} \tilde{\mathbf{y}}_\kappa^\top \mathbf{H}^{[l+1]} \tilde{\mathbf{y}}_\kappa = \phi(\tilde{\mathbf{y}}_\kappa)^\top \text{vecr}(\mathbf{H}^{[l+1]}), \quad (3.125)$$

$\kappa = k - M + 1, k - M + 2, \dots, k$, genutzt wird. Da $\mathbf{H}^{[l+1]}$ eine symmetrische Matrix ist, die aus den nicht-redundanten Elementen von $\hat{\mathbf{w}}^{[l+1]}$ gebildet wird, folgt aus (3.124) direkt, dass die Folge $\hat{\mathbf{w}}^{[l]}$, d. h. die Value Iteration, äquivalent zur Folge (3.119) ist. \square

Das folgende technische Lemma wird für den Beweis von Lemma 3.8 benötigt.

Lemma 3.7

Sei $\mathbf{H}^{[0]} = \mathbf{0}$, $\mathbf{R}^\top = \mathbf{R} \succ \mathbf{0}$ und $\mathbf{Q}^\top = \mathbf{Q} \succeq \mathbf{0}$. Dann ist

$$\mathbf{u}_{k+\kappa} = -\mathbf{K}^{[l]} \begin{bmatrix} \mathbf{x}_{k+\kappa}^\top & \mathbf{x}_{r,k+\kappa+1}^\top & \dots & \mathbf{x}_{r,k+\kappa+n_h}^\top \end{bmatrix}, \quad \forall l > 0, \quad (3.126)$$

wobei $\mathbf{K}^{[l]} = \mathbf{K}(\mathbf{H}^{[l]})$ analog zu (3.117) aus $\mathbf{H}^{[l]}$ berechnet wird, die eindeutige Lösung, die

$$\hat{\mathbf{Q}}_\kappa^{[l]} = \frac{1}{2} \tilde{\mathbf{y}}_{k+\kappa}^\top \mathbf{H}^{[l]} \tilde{\mathbf{y}}_{k+\kappa} \quad (3.127)$$

minimiert.

⁶⁶ Durch $\text{vecr}(\cdot)$ wird die symmetrische Matrix $\mathbf{H}^{[l]}$ in einen Vektor überführt, d. h. es gilt $\text{vecr}(\mathbf{H}^{[l]}) = \hat{\mathbf{w}}^{[l]}$.

Beweis:

Aufgrund von $\mathbf{H}^{[0]} = \mathbf{0} \succcurlyeq \mathbf{0}$, $\mathbf{R} \succcurlyeq \mathbf{0}$ und $\mathbf{Q} \succcurlyeq \mathbf{0}$ folgt aus (3.119), dass $\mathbf{H}^{[l]} \succcurlyeq \mathbf{0}$, $\forall l > 0$. Weiterhin ist wegen $\mathbf{R} \succcurlyeq \mathbf{0}$ ersichtlich, dass $\mathbf{h}_{uu}^{[l]} \succcurlyeq \mathbf{0}$, $\forall l > 0$, gilt. Weil aus $\frac{\partial \hat{Q}_k^{[l]}}{\partial \mathbf{u}_{k+\kappa}} = \mathbf{0}$ die Reglermatrix $\mathbf{K}^{[l]}$ folgt und sich wegen $\mathbf{h}_{uu}^{[l]} \succcurlyeq \mathbf{0}$ zudem $\frac{\partial^2 \hat{Q}_k^{[l]}}{\partial \mathbf{u}_{k+\kappa}^2} \succcurlyeq \mathbf{0}$ ergibt, gilt die Aussage von Lemma 3.7. \square

Im Folgenden werde der durch

$$F\left(\boldsymbol{\Omega}^{[l]}, \boldsymbol{\Gamma}^{[l]}\right) := \mathbf{G} + \gamma \mathbf{M}\left(\boldsymbol{\Gamma}^{[l]}\right)^{\top} \boldsymbol{\Omega}^{[l]} \mathbf{M}\left(\boldsymbol{\Gamma}^{[l]}\right) \quad (3.128)$$

definierte Operator verwendet. Somit gilt nach (3.119) $F\left(\mathbf{H}^{[l]}, -\mathbf{K}^{[l]}\right) = \mathbf{H}^{[l+1]}$.

Um zu beweisen, dass $\mathbf{H}^{[l]}$ nach Lemma 3.6 im Sinne von $\mathbf{0} \preceq \mathbf{H}^{[l]} \preceq \mathbf{Y}$ nach oben beschränkt ist, wird zunächst das nachfolgende technische Lemma benötigt, das eine Generalisierung von [Lan97, Lemma B.1.1] auf die Solltrajektorienabhängige Q-Function darstellt. Die Kenntnis der analytischen Struktur von \mathbf{H} (vgl. Satz 3.4) ist hierbei wesentlich für die Erweiterung auf den Solltrajektorienfall.

Lemma 3.8

Sei $\mathbf{W}^{[l]}$ eine beliebige Matrixfolge, wobei $\mathbf{W}^{[l]}$ dieselbe Dimension wie $\mathbf{K}^{[l]}$ besitze, zudem gelte $\mathbf{0} \preceq \mathbf{H}^{[0]} \preceq \mathbf{Z}^{[0]}$. Dann ergibt sich mit dem in (3.128) definierten Operator $F(\cdot)$ für die Matrixfolgen

$$\mathbf{Z}^{[l+1]} := F\left(\mathbf{Z}^{[l]}, \mathbf{W}^{[l]}\right) \quad (3.129)$$

und

$$\mathbf{H}^{[l+1]} = F\left(\mathbf{H}^{[l]}, -\mathbf{K}\left(\mathbf{H}^{[l]}\right)\right), \quad (3.130)$$

dass

$$\mathbf{0} \preceq \mathbf{H}^{[l+1]} \preceq \mathbf{Z}^{[l+1]} \quad (3.131)$$

gilt.

Beweis:

Der Beweis ist in Anhang A.2 gegeben. \square

Im nächsten Schritt wird die Beschränktheit von $\mathbf{H}^{[l]}$ im Sinne von $\mathbf{0} \preceq \mathbf{H}^{[l]} \preceq \mathbf{Y}$ gezeigt. Für den Regulationsfall, d. h. für eine Q-Function ohne Solltrajektorienvorgabe, wurde diese Beschränktheit von Landelius [Lan97, Lemma B.1.2] bewiesen. Im Gegensatz dazu wird im

Folgenden der allgemeinere Solltrajektorienregelungsfall, d. h. die durch (3.119) beschriebene Iteration, betrachtet.

Lemma 3.9

Sei (A, B) steuerbar und $H^{[l]}$ die durch (3.119) beschriebene Folge mit $H^{[0]} = \mathbf{0}$. Dann existiert eine Matrix Y , sodass $\mathbf{0} \preceq H^{[l]} \preceq Y$ gilt.

Beweis:

Siehe Anhang A.3. □

Nach dieser Vorarbeit kann schließlich das Hauptergebnis der Konvergenzanalyse formuliert werden. Dieses besagt, dass die Value Iteration gegen das optimale Gewicht w^* und somit gegen die optimale Reglermatrix K^* des Solltrajektorienfolgeregelungsproblems, das durch Problem 3.3 gegeben ist, konvergiert.

Satz 3.5

Sei $R^\top = R \succ \mathbf{0}$, $Q^\top = Q \succeq \mathbf{0}$, (A, B) steuerbar und $\hat{w}^{[0]} = \mathbf{0}$, d. h. $H^{[0]} = \mathbf{0}$. Dann führt im Fall einer erfüllten Anregungsbedingung (3.115) die durch (3.112) und (3.117) gegebene Value Iteration zu $H^{[l]} \rightarrow H$, d. h. zur Konvergenz von $\hat{w}^{[l]} \rightarrow w^*$ und $K^{[l]} \rightarrow K^*$.

Beweis:

Der Beweis erfolgt mittels vollständiger Induktion. Nach Lemma 3.6 ist die betrachtete Value Iteration äquivalent zur Folge $H^{[l]}$ (vgl. (3.119)). Mit $Z^{[0]} = H^{[0]}$ und

$$Z^{[l+1]} = F\left(Z^{[l]}, -K\left(H^{[l+1]}\right)\right) \quad (3.132)$$

folgt

$$\mathbf{0} \preceq H^{[l]} \preceq Z^{[l]} \quad (3.133)$$

aus Lemma 3.8. Aufgrund von $H^{[0]} = \mathbf{0}$ und mit G aus (3.120) folgt $H^{[1]} = G \succeq \mathbf{0}$ und somit der Induktionsanfang

$$H^{[1]} - Z^{[0]} \succeq \mathbf{0}. \quad (3.134)$$

Die Induktionsbehauptung sei durch

$$H^{[l]} - Z^{[l-1]} \succeq \mathbf{0} \quad (3.135)$$

gegeben. Es gilt

$$H^{[l+1]} - Z^{[l]} = \gamma M\left(-K\left(H^{[l]}\right)\right)^\top \left(H^{[l]} - Z^{[l-1]}\right) M\left(-K\left(H^{[l]}\right)\right) \succeq \mathbf{0}. \quad (3.136)$$

Daraus folgt

$$\mathbf{0} \preceq H^{[l]} \preceq Z^{[l]} \preceq H^{[l+1]}. \quad (3.137)$$

Da die Matrixfolge $\mathbf{H}^{[l]}$ nach Lemma 3.9 durch \mathbf{Y} nach oben beschränkt ist und wegen (3.137)

$$\mathbf{0} \preceq \mathbf{H}^{[l]} \preceq \mathbf{H}^{[l+1]} \quad (3.138)$$

gilt, existiert der Grenzwert $\mathbf{H}^{[\infty]}$, d. h. die Value Iteration konvergiert zu

$$\mathbf{H}^{[\infty]} = \mathbf{G} + \gamma \mathbf{M} \left(-\mathbf{K} \left(\mathbf{H}^{[\infty]} \right) \right)^\top \mathbf{H}^{[\infty]} \mathbf{M} \left(-\mathbf{K} \left(\mathbf{H}^{[\infty]} \right) \right). \quad (3.139)$$

Weiterhin minimiert $-\mathbf{K} \left(\mathbf{H}^{[\infty]} \right)$ nach Lemma 3.7 $\hat{Q}_0^{[\infty]}$. Daraus ergibt sich

$$\lim_{l \rightarrow \infty} \delta_k^{[l]} = r_k + \gamma \tilde{\mathbf{y}}_{k+1}^*{}^\top \mathbf{H}^{[\infty]} \tilde{\mathbf{y}}_{k+1}^* - \tilde{\mathbf{y}}_k^\top F \left(\mathbf{H}^{[\infty]}, -\mathbf{K} \left(\mathbf{H}^{[\infty]} \right) \right) \tilde{\mathbf{y}}_k = 0. \quad (3.140)$$

Somit ist die Bellman-Gleichung für $\mathbf{H}^{[\infty]}$ erfüllt und es folgt $\mathbf{H}^{[l]} \rightarrow \mathbf{H}^{[\infty]} = \mathbf{H}$, d. h. $\hat{\mathbf{w}}^{[l]} \rightarrow \mathbf{w}^*$ und $\mathbf{K}^{[l]} \rightarrow \mathbf{K}^*$. \square

3.3.4 Ergebnisse

Im Folgenden wird der vorgestellte ADP-basierte Solltrajektorienfolgeregler, der dem Sollzustandsverlauf auf einem gleitenden Vorausschauhorizont der Länge n_h optimal im Sinne von Problem 3.3 folgt, und der ohne explizite Kenntnis der Systemmatrizen \mathbf{A} und \mathbf{B} trainiert wird, simulativ ausgewertet. Die Ergebnisse werden zudem mit einer ADP-Methode aus der Literatur verglichen. Die Anzahl der Datentupel werde zunächst zu $M = 1,2h$ gesetzt, wobei h jeweils aus (3.107) folgt. Die Abbruchbedingung wird zu $e_{\hat{\mathbf{w}}} = 10^{-3}$ gewählt, die Diskontierung zunächst zu $\gamma = 0,9$ gesetzt und der Vorausschauhorizont, auf dem der Sollzustandsverlauf gegeben ist, zu $n_h = 10$ gewählt. Der Anfangszustand des Systems sei jeweils $\mathbf{x}_0 = \mathbf{0}$.

3.3.4.1 Simulationsbeispiele

Betrachtet werden im Folgenden zwei Simulationsbeispiele. Das erste System ist ein rotatorisches Feder-Masse-Dämpfer-System zweiter Ordnung, das zweite System ein lineares Einspurmodell eines vorderradgelenkten Fahrzeugs der Ordnung sechs. Beide Systemmodelle resultieren aus einer Tustin-Approximation eines zeitkontinuierlichen Modells mit einer Abtastzeit von 1 s.

System 1 – Feder-Masse-Dämpfer-System

Die zeitdiskrete Zustandsraumrepräsentation des rotatorischen Feder-Masse-Dämpfer-Systems ist durch

$$\mathbf{x}_{k+1} = \begin{bmatrix} 0,99 & 0,9 \\ -0,02 & 0,8 \end{bmatrix} \mathbf{x}_k + \begin{bmatrix} 0,01 \\ 0,02 \end{bmatrix} u_k \quad (3.141)$$

gegeben. Die Stellgröße u_k stellt dabei ein Drehmoment dar, welches auf das System wirkt. Weiterhin seien

$$\mathbf{Q} = \text{diag}(100, 0) \text{ und } \mathbf{R} = 1 \quad (3.142)$$

die Parameter des Gütefunktional (3.82). Somit liegt der Fokus darauf, den Systemzustand x_1 einem Sollwinkel $x_{1,\text{ref}} = \alpha_{\text{ref}}$, der auf einem gleitenden Vorausschauhorizont der Länge $n_h = 10$ gegeben ist, folgen zu lassen.

Für dieses Optimierungsproblem liegen für $\gamma = 0,9$ die Systempole des optimal geregelten Systems alle im Einheitskreis, da $(\mathbf{A} - \mathbf{B}\mathbf{K}_x^*)$ einen doppelten Eigenwert bei 0,7503 aufweist. Das hinreichende Stabilitätskriterium nach Satz 3.3 (vgl. Bemerkung 3.9) ist wegen $b_\gamma = 0,7067 < 0,9 = \gamma$ ebenfalls erfüllt, aus einer korrekt gelernten Q-Function lässt sich für dieses Szenario somit auch ohne explizite Kenntnis der Systemmatrizen \mathbf{A} und \mathbf{B} auf Stabilität des geregelten Systems schließen.

Während des Trainingsvorgangs wird für dieses Simulationsbeispiel für die erforderliche Anregung (vgl. (3.115)) der additive Rauschterm ξ_k der Stellgröße in (3.118) mit einer Varianz von $\mathbf{V} = 0,01$ gewählt. Außerdem wird während des Trainings mittelwertfreies weißes Gaußsches Rauschen $\xi_{\text{ref}} \sim \mathcal{N}(0; 0,01)$ zum Solltrajektorienverlauf addiert.

System 2 – Lineares Einspurmodell

Das zweite betrachtete Beispielsystem ist ein lineares Einspurmodell sechster Ordnung (vgl. [Fla16, Anhang B]). Dieses ist durch

$$\mathbf{x}_{k+1} = \begin{bmatrix} -0,633 & -0,035 & 0 & 0 & -1,4 \cdot 10^{-4} & -4,5 \cdot 10^{-3} \\ 2,445 & -0,771 & 0 & 0 & 5,1 \cdot 10^{-3} & 0,164 \\ 1,222 & 0,115 & 1 & 0 & 2,6 \cdot 10^{-3} & 0,082 \\ 15,890 & 0,794 & 20 & 1 & 0,024 & 0,774 \\ 6,197 & -0,323 & 0 & 0 & -0,925 & -1,587 \\ 3,099 & -0,161 & 0 & 0 & 0,038 & 0,206 \end{bmatrix} \mathbf{x}_k + \begin{bmatrix} -2,7 \cdot 10^{-4} \\ 9,9 \cdot 10^{-3} \\ 4,9 \cdot 10^{-3} \\ 0,047 \\ 0,145 \\ 0,073 \end{bmatrix} u_k \quad (3.143)$$

gegeben. Die physikalische Bedeutung der Stellgröße u_k und des Systemzustands \mathbf{x}_k ist in Anhang A.4 erläutert. Da im Folgenden die Abweichung von $x_{4,\text{ref}} = y_{\text{ref}}$ bestraft werden soll, werde

$$\mathbf{Q} = \text{diag}(0, 0, 0, 100, 0, 0) \text{ und } \mathbf{R} = 1 \quad (3.144)$$

gewählt.

Auch bei diesem Beispiel liegen für $\gamma = 0,9$ die Systempole des optimal geregelten Systems im Inneren des Einheitskreises. Konkret ergeben sich für $(\mathbf{A} - \mathbf{B}\mathbf{K}_x^*)$ doppelte Eigenwerte bei 0,4834 und 0,7136 sowie einfache Eigenwerte bei 0,0526 und 0,8675. Obwohl das optimal geregelte System somit auch für $\gamma = 0,9$ stabil ist, liefert das durch Satz 3.3 gegebene lediglich hinreichende Kriterium in diesem Fall wegen $b_\gamma \approx 1$ keine eindeutige Aussage. Daher soll für dieses Beispiel im Folgenden zusätzlich zum Fall $\gamma = 0,9$ noch der Fall $\gamma = 1$ betrachtet

werden. Für $\gamma = 1$ hat $(\mathbf{A} - \mathbf{BK}_x^*)$ doppelte Eigenwerte bei 0,4556 und 0,7136 sowie einfache Eigenwerte bei 0,0507 und 0,8675.

Zur Anregung (vgl. (3.115)) wird der additive Rauschterm ξ_k mit $\mathbf{V} = 0,03$ in (3.118) verwendet. Zudem wird während des Trainingsvorgangs mittelwertfreies weißes Gaußsches Rauschen $\xi_{\text{ref}} \sim \mathcal{N}(0; 0,01)$ zum Solltrajektorienverlauf addiert.

3.3.4.2 Auswertungsmethodik

Im Folgenden bezeichne α_{ref} die Sollwinkeltrajektorie für Beispielsystem 1 und y_{ref} die Solltrajektorie für Beispielsystem 2. Die im Rahmen dieser Arbeit vorgestellte ADP-Solltrajektorienfolgeregelungsmethode, die den exakten Solltrajektorienverlauf auf einem gleitenden Vorausschauhorizont der Länge n_h in die Q-Function integriert, soll mit Literaturmethoden, die annehmen, der Solltrajektorienverlauf werde durch $\mathbf{f}_{\mathbf{x}_r}(\mathbf{x}_{r,k}) = \mathbf{F}_{\text{ref}}\mathbf{x}_{r,k}$ beschrieben, verglichen werden. Um diesen Vergleich zu ermöglichen, wird für die während des Trainingsvorgangs verwendeten Solltrajektorien

$$\alpha_{\text{ref}} = y_{\text{ref}} = [1 \quad 0] \mathbf{x}_{r,k}, \quad (3.145)$$

mit

$$\mathbf{x}_{r,k+1} = \underbrace{\begin{bmatrix} 0,9801 & 0,1987 \\ -0,1987 & 0,9801 \end{bmatrix}}_{\mathbf{F}_{\text{ref}}} \mathbf{x}_{r,k}, \quad (3.146)$$

angenommen. Zu Vergleichszwecken wird sowohl ein ADP-Regler entsprechend Algorithmus 3.2 und ein ADP-Regler entsprechend [LLHW16], [KLM⁺14], nachfolgend als Exosystem-Methode bezeichnet, trainiert. Der Exosystem-Ansatz wird folglich bezüglich des erweiterten Systems

$$\tilde{\mathbf{x}}_{k+1} = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{F}_{\text{ref}} \end{bmatrix} \tilde{\mathbf{x}}_k + \begin{bmatrix} \mathbf{B} \\ \mathbf{0} \end{bmatrix} u_k, \quad (3.147)$$

mit $\tilde{\mathbf{x}}_k := [\mathbf{x}_k^\top \quad \mathbf{x}_{r,k}^\top]^\top$ trainiert.

Nach Abschluss des Trainingsprozesses wird der Solltrajektorienverlauf $\alpha_{\text{ref}} = y_{\text{ref}}$ deutlich variiert, um die Generalisierungsfähigkeit der ADP-Regler auf von (3.146) abweichende Solltrajektorienverläufe zu untersuchen. Die Bewertung der Ergebnisse geschieht dabei auf zwei Arten. Einerseits wird nach Beendigung des Trainingsvorgangs der Trajektorienverlauf, der sich unter Verwendung der gelernten ADP-Regler ergibt (bezeichnet durch α_{ADP} und y_{ADP} bzw. $\alpha_{\text{ADP,Exosystem}}$ und $y_{\text{ADP,Exosystem}}$), mit dem Trajektorienverlauf α_{opt} bzw. y_{opt} , der sich aus der optimalen analytischen Lösung nach Satz 3.4 ergibt, verglichen. Die mittlere quadratische Abweichung wird dabei mit α_{RMS} bzw. y_{RMS} bezeichnet. Andererseits werden die durch Algorithmus 3.2 gelernten Gewichte $\hat{\mathbf{w}}$ der in dieser Arbeit vorgestellten ADP-Methode

mit der optimalen Lösung \mathbf{w}^* , d. h. den zu Satz 3.4 und Lemma 3.5 korrespondierenden Gewichten, verglichen⁶⁷. Zugunsten einer Vergleichbarkeit für unterschiedliche Wertebereiche von $\hat{\mathbf{w}}$ wird der Absolutfehler jedes Gewichts bezüglich des größten absoluten Gewichts $\max_j |\{\mathbf{w}^*\}_j|$ normiert. Der Mittelwert dieser normierten Absolutfehler über alle h Gewichte wird durch

$$e_I = \frac{1}{h} \sum_{i=1}^h \frac{|\{\hat{\mathbf{w}}\}_i - \{\mathbf{w}^*\}_i|}{\max_j |\{\mathbf{w}^*\}_j|} \quad (3.148)$$

gegeben, der Maximalwert des normierten Absolutfehlers hingegen durch

$$e_{II} = \max_{i \in \{1, \dots, h\}} \frac{|\{\hat{\mathbf{w}}\}_i - \{\mathbf{w}^*\}_i|}{\max_j |\{\mathbf{w}^*\}_j|}. \quad (3.149)$$

3.3.4.3 Simulationsergebnisse

Die Trajektorienabweichungen α_{RMS} (System 1) und y_{RMS} (System 2) von der optimalen Lösung sowie die Schätzfehler e_I und e_{II} der Gewichte $\hat{\mathbf{w}}$ sind in Tabelle 3.3 gegeben⁶⁸.

Die resultierenden Trajektorienverläufe des optimalen Reglers sowie der beiden ADP-Ansätze sind in Abbildung 3.7 (Beispielsystem 1) und Abbildung 3.8 (Beispielsystem 2 mit $\gamma = 0,9$) gegeben. Für $\gamma = 1$ ergibt sich für System 2 ein nahezu identischer Verlauf zur Diskontierung $\gamma = 0,9$, dieser ist der Vollständigkeit halber in Anhang A.5 (Abbildung A.2) gegeben. Die vertikale, punkt-gestrichelte Linie gibt hierbei den Zeitschritt an, ab dem $k \geq M$ gilt und somit die Value Iteration nach Algorithmus 3.2 durchgeführt wird. Diese erreicht nach 42 (System 1), 43 (System 2, $\gamma = 0,9$), bzw. 67 (System 2, $\gamma = 1$) Iterationen die durch $e_{\hat{\mathbf{w}}}$ beschriebene Abbruchbedingung. Der Solltrajektorienverlauf α_{ref} bzw. y_{ref} ist in Grau dargestellt. Die schwarze, gestrichelte Linie zeigt den Verlauf der optimalen Lösungen α_{opt} bzw. y_{opt} . In Rot ist der Verlauf gezeigt, der sich für die in dieser Arbeit vorgestellte ADP-Methode ergibt, wohingegen in Blau die Ergebnisse der Vergleichsmethode gezeigt sind.

Die Gewichtsfehler e_I und e_{II} in Tabelle 3.3 lassen erkennen, dass die optimalen Q-Function-Gewichte $\hat{\mathbf{w}}$ bei beiden Systemen erfolgreich gelernt werden. Die Abnahme der Fehlermaße e_I und e_{II} während der Value Iteration nach Algorithmus 3.2 ist in Abbildung 3.9 (System 1), Abbildung 3.10 (System 2 mit $\gamma = 0,9$) und Abbildung A.3 (System 2 mit $\gamma = 1$) gezeigt.

⁶⁷ Es sei anzumerken, dass die zu Vergleichszwecken verwendete Exosystem-Methode nicht \mathbf{w}^* aus Satz 3.4 und Lemma 3.5 lernt, sondern die optimalen Gewichte, welche mit (3.147) korrespondieren. Daher werden die nachfolgend vorgestellten Bewertungsmaße nur für die in dieser Arbeit vorgestellte ADP-Solltrajektorienregelungsmethode berechnet, um zu bewerten, ob ein erfolgreicher Trainingsvorgang vorliegt.

⁶⁸ An dieser Stelle sei anzumerken, dass die Exosystem-Methode, welche als Vergleichsalgorithmus verwendet wird, aufgrund der Wahl von (3.146) für $\gamma = 1$ ungeeignet ist, da der durch \mathbf{F}_{ref} beschriebene Solltrajektorienverlauf nicht abklingt, da \mathbf{F}_{ref} zwei Eigenwerte auf dem Einheitskreis der komplexen Ebene aufweist.

System 1			
	ADP	ADP, Exosystem	
α_{RMS}	$2,1 \cdot 10^{-3}$	1,55	
e_{I}	$6,2 \cdot 10^{-5}$	-	
e_{II}	$2,1 \cdot 10^{-3}$	-	

System 2			
	ADP ($\gamma = 0,9$)	ADP ($\gamma = 1$)	ADP, Exosystem ($\gamma = 0,9$)
y_{RMS}	$6,5 \cdot 10^{-5}$	$5,2 \cdot 10^{-4}$	1,00
e_{I}	$9,5 \cdot 10^{-6}$	$4,2 \cdot 10^{-5}$	-
e_{II}	$1,0 \cdot 10^{-3}$	$4,5 \cdot 10^{-3}$	-

Tabelle 3.3: Trajektorienfolgefehler und Gewichtsfehler der betrachteten Simulationsbeispiele.

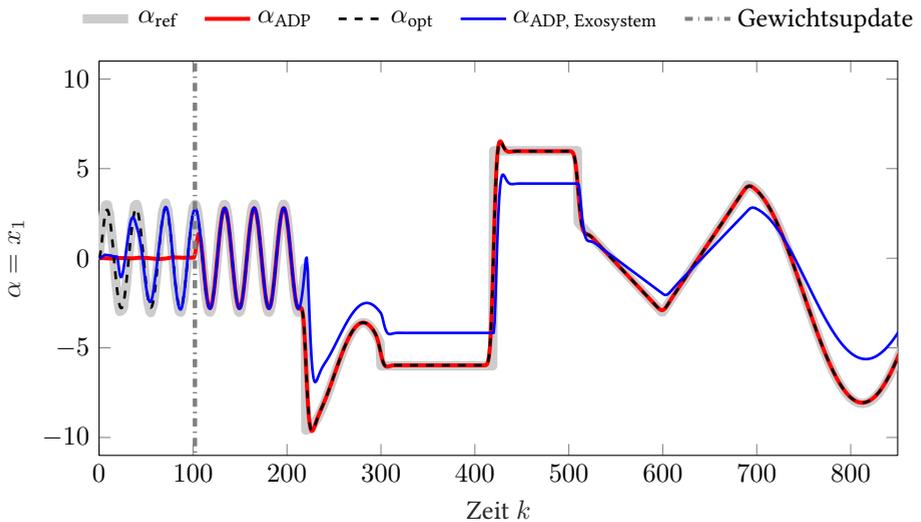


Abbildung 3.7: Ergebnis der ADP-Solltrajektorienregelung für System 1 (rotatorisches Feder-Masse-Dämpfer-System zweiter Ordnung).

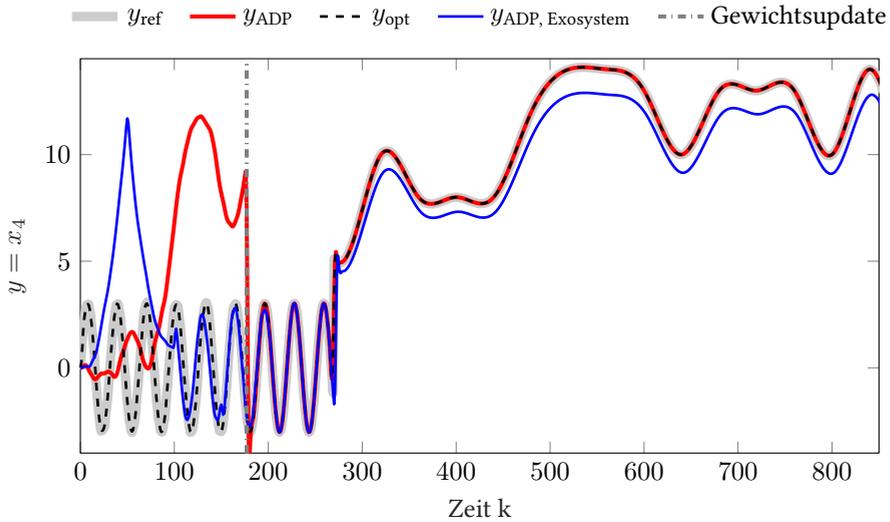


Abbildung 3.8: Ergebnis der ADP-Solltrajektorienregelung für System 2 (lineares Einspurmodell sechster Ordnung) für $\gamma = 0,9$.

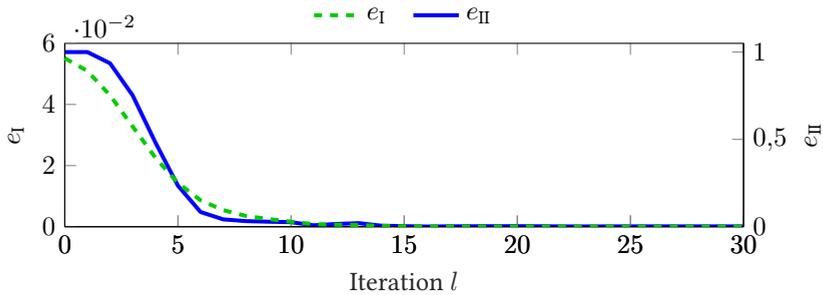


Abbildung 3.9: Gewichtsfehlerverlauf während der ersten 30 Iterationen des Lernvorgangs für System 1 (rotatorisches Feder-Masse-Dämpfer-System zweiter Ordnung). Hierbei stellt e_I (3.148) den Mittelwert und e_{II} (3.149) das Maximum des elementweisen absoluten Fehlers von \hat{w} , jeweils durch $\max_j |\{w^*\}_j|$ normiert, dar.

3.3.4.4 Einfluss von Messrauschen

Bisher wurde, wie in der ADP-Literatur üblich, angenommen, dass exakte Messungen des Systemzustands x_k vorliegen. Wenngleich der Fokus der vorliegenden Arbeit auf einer Formulierung einer Q-Function, die den Solltrajektorienverlauf einbezieht, und nicht auf robusten

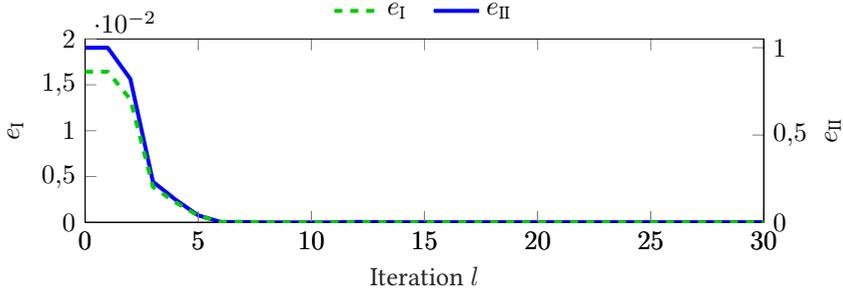


Abbildung 3.10: Gewichtsfehlerverlauf während des Lernvorgangs für System 2 (lineares Einspurmodell sechster Ordnung) für $\gamma = 0,9$. Hierbei stellt e_I (3.148) den Mittelwert und e_{II} (3.149) das Maximum des elementweisen absoluten Fehlers von $\hat{\mathbf{w}}$, jeweils durch $\max_j \{|\mathbf{w}^*\}_j$ normiert, dar.

ADP-Methoden liegt⁶⁹, soll im Folgenden kurz der Einfluss unbekanntes Messrauschens⁷⁰ ν_k untersucht werden. Hierzu werde angenommen, dass dem ADP-basierten Regelungsansatz lediglich ein gemessener Systemzustand

$$\tilde{\mathbf{x}}_k = \mathbf{x}_k + \nu_k \quad (3.150)$$

zur Verfügung steht, wobei $\nu_k \in \mathbb{R}^n$ mittelwertfreies weißes Gaußsches Rauschen darstellt.

In dieser Situation werde der Regler für System 1 (3.141) mit $\gamma = 0,9$, $M = 10h$ (mit h wie in (3.107)) und für $n_h = 10$ trainiert. Während des Trainings werden hierbei unterschiedliche Signal-Rausch-Verhältnisse $\text{SNR}_x(\mathbf{x}_k, \nu_k)$ und $\text{SNR}_{x_r}(\mathbf{x}_{r,k}, \nu_k)$ des Zustands \mathbf{x}_k und Sollzustands $\mathbf{x}_{r,k}$ im Vergleich zum Messrauschen ν_k betrachtet. Die Abweichung des gelernten Reglers \mathbf{K} vom optimalen Regler \mathbf{K}^* ist in Tabelle 3.4 gegeben. Die resultierenden Trajektorienverläufe sowie deren Standardabweichungen ($\text{std}(\cdot)$) für unterschiedliche Rauschamplituden sind in Abbildung 3.11 gezeigt. Je stärker das Anregungsrauschen im Vergleich zum Messrauschen dominiert, desto ähnlicher ist der gelernte Regler dem optimalen Regler und desto besser folgt der Winkel α der Solltrajektorie.

⁶⁹ Für erste Ansätze zu robusten ADP-Methoden, welche Gegenstand der aktuellen Forschung sind, sei auf [BJD10], [JJ14c], [NLW⁺19], [QZLY19], [WHL17], [WY18] und [ZCZL11] verwiesen.

⁷⁰ Dabei ist das unbekanntes Messrauschen nicht mit dem Anregungsrauschen ξ_k in (3.118), welches bekannt ist und daher korrekt in den TD-Fehler einbezogen wird, zu verwechseln.

$\text{SNR}_x = \text{SNR}_{x_r}$	$\ \mathbf{K}^* - \mathbf{K}\ _2$	$\frac{\ \mathbf{K}^* - \mathbf{K}\ _2}{\ \mathbf{K}^*\ _2}$
10 dB	6,61	0,40
20 dB	1,61	$9,66 \cdot 10^{-2}$
30 dB	0,39	$2,37 \cdot 10^{-2}$
40 dB	0,15	$8,84 \cdot 10^{-3}$
50 dB	$4,75 \cdot 10^{-2}$	$2,84 \cdot 10^{-3}$
60 dB	$1,36 \cdot 10^{-2}$	$8,13 \cdot 10^{-4}$

Tabelle 3.4: Abweichungen der gelernten Regler vom optimalen Regler für unterschiedlich starkes Messrauschen. Die angegebenen Fehler stellen über 20 Trainingsdurchgänge gemittelte Werte dar.

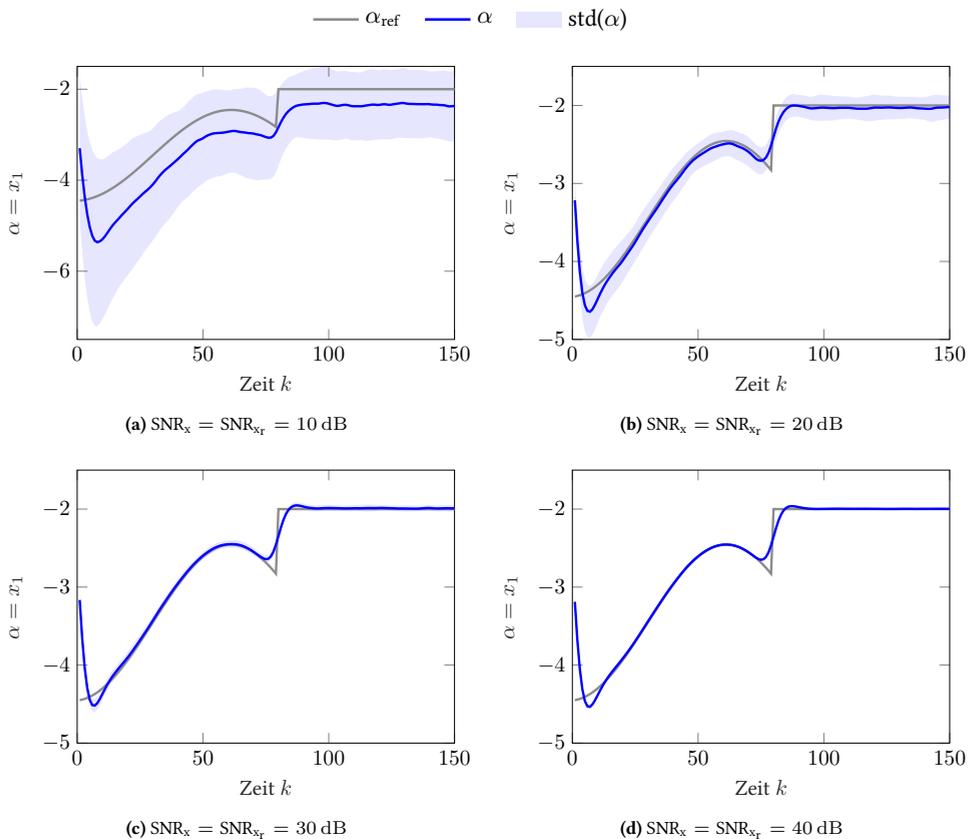


Abbildung 3.11: Trajektorienfolgeergebnisse für ADP-Regler, die bei unterschiedlich starkem Messrauschen trainiert und validiert wurden.

3.3.5 Diskussion der Simulationsergebnisse

Für den in dieser Arbeit vorgestellten ADP-Solltrajektorienfolgeregler, der den Sollzustandsverlauf auf einem gleitenden Vorausschauhorizont explizit einbezieht, wurde die Konvergenz gegen die optimale Lösung unter einer geeigneten Anregung nach (3.115) in Abschnitt 3.3.3.3 formal gezeigt. Insbesondere offenbaren Abbildung 3.9, Abbildung 3.10 und Abbildung A.3, dass die geschätzte solltrajektorienabhängige Q-Function \hat{Q}_0 für die untersuchten Simulationsbeispiele innerhalb weniger Iterationen gegen die optimale Q-Function Q_0 konvergiert. Aufgrund des optimalen Verhaltens entsprechend des durch (3.82) gegebenen Gütefunktional und des Vorausschauhorizonts der Länge n_h weist der gelernte Regler des Weiteren prädiktives Verhalten auf und reagiert beispielsweise auf Sprünge bereits vor der Veränderung des aktuellen Sollzustands. Dies ist in Abbildung 3.12, welche die Detailansicht eines Sprungs aus Abbildung 3.7 zeigt, zu sehen.

Die Betrachtung von α_{RMS} (System 1) bzw. y_{RMS} (System 2) sowie die Abbildungen 3.7, 3.8 und 3.12 verdeutlichen zudem, dass der in dieser Arbeit präsentierte ADP-Solltrajektorienfolgeregelungsansatz erfolgreich und ohne weitere Modifikation auch Solltrajektorien, die von den während des Trainingsvorgangs verwendeten Trajektorienverläufen stark abweichen, folgen kann. Diese fundamentale Eigenschaft, welche die Verwendung flexibler und während des Trainingsvorgangs nicht vorliegender Solltrajektorien ermöglicht, ist insbesondere auf die explizite Abhängigkeit der gelernten Q-Function von einem beliebigen Solltrajektorienverlauf auf einem gleitenden Vorausschauhorizont der Länge n_h zurückzuführen. Im Gegensatz dazu sind existierende ADP-Methoden, die annehmen, der Sollzustandsverlauf folge der durch f_{x_r} beschriebenen Dynamik (beispielsweise [LLHW16] und [KLM⁺14]), zwar gut geeignet, solange diese Annahme erfüllt ist (wie dies beispielsweise in Abbildung 3.7 für $100 < k < 200$

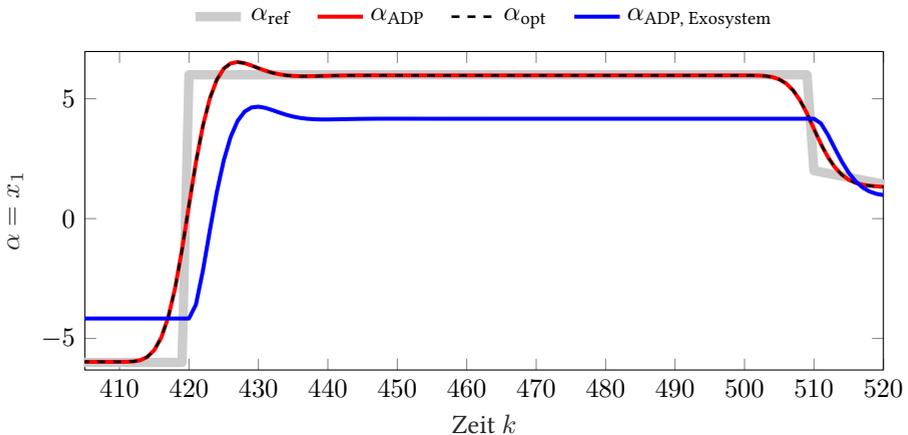


Abbildung 3.12: Detailansicht von Abbildung 3.7 (System 1), um das prädiktive Verhalten des vorgestellten ADP-Ansatzes aufgrund der Verwendung des gleitenden Vorausschauhorizonts der Länge n_h zu visualisieren.

der Fall ist), jedoch nimmt deren Güte ab, sobald der Solltrajektorienverlauf von (3.145)–(3.146) abweicht. Dieses Verhalten ist wenig verwunderlich, da bei diesen Reglern implizit die während des Trainingsvorgangs zugrunde liegende Exosystemdynamik auf die gelernte Q-Function Einfluss nimmt.

Besonders erwähnenswert ist auch, dass die genaue Kenntnis der Struktur der Matrix \mathbf{H}_K als Ergebnis von Satz 3.4 und somit auch der Matrix \mathbf{H} signifikant zu einer effizienten Wahl geeigneter Basisfunktionen beiträgt und eine wesentliche Reduktion der Anzahl h der zu lernenden Gewichte ermöglicht. Würde beispielsweise nur angenommen, \mathbf{H} sei symmetrisch, so müssten bei den gezeigten Simulationsbeispielen $h = 325$ (System 1) bzw. $h = 2701$ (System 2) Gewichte geschätzt werden⁷¹. Unter Beachtung von Lemma 3.4 reduziert sich diese Anzahl um rund 24 % (System 1) bzw. 26 % (System 2). Nutzung der dünnbesetzten Matrix \mathbf{Q} nach Lemma 3.5 reduziert die zu lernenden Gewichte schließlich sogar um etwa 74 % (System 1) bzw. 95 % (System 2) im Vergleich zur reinen Berücksichtigung der Symmetrie von \mathbf{H} (vgl. Tabelle 3.2 für den Fall $n_h = 10$). Wie (3.107) zeigt, beeinflusst die Wahl der Länge n_h des Vorausschauhorizonts direkt die Anzahl h der zu lernenden Gewichte. Für den Fall, dass der Solltrajektorienverlauf über einen großen Vorausschauhorizont bekannt ist, kann somit zwar ein potenziell besseres prädiktives Verhalten erreicht werden, die Komplexität des ADP-Problems erhöht sich jedoch entsprechend. Somit ist bei der Wahl einer geeigneten Vorausschauhorizontlänge n_h ein von der spezifischen Anwendung abhängiger Kompromiss einzugehen.

3.4 Resümee zur zeitdiskreten ADP-basierten Solltrajektorienfolgeregelung

Zusammenfassend wurde in Kapitel 3 erstmalig formal der Begriff zeitdiskreter ADP-kompatibler Solltrajektorienendarstellungen definiert. Zudem wurden zwei neuartige, flexible und effiziente ADP-basierte Solltrajektorienfolgeregelungsansätze präsentiert und analysiert. Ohne exakte Kenntnis der Systemdynamik können so Regelgesetze erlernt werden, die generalisieren, indem während des Trainingsvorgangs nicht gesehene Solltrajektorienverläufe optimal gefolgt werden kann. Dies stellt einen wesentlichen Vorteil gegenüber bestehenden Vergleichsmethoden dar. Des Weiteren wird durch die vorgestellten Ansätze auch der zukünftige Verlauf der Solltrajektorie berücksichtigt, wodurch die resultierenden Regler prädiktives Verhalten aufweisen.

Die beiden neuen, zeitdiskreten Solltrajektorienendarstellungen sind dabei ADP-kompatibel nach Definition 3.1 (vgl. Bemerkung 3.2 und Bemerkung 3.7). Bei der in Abschnitt 3.2 präsentierten lokalen Beschreibung des Solltrajektorienverlaufs durch den Parameter \mathbf{Z}_k und Basisfunktionsvektor $\boldsymbol{\rho}$ nach (3.15) entscheidet insbesondere die Wahl von $\boldsymbol{\rho}$ über die Approximationsfähigkeit des Sollzustandsverlaufs und die Anzahl h der zu lernenden Gewichte. Wenngleich eine lokale Approximation des Solltrajektorienverlaufs, wie beispielsweise die in

⁷¹ Selbst in diesem Fall wird schon das Wissen genutzt, dass die Q-Function quadratisch ist.

Abschnitt 3.2.5.1 gezeigten kubischen Polynome, gewisse Einschränkungen an die Vorausschaufähigkeit von Solltrajektorienreglern implizieren kann, so bietet eine an die Problemstellung angepasste Wahl der Solltrajektorienendarstellung im aktuellen Zeitschritt k einen entscheidenden Kompromiss zwischen Flexibilität und einer möglichst kompakten Darstellung mit handhabbarer Dimension h des Gewichtsvektors \hat{w} . Dieser neuartige Ansatz erlaubt somit erstmals allgemein die Entwicklung vorausschauender, flexibler, kompakter, modellfreier ADP-basierter Solltrajektorienfolgeregler.

Zudem wird in Abschnitt 3.3 der Fall betrachtet, direkt die Sollzustände, die auf einem gleitenden Vorausschauhorizont der Länge n_h gegeben sind, in die Q-Function zu integrieren, anstatt eine parametrische Beschreibung des Solltrajektorienverlaufs zu verwenden. Dieser Ansatz erlaubt einerseits eine völlig beliebige und exakte Vorgabe von Sollzuständen auf dem gleitenden Vorausschauhorizont, korreliert andererseits jedoch direkt mit einer quadratischen Zunahme der Dimension h des Gewichtsvektors \hat{w} bei steigender Vorausschauhorizontlänge n_h (vgl. beispielsweise (3.107)). Bei einer sehr großen Vorausschauhorizontlänge n_h stößt dieser Ansatz aufgrund der damit verbundenen unverhältnismäßigen Anzahl zu schätzender Critic-Gewichte somit an seine Grenzen. Beispielsweise zeigt Tabelle 3.2, dass für den Fall $n_h = 100$ selbst bei einem System zweiter Ordnung bereits Tausende Gewichte gelernt werden müssten. Dies geht neben einer höheren Anzahl benötigter Messdaten mit einer erhöhten Rechenzeit, einer potenziell erschwerten Erfüllung der Anregungsbedingung (3.115) und möglichen numerischen Problemen einher. Letztlich ist im Einzelfall abzuwägen, ob die direkte Verwendung der Sollzustände wie in Abschnitt 3.3 oder eine parametrische Approximation wie in Abschnitt 3.2 für eine konkrete Problemstellung von Vorteil ist. Aus theoretischer Sicht sei der Vollständigkeit halber noch angemerkt, dass die exakte Verwendung des Solltrajektorienverlaufs auf einem gleitenden Vorausschauhorizont der Länge n_h nach Abschnitt 3.3 letztlich als Spezialfall der parametrischen Darstellung aus Abschnitt 3.2 interpretiert werden kann. Hierbei entspricht \mathbf{Z}_k den unveränderten Sollzuständen auf dem Vorausschauhorizont der Länge n_h und $\rho(\kappa)$ fungiert als Schieberegister.

Schließlich sei noch zu betonen, dass die im vorliegenden Kapitel vorgestellten Methoden zur ADP-kompatiblen Solltrajektorienrepräsentation grundsätzlich unabhängig von der konkreten Wahl des ADP-Algorithmus sind. Während in Abschnitt 3.2 beispielhaft eine Policy Iteration und in Abschnitt 3.3 eine Value Iteration verwendet wurde, können die vorgestellten ADP-kompatiblen Solltrajektorienrepräsentationen ebenso im Rahmen von Actor-Critic-Ansätzen angewandt werden. Hierzu sei exemplarisch auf das Anwendungsbeispiel in Abschnitt 6.1 verwiesen. Ein weiterer Vorteil der vorgestellten ADP-kompatiblen Solltrajektorienrepräsentationen ist zudem, dass, wann immer bei den vorgestellten Ansätzen eine Off-Policy-Methode zum Einsatz kommt, vom System aufgezeichnete Daten $\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}$ während des Trainingsvorgangs wiederverwendet werden können, um die Dateneffizienz zu steigern (vgl. Abschnitt 2.1.4.4).

Die vorgestellten, neuartigen, flexiblen Mechanismen für ADP-basierte Solltrajektorienfolgeregler beantworten somit für den Fall zeitdiskreter Problemstellungen die in Abschnitt 2.4.1 formulierte Forschungsfrage 1 nach geeigneten Solltrajektorienrepräsentationen und ihrer effizienten Integration in modellfreie ADP-Ansätze. Im nächsten Kapitel wird, auf den Erkenntnissen des

vorliegenden Kapitels aufbauend, der zeitkontinuierliche ADP-basierte Solltrajektorienfolge-
regelungsfall betrachtet.

4 Zeitkontinuierliche ADP-basierte Solltrajektorienfolgeregelung

Nachdem im vorigen Kapitel zeitdiskrete ADP-kompatible Solltrajektorienendarstellungen betrachtet wurden, wird in diesem Kapitel die in Abschnitt 2.4.1 formulierte Forschungsfrage 1 für den zeitkontinuierlichen Fall beantwortet. Zunächst wird die Definition ADP-kompatibler Sollzustandsverläufe (vgl. Abschnitt 3.1) auf den zeitkontinuierlichen Fall übertragen. Anschließend wird eine ADP-kompatible zeitkontinuierliche Solltrajektorienendarstellung⁷² präsentiert, bei welcher der Solltrajektorienverlauf mithilfe eines von außen vorgebbaren Parameters definiert wird. Durch die explizite Abhängigkeit der Value Function nicht nur vom Systemzustand, sondern auch von diesem Trajektorienparameter, resultiert eine hohe Flexibilität der möglichen Solltrajektorienverläufe. Darauf aufbauend wird für linear-quadratische zeitkontinuierliche optimale Trajektorienfolgeregelungsprobleme der Einfluss eines globalen Diskontierungsfaktors untersucht sowie eine Problemstellung mit teilweiser Dämpfung formuliert. Für diese Problemstellung lassen sich anschließend Aussagen über die Stabilität der optimalen Lösung treffen. Die vorgestellte Solltrajektorienendarstellung wird in einen bestehenden ADP-Ansatz integriert, sodass der resultierende optimale Trajektorienfolgeregler ohne Kenntnis der Systemdynamik aus aufgezeichneten Daten erlernt werden kann. Simulationsergebnisse für verschiedene Hyperparameter der neuartigen ADP-kompatiblen Solltrajektorienendarstellung sowie eine anschließende Diskussion schließen das vorliegende Kapitel ab.

4.1 Definition ADP-kompatibler zeitkontinuierlicher Trajektorien

Analog zu Abschnitt 3.1 lässt sich auch für den zeitkontinuierlichen Fall eine ADP-kompatible Solltrajektorienendarstellung formulieren.

⁷² Dieser Abschnitt basiert auf einem im Rahmen der vorliegenden Dissertation entstandenen Konferenzbeitrag [BKIH21].

**Definition 4.1 (Zeitkontinuierliche ADP-kompatible Solltrajektorien-
darstellung)**

Eine mit dem ADP-Formalismus kompatible Solltrajektorien-
darstellung mit dem Sollzu-
stand $\mathbf{x}_r(t) \in \mathcal{X}$ ist durch zeitinvariante Funktionen $\mathbf{f}_{\mathbf{x}_r, \zeta} : \mathcal{Z} \rightarrow \mathcal{X}$ und $\mathbf{f}_\zeta : \mathcal{Z} \rightarrow \mathbb{R}^{n_\zeta}$
mit

$$\mathbf{x}_r(t) = \mathbf{f}_{\mathbf{x}_r, \zeta}(\zeta(t)) \quad (4.1)$$

und

$$\dot{\zeta}(t) = \mathbf{f}_\zeta(\zeta(t)) \quad (4.2)$$

charakterisiert, wobei $\zeta \in \mathcal{Z} \subset \mathbb{R}^{n_\zeta}, \forall t, n_\zeta < \infty$. Zudem sei $\mathbf{f}_{\mathbf{x}_r, \zeta}(\cdot)$ stetig und $\mathbf{f}_\zeta(\cdot)$
Lipschitz-stetig auf \mathcal{Z} .

Basierend auf Definition 4.1 lässt sich die folgende formale Aussage treffen.

Proposition 4.1

Die Solltrajektorien-
darstellung sei nach Definition 4.1 ADP-kompatibel. Für ein durch
die zeitkontinuierliche Systemdynamik (2.10) beschriebenes System⁷³ und ein Lipschitz-
stetiges Regelgesetz $\boldsymbol{\mu}(\cdot)$ gilt: Sind die durch

$$\int_t^\infty r(\mathbf{x}(\tau), \mathbf{x}_r(\tau), \boldsymbol{\mu}(\mathbf{x}(\tau), \mathbf{x}_r(\tau))) d\tau \quad (4.3)$$

gegebenen Gesamtkosten (vgl. (2.20)), die von der Solltrajektorie $\mathbf{x}_r(\tau), \tau \geq t$, nach
Definition 4.1 abhängen, endlich⁷⁴, dann können sie durch eine Value Function der Form
 $V^\mu(\mathbf{x}(t), \zeta(t))$ beschrieben werden.

Beweis:

Sei

$$V^\mu(\mathbf{x}(t), \zeta(t)) := \int_t^\infty r(\mathbf{x}(\tau), \mathbf{x}_r(\tau), \boldsymbol{\mu}(\mathbf{x}(\tau), \mathbf{x}_r(\tau))) d\tau \quad (4.4a)$$

$$= \int_t^\infty r(\mathbf{x}(\tau), \mathbf{f}_{\mathbf{x}_r, \zeta}(\zeta(\tau)), \boldsymbol{\mu}(\mathbf{x}(\tau), \mathbf{f}_{\mathbf{x}_r, \zeta}(\zeta(\tau)))) d\tau. \quad (4.4b)$$

Aufgrund von $\dot{\zeta}(\tau) = \mathbf{f}_\zeta(\zeta(\tau))$ mit $\mathbf{f}_\zeta(\cdot)$ Lipschitz-stetig existiert ausgehend von $\zeta(t)$ nach
dem Satz von Picard-Lindelöf [BSMM08, S. 673] $\forall \tau \geq t$ eine eindeutige Lösung für $\zeta(\tau) \in \mathcal{Z}$.
Somit legen $\zeta(t)$ und $\mathbf{f}_\zeta(\cdot) \forall \tau \geq t$ den Referenzparameter $\zeta(\tau)$ in (4.4b) eindeutig fest.

⁷³ Ebenso gilt die Aussage von Proposition 4.1 auch für allgemeinere nichtlineare Systeme der Form
 $\dot{\mathbf{x}}(t) = \mathbf{F}(\mathbf{x}(t), \mathbf{u}(t))$, aus Gründen der Einheitlichkeit wurde hier jedoch der eingangsauffine Fall formuliert.

⁷⁴ Für sinnvoll gestellte Probleme schließt dies neben der Stabilisierbarkeit des Systems und einem zulässigen
Regelgesetz $\boldsymbol{\mu}$ auch ein, dass durch die Wahl des Solltrajektorienverlaufs oder im Falle einer Erweiterung auf ein
durch γ diskontiertes Gütefunktional die Gesamtkosten endlich bleiben.

Analoges gilt ausgehend von $\mathbf{x}(t)$ und $\zeta(t)$ für den Systemzustand $\mathbf{x}(\tau)$, $\forall \tau \geq t$, für den durch

$$\dot{\mathbf{x}}(\tau) = \mathbf{f}(\mathbf{x}(\tau)) + \mathbf{g}(\mathbf{x}(\tau))\boldsymbol{\mu}(\mathbf{x}(\tau), \mathbf{f}_{\mathbf{x},\zeta}(\zeta(\tau))) \quad (4.5)$$

geschlossenen Regelkreis. Daraus folgt, dass eine explizite Abhängigkeit der Value Function (4.4) von $\mathbf{x}(t)$ und $\zeta(t)$ zur Beschreibung der Gesamtkosten geeignet ist. \square

Eine ADP-kompatible Solltrajektorien-darstellung stellt somit sicher, dass die Gesamtkosten-repräsentation in Form einer Value Function lediglich eine explizite Abhängigkeit von $\mathbf{x}(t)$ und $\zeta(t)$ aufweist. Die in den Abschnitten 2.2.2 und 2.2.3 diskutierten zeitkontinuierlichen ADP-Trajektorienfolgeregelungsmethoden erfüllen zwar die Kompatibilitätsforderung nach Definition 4.1⁷⁵, ermöglichen jedoch keine externe Vorgabe eines Solltrajektorienverlaufs. Nachfolgend wird eine nach Definition 4.1 kompatible Solltrajektorien-darstellung vorgestellt, die zudem den zukünftigen Referenzzustandsverlauf berücksichtigt.

4.2 Zeitkontinuierliche ADP-kompatible parametrisierte Referenztrajektorie

Im nächsten Unterabschnitt wird zunächst eine ADP-kompatible Solltrajektorien-darstellung präsentiert. Der Einfluss eines global diskontierten Gütemaßes auf den resultierenden Optimal-regler wird anschließend in Abschnitt 4.2.2 untersucht. Da es hierbei potenziell zu Instabilität des geschlossenen Regelkreises kommen kann, wird in Abschnitt 4.2.3 ein Optimierungsproblem mit gedämpfter Referenzdynamik als Alternative zur Verwendung global diskontierter Gütemaße vorgestellt und analysiert. Für dieses Optimierungsproblem wird in Abschnitt 4.2.4 ein ADP-basierter Regelungsansatz entworfen. Schließlich werden in Abschnitt 4.2.5 Simulationsergebnisse gezeigt, welche die Vorzüge der neuartigen ADP-kompatiblen Solltrajektorien-darstellung verdeutlichen.

4.2.1 Solltrajektorien-darstellung

Betrachtet werde die Sollzustandstrajektorie

$$\mathbf{x}_r(t) = [x_{r,1}(t) \quad x_{r,2}(t) \quad \dots \quad x_{r,n}(t)]^T. \quad (4.6)$$

Da auch in diesem Abschnitt zugunsten der Analysierbarkeit der LQ-Fall betrachtet werden soll, wird für die Elemente $x_{r,i}(t)$, $\forall i = 1, \dots, n$, eine lineare Dynamik angesetzt.

⁷⁵ Bei Methoden, bei denen ein stationärer Sollzustand vorgegeben wird, stellt $\mathbf{f}_{\mathbf{x},\zeta}(\cdot)$ eine Identitätsabbildung dar und es gilt $\mathbf{f}_{\zeta}(\cdot) = \mathbf{0}$.

Konkret wird der Ansatz

$$x_{r,i}(t) = \sum_{j=1}^{J_i} \left(\sum_{k=0}^{\nu_{ij}-1} c_{ijk} t^k \right) e^{\lambda_{ij} t} = \sum_{j=1}^{J_i} x_{r,ij}(t) \quad (4.7)$$

mit den Parametern $c_{ijk} \in \mathbb{C}$ und den Hyperparametern $\lambda_{ij} \in \mathbb{C}$ und $\nu_{ij} \in \mathbb{N}$ gewählt. Hierbei bezeichnet J_i die Anzahl verschiedener Eigenwerte und ν_{ij} die algebraische Vielfachheit des Eigenwerts λ_{ij} . Die Parameter $c_{ijk} \in \mathbb{C}$ sind zunächst frei wählbar. Sie bestimmen den konkreten Verlauf der Sollzustandstrajektorie $x_r(t)$. Somit ermöglichen die Parameter c_{ijk} die direkte Beeinflussung des Solltrajektorienverlaufs von außen (vgl. Abbildung 2.4), im Gegensatz zu der in der Literatur häufig verwendeten globalen Referenztrajektorie ohne äußere Beeinflussungsmöglichkeit nach Abbildung 2.2 und Tabelle 2.2. Verglichen mit den in Abschnitt 2.2.3 erwähnten Methoden, die lediglich den aktuellen Sollzustand oder die aktuelle Abweichung von diesem verwenden (vgl. Abbildung 2.3), wird durch den hier vorgestellten neuen Ansatz außerdem der Verlauf der Solltrajektorie berücksichtigt. Es sei weiterhin erwähnt, dass die aus der Literatur bekannte stationäre Sollzustandsvorgabe einen Spezialfall des neu vorgestellten Ansatzes darstellt (vgl. Abschnitt 4.2.5.2). Der auf (4.6) und (4.7) basierende Ansatz wird nachfolgend analysiert und dessen ADP-Kompatibilität im Sinne von Definition 4.1 gezeigt.

Die Darstellung nach (4.7) entspricht der Summation von J_i Lösungen $x_{r,ij}$ homogener linearer Differenzialgleichungen⁷⁶

$$\underbrace{a_{ij\nu_{ij}}}_{:=1} x_{r,ij}^{(\nu_{ij})}(t) + \dots + a_{ij1} \dot{x}_{r,ij} + a_{ij0} x_{r,ij} = 0 \quad (4.8)$$

mit ν_{ij} -fachem Eigenwert λ_{ij} [MBK19, Kapitel 4.7], deren charakteristische Polynome durch

$$p_{ij}(\lambda) = (\lambda - \lambda_{ij})^{\nu_{ij}} = \sum_{k=0}^{\nu_{ij}} \underbrace{\binom{\nu_{ij}}{k} (-\lambda_{ij})^{\nu_{ij}-k}}_{=: a_{ijk}} \lambda^k \quad (4.9)$$

(vgl. Binomischer Lehrsatz [AE10, Satz 8.4]) gegeben sind. Während die Parameter c_{ijk} , die den Solltrajektorienverlauf beschreiben, zur Laufzeit variiert werden können, werden die Hyperparameter λ_{ij} und ν_{ij} vorab festgelegt⁷⁷.

Unter der folgenden Annahme ist die Reellwertigkeit des Solltrajektorienverlaufs $x_r(t)$ gewährleistet.

⁷⁶ Die Notation $x_{r,ij}^{(k)}(t)$ bezeichnet die k -te zeitliche Ableitung von $x_{r,ij}(t)$.

⁷⁷ Wie sich später in Bemerkung 4.3 zeigen wird, stellt dies ADP-Kompatibilität gemäß Definition 4.1 sicher.

Annahme 4.1

Für alle reellwertigen Eigenwerte $\lambda_{ij} \in \mathbb{R}$ sei $c_{ijk} \in \mathbb{R} \forall k \in \{0, \dots, \nu_{ij} - 1\}$. Zudem gelte für jeden komplexwertigen Eigenwert $\lambda_{ij} \in \mathbb{C}$ und den dazu konjugiert-komplexen Eigenwert λ_{ij}^* , gekennzeichnet durch den Index j^* (d. h. $\lambda_{ij}^* =: \lambda_{ij^*}$, $\nu_{ij^*} = \nu_{ij}$ und $c_{ij^*k} = c_{ijk}^*$).

Lemma 4.1

Unter Annahme 4.1 folgt für Sollzustandstrajektorien $\mathbf{x}_r(t)$ nach (4.6) und (4.7) Reellwertigkeit, d. h. $\mathbf{x}_r(t) \in \mathbb{R}^n$.

Beweis:

Siehe Anhang B.1. □

Bemerkung 4.1

Für mehrfache Eigenwerte λ_{ij} bei null entspricht der Sollzustandsverlauf $\mathbf{x}_r(t)$ einer Polynomfunktion und reelle Eigenwerte ungleich null resultieren in Exponentialfunktionen. Aus (4.7) folgt für reelle Eigenwerte λ_{ij} mit $c_{ijk} \in \mathbb{R} (\forall k \in \{0, \dots, \nu_{ij} - 1\})$ nach Annahme 4.1 für jedes der J_i Summenglieder

$$x_{r,ij}(t) = \left(c_{ij0} + c_{ij1}t + \dots + c_{ij(\nu_{ij}-1)}t^{(\nu_{ij}-1)} \right) e^{\lambda_{ij}t}. \quad (4.10)$$

Konjugiert-komplexe Polpaare nach Annahme 4.1 führen zu harmonischen Funktionen, deren Amplitude und Phase durch $c_{ij0} \in \mathbb{C}$ eingestellt werden können, da für $\lambda_{ij} = \lambda_R + j\lambda_I = \lambda_{ij}^*$ und $c_{ij0} = c_R + jc_I = c_{ij0}^*$ gilt:

$$\begin{aligned} x_{r,ij}(t) + x_{r,ij^*}(t) &= (c_R + jc_I) e^{(\lambda_R + j\lambda_I)t} + (c_R - jc_I) e^{(\lambda_R - j\lambda_I)t} \\ &= c_R e^{\lambda_R t} (e^{j\lambda_I t} + e^{-j\lambda_I t}) + jc_I e^{\lambda_R t} (e^{j\lambda_I t} - e^{-j\lambda_I t}) \\ &= 2e^{\lambda_R t} (c_R \cos(\lambda_I t) - c_I \sin(\lambda_I t)) \\ &\stackrel{[\text{Apo67, S. 334}]}{=} 2e^{\lambda_R t} \operatorname{sgn}(c_R) \sqrt{c_R^2 + c_I^2} \cos \left(\lambda_I t + \arctan \left(\frac{c_I}{c_R} \right) \right). \end{aligned} \quad (4.11)$$

Mit

$$\zeta_{ij}(t) := \left[x_{r,ij}(t) \quad \dot{x}_{r,ij}(t) \quad \dots \quad x_{r,ij}^{(\nu_{ij}-1)}(t) \right]^T \quad (4.12)$$

lässt sich die gewöhnliche Differenzialgleichung (4.8) durch

$$\dot{\zeta}_{ij}(t) = \mathbf{D}_{ij} \zeta_{ij}(t) \quad (4.13)$$

darstellen, wobei

$$D_{ij} := \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \\ -a_{ij0} & -a_{ij1} & -a_{ij2} & \dots & -a_{ij(\nu_{ij}-1)} \end{bmatrix} \quad (4.14)$$

gilt. Werden die Vektoren $\zeta_{ij}(t)$ zu

$$\zeta_i(t) := [\zeta_{i1}^T(t) \quad \dots \quad \zeta_{iJ_i}^T(t)]^T \quad (4.15)$$

zusammengefasst, folgt

$$\mathbf{x}_r(t) = \underbrace{\begin{bmatrix} \mathbf{s}_1^T & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{s}_2^T & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{s}_n^T \end{bmatrix}}_{=: \mathbf{S}} \underbrace{\begin{bmatrix} \zeta_1(t) \\ \zeta_2(t) \\ \vdots \\ \zeta_n(t) \end{bmatrix}}_{=: \zeta(t)}, \quad (4.16)$$

wobei

$$\mathbf{s}_i^T := [\mathbf{s}_{i1}^T \quad \dots \quad \mathbf{s}_{iJ_i}^T] \quad (4.17)$$

und

$$\mathbf{s}_{ij} := [1 \quad 0 \quad \dots \quad 0]^T \in \mathbb{R}^{\nu_{ij}}, \forall i, j, \quad (4.18)$$

gilt. Zudem folgt

$$\dot{\zeta}(t) = \mathbf{D}\zeta(t) \quad (4.19)$$

mit

$$\mathbf{D} := \text{diag}(\mathbf{D}_{11}, \dots, \mathbf{D}_{1J_1}, \dots, \mathbf{D}_{nJ_n}). \quad (4.20)$$

Für feste Hyperparameter λ_{ij} und ν_{ij} ist der Verlauf der Solltrajektorie $\mathbf{x}_r(t)$ durch

$$\mathbf{c} := [\mathbf{c}_{11}^T \quad \dots \quad \mathbf{c}_{1J_1}^T \quad \dots \quad \mathbf{c}_{nJ_n}^T]^T \quad (4.21)$$

mit

$$\mathbf{c}_{ij} := [c_{ij0} \quad c_{ij1} \quad \dots \quad c_{ij(\nu_{ij}-1)}]^T \quad (4.22)$$

parametriert.

Im Folgenden wird gezeigt, dass eine Bijektion zwischen \mathbf{c} und $\zeta(0)$ existiert, die Parametrierung der Solltrajektorie also eindeutig durch $\zeta(0)$ codiert ist. Hierzu wird zunächst das folgende Hilfslemma benötigt.

Lemma 4.2

Die l -te zeitliche Ableitung von $x_{r,ij}(t)$ ist durch

$$x_{r,ij}^{(l)}(t) = \sum_{m=0}^l e^{\lambda_{ij}t} \binom{l}{m} \lambda_{ij}^{l-m} \sum_{k=m}^{\nu_{ij}-1} \frac{k!}{(k-m)!} c_{ijk} t^{k-m} \quad (4.23)$$

gegeben.

Beweis:

Siehe Anhang B.2. □

Proposition 4.2

Es existiert eine bijektive Abbildung zwischen c und $\zeta(0)$ der Form

$$\zeta(0) = \Gamma c \quad (4.24)$$

mit $\Gamma \in \mathbb{R}^{n_\zeta \times n_c}$.

Beweis:

Aus Lemma 4.2 folgt für $t = 0$

$$x_{r,ij}^{(l)}(0) = \sum_{m=0}^l \binom{l}{m} \lambda_{ij}^{l-m} m! c_{ijm} = \sum_{m=0}^l \frac{l!}{(l-m)!} \lambda_{ij}^{l-m} c_{ijm}. \quad (4.25)$$

Aus (4.12) ergibt sich der Initialzustand $\zeta_{ij}(0)$ zu

$$\zeta_{ij}(0) = \Gamma_{ij} c_{ij}, \quad (4.26)$$

wobei

$$(\Gamma_{ij})_{lm} := \begin{cases} \frac{(l-1)!}{(l-m)!} \lambda_{ij}^{l-m}, & m \leq l \\ 0, & m > l \end{cases}, \quad l, m = 1, \dots, \nu_{ij}, \quad (4.27)$$

gilt. Mithilfe des Laplaceschen Entwicklungssatzes [BHW12, Satz 3.21] folgt, da es sich bei Γ_{ij} um eine untere Dreiecksmatrix handelt,

$$\det(\Gamma_{ij}) = \prod_{l=1}^{\nu_{ij}} (\Gamma_{ij})_{ll} = \prod_{l=0}^{\nu_{ij}-1} l! > 0. \quad (4.28)$$

Mit

$$\Gamma := \text{diag}(\Gamma_{11}, \dots, \Gamma_{1J_1}, \dots, \Gamma_{nJ_n}) \quad (4.29)$$

ergibt sich nach [Pow11]

$$\det(\Gamma) = \prod_{i=1}^n \prod_{j=1}^{J_i} \Gamma_{ij} > 0. \quad (4.30)$$

Somit ist Γ regulär und stellt nach [Axl15, Theorem 10.24] eine Bijektion dar. \square

Die nachfolgende Bemerkung begründet, weshalb $\zeta(t)$ in (4.16) und D in (4.20) als reellwertig angenommen werden können.

Bemerkung 4.2

Zwar stellt Lemma 4.1 $\mathbf{x}_r(t) \in \mathbb{R}^n$ sicher, dennoch können durch konjugiert-komplexe Polpaare λ_{ij} und $\lambda_{ij}^ = \lambda_{ij}^*$ komplexwertige Zustände ζ_{ij} und ζ_{ij}^* sowie komplexwertige Einträge in den Elementen von D_{ij} bzw. D_{ij}^* (vgl. (4.14)) resultieren. Durch eine eindeutige Zustandstransformation lassen sich jedoch stets reellwertige Ersatzzustände mit einer reellen Dynamik erzeugen. Dieses Vorgehen ist in Anhang B.3 skizziert. Daher wird o. B. d. A. stets $\zeta(t)$ in (4.16) und D in (4.20) als reellwertig angenommen.*

Schließlich sei zu betonen, dass die soeben vorgestellte Referenztrajektorienendarstellung ADP-kompatibel ist. Dies ist in der nachfolgenden Bemerkung zusammengefasst.

Bemerkung 4.3

Die in diesem Abschnitt vorgestellte Referenztrajektorienendarstellung ist ADP-kompatibel im Sinne von Definition 4.1 mit

$$\mathbf{f}_{\mathbf{x}_r, \zeta}(\zeta(t)) := S\zeta(t) \quad (4.31)$$

(vgl. (4.1)) und

$$\mathbf{f}_{\zeta}(\zeta(t)) := D\zeta(t) \quad (4.32)$$

(vgl. (4.2)). Die Wahl der Hyperparameter λ_{ij} und ν_{ij} definiert die Klasse der Referenztrajektorienapproximation (vgl. Bemerkung 4.1). Zur Laufzeit kann der konkrete Solltrajektorienverlauf durch Vorgabe des Parameters \mathbf{c} schließlich von außen beeinflusst werden.

4.2.2 Trajektorienfolgeregelung mit global diskontiertem Gütemaß

Bevor ein ADP-basierter Solltrajektorienfolgeregler entworfen wird, soll untersucht werden, unter welchen Bedingungen und für welche Gütemaße unter Verwendung der in Abschnitt 4.2.1 vorgestellten ADP-kompatiblen Referenztrajektorienendarstellung stabilisierende optimale Lösungen existieren.

In Anlehnung an die Eingangs-Zustands-Stabilität [Kha02, Definition 4.7] und analog zur zeitdiskreten Referenz-Zustands-Stabilität (vgl. Definition 3.5) wird hierfür zunächst der Begriff der zeitkontinuierlichen Referenz-Zustands-Stabilität definiert.

Definition 4.2 (Zeitkontinuierliche Referenz-Zustands-Stabilität)

Ein geregeltes System ist Referenz-Zustands-stabil, wenn für endliche Solltrajektorienparameter

$$\|\mathbf{c}\|_2 < \infty$$

beschränkte Systemzustände

$$\|\mathbf{x}(t)\|_2 < \infty, \forall t \geq 0,$$

resultieren.

Bei der Wahl eines geeigneten Gütemaßes, das einerseits die Abweichung des Systemzustands $\mathbf{x}(t)$ vom Sollzustand $\mathbf{x}_r(t)$ und andererseits die durch die Stellgröße $\mathbf{u}(t)$ aufgebrachte Stellenergie bestraft, ist ein schlecht gestelltes, d. h. mit unendlichen Gesamtkosten verbundenes, Problem zu vermeiden (vgl. [XZLJ16]). Unendlich hohe Gesamtkosten entstünden beispielsweise dann, wenn das Halten des Systemzustands in einem Sollzustand mit Stellenergie verbunden wäre, und Stellenergie sowie die Abweichung vom Sollzustand im Gütemaß mit unendlichem Optimierungshorizont bestraft würden. Daher ist die Verwendung eines globalen Diskontierungsfaktors γ , wobei für $\gamma > 0$ in der Zukunft liegende Kosten gedämpft werden, ein in der RL- und ADP-Literatur häufig verwendeter Ansatz (vgl. Tabelle 2.2). Für

$$\dot{\tilde{\mathbf{x}}}(t) = \frac{d}{dt} \begin{bmatrix} \mathbf{x}(t) \\ \boldsymbol{\zeta}(t) \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{x}(t) \\ \boldsymbol{\zeta}(t) \end{bmatrix} + \begin{bmatrix} \mathbf{B} \\ \mathbf{0} \end{bmatrix} \mathbf{u}(t) =: \tilde{\mathbf{A}}\tilde{\mathbf{x}}(t) + \tilde{\mathbf{B}}\mathbf{u}(t) \quad (4.33)$$

werde zunächst das folgende Problem betrachtet⁷⁸.

⁷⁸ Grundsätzlich lässt sich die in Abschnitt 4.2 vorgestellte ADP-kompatible Solltrajektorienrepräsentation auch auf nichtlineare Systeme anwenden. Für eingangsaffine Systeme ergibt sich beispielsweise die erweiterte Systemdynamik $\dot{\tilde{\mathbf{x}}}(t) := \begin{bmatrix} \mathbf{f}(\mathbf{x}(t)) \\ \mathbf{D}\boldsymbol{\zeta}(t) \end{bmatrix} + \begin{bmatrix} \mathbf{g}(\mathbf{x}(t)) \\ \mathbf{0} \end{bmatrix} \mathbf{u}(t) =: \tilde{\mathbf{f}}(\tilde{\mathbf{x}}(t)) + \tilde{\mathbf{g}}(\tilde{\mathbf{x}}(t))\mathbf{u}(t)$ sowie die Value Function $V^\mu(\mathbf{x}, \boldsymbol{\zeta}) = \int_t^\infty e^{-\gamma(\tau-t)} (q(\mathbf{x}, \boldsymbol{\zeta}) + \boldsymbol{\mu}^\top \mathbf{R}\boldsymbol{\mu}) d\tau$. Jedoch ist die Form der Value Function V^μ in diesem Fall im Allgemeinen unbekannt und die Wahl geeigneter Basisfunktionen zur Funktionsapproximation (vgl. Abschnitt 2.1.3) nach wie vor ungelöst [WHL17]. Zwar erlaubt beispielsweise der Satz von Stone-Weierstraß Aussagen zur Approximierbarkeit stetiger Funktionen auf einer kompakten Menge durch genügend viele Basisfunktionen (vgl. [Che78, S. 226], [Wei85]), jedoch kann die Anzahl benötigter Basisfunktionen gegebenenfalls hoch sein und die lediglich lokale Gültigkeit muss berücksichtigt werden.

Problem 4.1 (Globale Diskontierung)

$$\min_{\mu(\mathbf{x}(\tau), \zeta(\tau))} \int_t^{\infty} e^{-\gamma(\tau-t)} \left(\tilde{\mathbf{x}}^\top(\tau) \tilde{\mathbf{Q}} \tilde{\mathbf{x}}(\tau) + \boldsymbol{\mu}^\top(\mathbf{x}(\tau), \zeta(\tau)) \mathbf{R} \boldsymbol{\mu}(\mathbf{x}(\tau), \zeta(\tau)) \right) d\tau \quad (4.34)$$

$$u. d. N. \dot{\tilde{\mathbf{x}}}(\tau) = \tilde{\mathbf{A}} \tilde{\mathbf{x}}(\tau) + \tilde{\mathbf{B}} \boldsymbol{\mu}(\mathbf{x}(\tau), \zeta(\tau)) \quad (4.35)$$

mit

$$\tilde{\mathbf{Q}} := [\mathbf{I} \quad -\mathbf{S}]^\top \mathbf{Q} [\mathbf{I} \quad -\mathbf{S}]. \quad (4.36)$$

Mit (4.16) ergibt sich

$$\begin{aligned} \tilde{\mathbf{x}}^\top(\tau) \tilde{\mathbf{Q}} \tilde{\mathbf{x}}(\tau) &= [\mathbf{x}(\tau) \quad \zeta(\tau)] \begin{bmatrix} \mathbf{I} \\ -\mathbf{S} \end{bmatrix} \mathbf{Q} [\mathbf{I} \quad -\mathbf{S}] \begin{bmatrix} \mathbf{x}(\tau) \\ \zeta(\tau) \end{bmatrix} \\ &= (\mathbf{x}(\tau) - \mathbf{x}_r(\tau))^\top \mathbf{Q} (\mathbf{x}(\tau) - \mathbf{x}_r(\tau)), \end{aligned} \quad (4.37)$$

d. h. in (4.34) wird die Abweichung des Zustands $\mathbf{x}(\tau)$ von $\mathbf{x}_r(\tau)$ bestraft. Problem 4.1 kann in das folgende äquivalente Ersatzproblem überführt werden. Die Äquivalenz ist in Anhang B.4 gezeigt.

Problem 4.2 (Äquivalente Formulierung zur globalen Diskontierung)

$$\min_{\mu(\mathbf{x}(\tau), \zeta(\tau))} \int_t^{\infty} \left(\tilde{\mathbf{x}}^\top(\tau) \tilde{\mathbf{Q}} \tilde{\mathbf{x}}(\tau) + \boldsymbol{\mu}^\top(\mathbf{x}(\tau), \zeta(\tau)) \mathbf{R} \boldsymbol{\mu}(\mathbf{x}(\tau), \zeta(\tau)) \right) d\tau \quad (4.38)$$

$$u. d. N. \dot{\tilde{\mathbf{x}}}(\tau) = \left(\tilde{\mathbf{A}} - \frac{\gamma}{2} \mathbf{I} \right) \tilde{\mathbf{x}}(\tau) + \tilde{\mathbf{B}} \boldsymbol{\mu}(\mathbf{x}(\tau), \zeta(\tau)) \quad (4.39)$$

mit $\tilde{\mathbf{Q}}$ wie in (4.36).

Somit führt die Verwendung eines globalen Diskontierungsfaktors γ zum gleichen Effekt, als würden bei einem nicht-diskontierten Optimierungsproblem alle Eigenwerte der erweiterten Systemmatrix $\tilde{\mathbf{A}}$ um $\frac{\gamma}{2}$ nach links verschoben. Diese Verschiebung der Eigenwerte findet jedoch nur virtuell statt, d. h. aus Sicht des zu minimierenden global diskontierten Gütefunctionals wirken die Systemeigenwerte verschoben, das reale System bleibt jedoch nach wie vor durch (4.33) gegeben. Dies kann potenziell dazu führen, dass eine für das Ersatzproblem 4.2 stabilisierende und optimale Lösung

$$\boldsymbol{\mu}^*(\mathbf{x}(t), \zeta(t)) = -\tilde{\mathbf{K}}^* \tilde{\mathbf{x}}(t) = -[\mathbf{K}_x^* \quad \mathbf{K}_\zeta^*] \begin{bmatrix} \mathbf{x}(t) \\ \zeta(t) \end{bmatrix} \quad (4.40)$$

gefunden wird, die zwar aufgrund der Äquivalenz zu Problem 4.1 auch (4.34) minimiert, jedoch mit einem instabilen Systemverhalten des Systems (4.35) und damit auch des Systems

$$\dot{\mathbf{x}}(t) = (\mathbf{A} - \mathbf{BK}_x^*) \mathbf{x}(t) \quad (4.41)$$

einhergeht. Dies wird anhand des nachfolgenden Beispiels illustriert.

Beispiel 4.1 *Betrachtet werde das stabilisierbare System zweiter Ordnung mit der Systemdynamik*

$$\dot{\mathbf{x}}(t) = \begin{bmatrix} 0 & 1 \\ -1,01 & 0,2 \end{bmatrix} \mathbf{x}(t) + \begin{bmatrix} 1 \\ 0 \end{bmatrix} u(t). \quad (4.42)$$

Für eine ADP-kompatible Solltrajektorienarstellung nach (4.7) werde ein einzelner Eigenwert bei $\lambda_{11} = 0$ gewählt, wodurch die Vorgabe eines konstanten Sollzustands für den Zustand $x_1(t)$ ermöglicht wird. Dies führt nach (4.39) auf die effektive erweiterte Systemdynamik

$$\dot{\tilde{\mathbf{x}}}(t) = \underbrace{\begin{bmatrix} -\frac{\gamma}{2} & 1 & 0 \\ -1,01 & 0,2 - \frac{\gamma}{2} & 0 \\ 0 & 0 & -\frac{\gamma}{2} \end{bmatrix}}_{=(\tilde{\mathbf{A}} - \frac{\gamma}{2} \mathbf{I})} \tilde{\mathbf{x}}(t) + \underbrace{\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}}_{\tilde{\mathbf{B}}} u(t). \quad (4.43)$$

Mit

$$\mathbf{Q} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{S} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \mathbf{R} = 100 \quad \text{und} \quad \gamma = 1 \quad (4.44)$$

existiert eine optimale Lösung für Problem 4.2 und das optimale Regelgesetz ergibt sich zu

$$\boldsymbol{\mu}(\tilde{\mathbf{x}}(t)) = - \underbrace{\begin{bmatrix} 0,0067 & 0,0016 & -0,0044 \end{bmatrix}}_{=\tilde{\mathbf{K}}^*} \tilde{\mathbf{x}}(t). \quad (4.45)$$

Die Systempole des mittels (4.45) geschlossenen Regelkreises, d. h. die Eigenwerte von $\tilde{\mathbf{A}} - \tilde{\mathbf{B}}\tilde{\mathbf{K}}^*$ ergeben sich zu $0,0966 \pm j0,9988; 0$. Dies führt selbst bei der Vorgabe $x_{r,1}(t) = 0$ zu einem instabilen System, wie in Abbildung 4.1 für $\tilde{\mathbf{x}}(0) = [1 \ 1 \ 0]^T$ veranschaulicht ist. Die mögliche Instabilität ist somit auf die globale Diskontierung und nicht auf die Solltrajektorienarstellung zurückzuführen. Insbesondere ist das System für dieses Beispiel auch nicht Referenz-Zustands-stabil nach Definition 4.2.

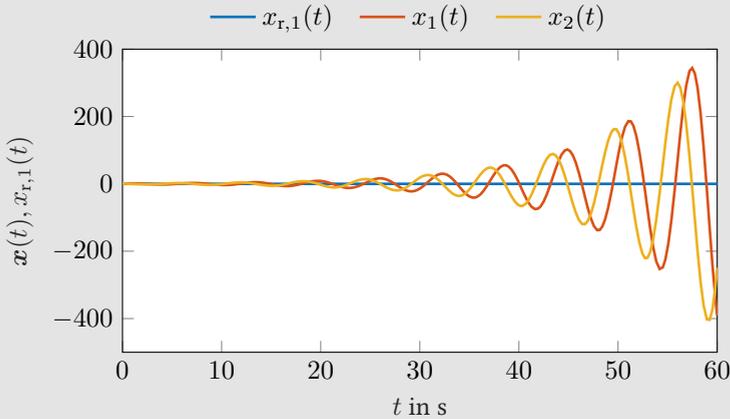


Abbildung 4.1: Verlauf der Systemtrajektorie für das betrachtete Beispielsystem bei globaler Diskontierung.

Da durch die Verwendung einer globalen Diskontierung $\gamma > 0$ und die damit verbundene Dämpfung der Kosten in (4.38), wie in Beispiel 4.2 ersichtlich, ein instabiles Regelgesetz resultieren kann, wird im nächsten Abschnitt die Verwendung eines Diskontierungsfaktors, der lediglich die Dynamik D dämpft, vorgestellt.

4.2.3 Trajektorienfolge mit gedämpfter Referenzdynamik

Um trotz einer gegebenenfalls nicht asymptotisch stabilen Solltrajektorien dynamik D nach (4.19) und ohne Verwendung einer globalen Diskontierung wie in Problem 4.1 ein Optimierungsproblem mit endlichen Kosten zu definieren, wird die Verwendung einer teilweisen Dämpfung der erweiterten Systemdynamik gemäß der nachfolgenden Problemstellung betrachtet.

Problem 4.3 (Teilweise Dämpfung)

$$\min_{\mu(x(\tau), \zeta(\tau))} \int_t^{\infty} \left(\tilde{x}^\top(\tau) \tilde{Q} \tilde{x}(\tau) + \mu^\top(x(\tau), \zeta(\tau)) R \mu(x(\tau), \zeta(\tau)) \right) d\tau \quad (4.46)$$

$$\text{u. d. N. } \dot{\tilde{x}}(\tau) = \tilde{A}' \tilde{x}(\tau) + \tilde{B} \mu(x(\tau), \zeta(\tau)) \quad (4.47)$$

mit

$$\tilde{A}' := \begin{bmatrix} A & \mathbf{0} \\ \mathbf{0} & D - \frac{\gamma}{2} I \end{bmatrix} =: \begin{bmatrix} A & \mathbf{0} \\ \mathbf{0} & D' \end{bmatrix} \quad (4.48)$$

und \tilde{Q} wie in (4.36).

Bemerkung 4.4

Für den mit der gedämpften Dynamik $D' = D - \frac{\gamma}{2}I$ propagierten Trajektorienparameter $\zeta(t)$ folgt mit

$$f_{\zeta}(\zeta(t)) := D'\zeta(t) \quad (4.49)$$

analog zu Bemerkung 4.3 ADP-Kompatibilität.

Zudem werden die folgenden Standardannahmen für LQ-Optimierungsprobleme (vgl. [AM89, Abschnitt 3.2], [LVS12, Abschnitt 3.4], [Kuč73]) vorausgesetzt.

Annahme 4.2

1. (A, B) sei stabilisierbar.
2. Seien $Q \succeq 0$ und $R \succ 0$.
3. (A, \sqrt{Q}) sei detektierbar, wobei $\sqrt{Q}^T \sqrt{Q} = Q$ gilt.

Im Folgenden wird die Stabilität des geschlossenen Regelkreises für die nach Problem 4.3 optimale Lösung untersucht. Bevor durch Satz 4.1 die Hauptaussage des vorliegenden Abschnitts folgt, werden zunächst einige hierfür benötigte Zusammenhänge gegeben.

Lemma 4.3

Seien

$$M_1 \in \mathbb{R}^{n_1 \times m_1}, \quad n_1 \leq m_1, \quad (4.50)$$

$$M_2 \in \mathbb{R}^{n_2 \times m_2}, \quad n_2 \leq m_2 \quad (4.51)$$

Matrizen mit Maximalrang, d. h. es gelte

$$\text{Rang}(M_1) = \text{Dim}(\text{Bild}(M_1)) = n_1, \quad (4.52)$$

$$\text{Rang}(M_2) = \text{Dim}(\text{Bild}(M_2)) = n_2. \quad (4.53)$$

Dann weist

$$M = \begin{bmatrix} M_1 & M_3 \\ \mathbf{0} & M_2 \end{bmatrix} \in \mathbb{R}^{(n_1+n_2) \times (m_1+m_2)} \quad (4.54)$$

für beliebige $M_3 \in \mathbb{R}^{n_1 \times m_2}$ den Maximalrang

$$\text{Rang}(M) = n_1 + n_2 \quad (4.55)$$

auf.

Beweis:

Der Beweis ist in Anhang B.5 gegeben. □

Im weiteren Verlauf sei

$$\gamma_{\min} := 2 \max(0, \lambda_{\text{Re}}^+(D)), \quad (4.56)$$

wobei $\lambda_{\text{Re}}^+(D)$ dem Eigenwert von D mit größtem Realteil⁷⁹ entspricht. Damit lässt sich das folgende Lemma formulieren.

Lemma 4.4

Unter Annahme 4.2 ist das durch (\tilde{A}', \tilde{B}) beschriebene System mit \tilde{A}' wie in (4.48) und $\tilde{B} = \begin{bmatrix} B \\ 0 \end{bmatrix}$ stabilisierbar, falls $\gamma > \gamma_{\min}$ gilt.

Beweis:

Siehe Anhang B.6. □

Lemma 4.5

Seien M_1, M_2 wie in (4.50)–(4.53). Dann weist

$$M = \begin{bmatrix} M_1 & \mathbf{0} \\ M_3 & M_2 \end{bmatrix} \in \mathbb{R}^{(n_1+n_2) \times (m_1+m_2)} \quad (4.57)$$

für beliebige $M_3 \in \mathbb{R}^{n_2 \times m_1}$ den Maximalrang

$$\text{Rang}(M) = n_1 + n_2 \quad (4.58)$$

auf.

Beweis:

Der Beweis erfolgt analog zum Beweis von Lemma 4.3. □

Lemma 4.6

Unter Annahme 4.2 ist $(\tilde{A}', \sqrt{\tilde{Q}})$ detektierbar, falls $\gamma > \gamma_{\min}$ gilt. Hierbei sei \tilde{A}' wie in (4.48) und $\sqrt{\tilde{Q}}^T \sqrt{\tilde{Q}} = \tilde{Q}$ mit \tilde{Q} aus (4.36).

⁷⁹ Nach Bemerkung 4.2 ist D reellwertig.

Beweis:Siehe Anhang B.7. □**Satz 4.1**

Sei Annahme 4.2 erfüllt, zudem gelte $\gamma > \gamma_{\min} = 2 \max(0, \lambda_{\text{Re}}^+(\mathbf{D}))$. Dann existiert eine eindeutige, symmetrische, positiv semidefinite Lösung $\tilde{\mathbf{P}}^*$ der algebraischen Riccati-Gleichung

$$\tilde{\mathbf{Q}} + \tilde{\mathbf{P}}^* \tilde{\mathbf{A}}' + \tilde{\mathbf{A}}'^{\top} \tilde{\mathbf{P}}^* - \tilde{\mathbf{P}}^* \tilde{\mathbf{B}} \mathbf{R}^{-1} \tilde{\mathbf{B}}^{\top} \tilde{\mathbf{P}}^* = \mathbf{0}. \quad (4.59)$$

Zudem ist

$$\boldsymbol{\mu}^*(\mathbf{x}(t), \zeta(t)) = -\tilde{\mathbf{K}}^* \tilde{\mathbf{x}}(t) = -\begin{bmatrix} \mathbf{K}_x^* & \mathbf{K}_\zeta^* \end{bmatrix} \begin{bmatrix} \mathbf{x}(t) \\ \zeta(t) \end{bmatrix} \quad (4.60)$$

mit

$$\tilde{\mathbf{K}}^* = -\mathbf{R}^{-1} \tilde{\mathbf{B}}^{\top} \tilde{\mathbf{P}}^* \quad (4.61)$$

optimal bezüglich Problem 4.3 und

$$\dot{\tilde{\mathbf{x}}}(t) = (\tilde{\mathbf{A}}' - \tilde{\mathbf{B}} \tilde{\mathbf{K}}^*) \tilde{\mathbf{x}}(t) \quad (4.62)$$

und

$$\dot{\mathbf{x}}(t) = (\mathbf{A} - \mathbf{B} \mathbf{K}_x^*) \mathbf{x}(t) \quad (4.63)$$

sind global asymptotisch stabil.

Beweis:Nach Annahme 4.2 gilt $\mathbf{R} \succ \mathbf{0}$ und $\mathbf{Q} \succeq \mathbf{0}$. Somit folgt aufgrund von

$$\tilde{\mathbf{x}}^{\top} \tilde{\mathbf{Q}} \tilde{\mathbf{x}} = (\begin{bmatrix} \mathbf{I} & -\mathbf{S} \end{bmatrix} \tilde{\mathbf{x}})^{\top} \mathbf{Q} (\begin{bmatrix} \mathbf{I} & -\mathbf{S} \end{bmatrix} \tilde{\mathbf{x}}) \geq 0, \quad \forall \tilde{\mathbf{x}}, \quad (4.64)$$

direkt $\tilde{\mathbf{Q}} \succeq \mathbf{0}$. Zudem ist $(\tilde{\mathbf{A}}', \tilde{\mathbf{B}})$ stabilisierbar (Lemma 4.4) und $(\tilde{\mathbf{A}}', \sqrt{\tilde{\mathbf{Q}}})$ detektierbar(Lemma 4.6). Existenz und Eindeutigkeit von $\tilde{\mathbf{P}}^*$ folgen daher direkt aus [AM89, Abschnitt 3.2] und [Kuč73, Theorem 5], ebenso die global asymptotische Stabilität von (4.62).

Schließlich wird global asymptotische Stabilität von (4.63) gezeigt. Wegen

$$\tilde{\mathbf{A}}' - \tilde{\mathbf{B}} \tilde{\mathbf{K}}^* = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}' \end{bmatrix} - \begin{bmatrix} \mathbf{B} \\ \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{K}_x^* & \mathbf{K}_\zeta^* \end{bmatrix} = \begin{bmatrix} \mathbf{A} - \mathbf{B} \mathbf{K}_x^* & -\mathbf{B} \mathbf{K}_\zeta^* \\ \mathbf{0} & \mathbf{D}' \end{bmatrix} \quad (4.65)$$

folgt mit [Pow11]

$$\det(\lambda \mathbf{I} - (\tilde{\mathbf{A}}' - \tilde{\mathbf{B}}\tilde{\mathbf{K}}^*)) = \det(\lambda \mathbf{I} - (\mathbf{A} - \mathbf{B}\mathbf{K}_x^*)) \det(\lambda \mathbf{I} - \mathbf{D}'). \quad (4.66)$$

Eigenwerte von $\mathbf{A} - \mathbf{B}\mathbf{K}_x^*$ sind somit zugleich auch Eigenwerte von $(\tilde{\mathbf{A}}' - \tilde{\mathbf{B}}\tilde{\mathbf{K}}^*)$, weshalb global asymptotische Stabilität von (4.63) resultiert. \square

Korollar 4.1

Das System

$$\dot{\mathbf{x}}(t) = (\mathbf{A} - \mathbf{B}\mathbf{K}_x^*) \mathbf{x}(t) - \mathbf{B}\mathbf{K}_\zeta^* \Gamma \mathbf{c} \quad (4.67)$$

ist Referenz-Zustands-stabil nach Definition 4.2.

Beweis:

Da (4.63) nach Satz 4.1 global asymptotisch stabil ist, ist (4.67), wenn $(\mathbf{A} - \mathbf{B}\mathbf{K}_x^*)$ als Systemmatrix und $-\mathbf{B}\mathbf{K}_\zeta^* \Gamma$ als Eingangsmatrix interpretiert wird, Eingangs-Ausgangs-stabil [Lun20b, S. 59]. Für endliche Solltrajektorienparameter \mathbf{c} folgen somit stets beschränkte Systemzustände und (4.67) ist Referenz-Zustands-stabil. \square

Dies zeigt insbesondere, dass, obwohl durch eine äußere Vorgabe von \mathbf{c} eine Neuinitialisierung von ζ erfolgt (vgl. Proposition 4.2), die Stabilität des Systemzustands $\mathbf{x}(t)$ durch die beliebige Vorgabe endlicher Solltrajektorienparameter \mathbf{c} nicht gefährdet wird. Der in (4.67) verwendete Regler

$$\boldsymbol{\mu}(\mathbf{x}(t), \mathbf{c}) = -\mathbf{K}_x^* \mathbf{x}(t) - \mathbf{K}_\zeta^* \Gamma \mathbf{c} \quad (4.68)$$

liefert dabei die nach Problem 4.3 optimale Stellgröße, wobei durch Wahl von $\Gamma \mathbf{c}$ der Initialwert $\zeta(t)$ in (4.46) eingestellt wird. Mithilfe von \mathbf{c} kann somit jederzeit der lokale Verlauf der Solltrajektorie beeinflusst werden.

Das nachfolgende Beispiel greift das in Beispiel 4.2 bereits betrachtete System erneut auf und veranschaulicht, dass die in Problem 4.3 vorgestellte teilweise Dämpfung zu einem geregelten System führt, welches Referenz-Zustands-stabil ist.

Beispiel 4.2 Betrachtet werde das durch (4.42) gegebene System, zudem sei $\mathbf{D} = 0$ und somit $\mathbf{D}' = -\frac{\gamma}{2} \mathbf{I}$. Mit \mathbf{Q} , \mathbf{R} , \mathbf{S} und γ wie in (4.44) folgt das optimale Regelgesetz

$$\boldsymbol{\mu}(\tilde{\mathbf{x}}(t)) = - \underbrace{\begin{bmatrix} 0,4245 & -0,0843 & -0,0022 \end{bmatrix}}_{=\tilde{\mathbf{K}}^*} \tilde{\mathbf{x}}(t) \quad (4.69)$$

für Problem 4.3. Die Eigenwerte von $\mathbf{A} - \mathbf{B}\mathbf{K}_x^*$ ergeben sich zu $-0,1122 \pm j0,9988$ und Referenz-Zustands-Stabilität folgt nach Korollar 4.1. Für $\tilde{\mathbf{x}}(0) = [1 \ 1 \ 0]^\top$ ist der sich ergebende Zustandsverlauf in Abbildung 4.2 gezeigt.

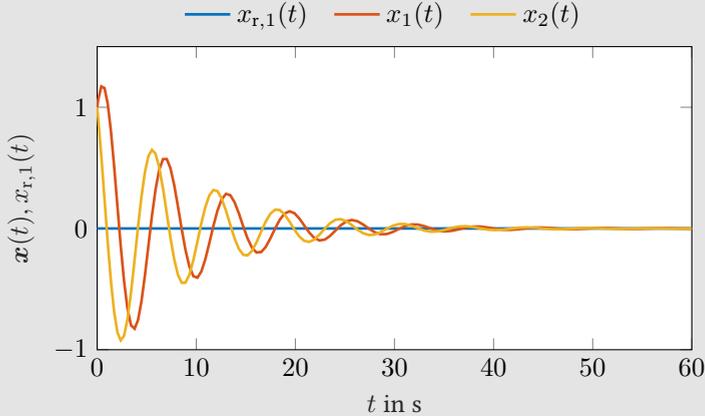


Abbildung 4.2: Verlauf der Systemtrajektorie für das betrachtete Beispielsystem bei teilweiser Dämpfung.

Für den Rest dieses Kapitels wird daher die in Problem 4.3 vorgestellte teilweise Dämpfung der erweiterten Systemdynamik verwendet.

Die nachfolgende Bemerkung zeigt jedoch, dass auch die Nutzung einer global diskontierten Gütefunktion prinzipiell möglich wäre, solange die Diskontierung nicht zu stark ist, d. h. γ nicht zu groß gewählt wird⁸⁰.

Bemerkung 4.5

Eine Möglichkeit, unter Annahme 4.2 stabilisierende Lösungen für LQ-Solltrajektorienfolgeregelungsprobleme zu gewährleisten, ist nach Satz 4.1 durch die Verwendung einer teilweise gedämpften erweiterten Systemdynamik nach Problem 4.3 gegeben.

Eine Alternative dazu ist die Verwendung einer global diskontierten Gütefunktion nach Problem 4.1 unter Beachtung zusätzlicher Nebenbedingungen bei der Wahl des Diskontierungsfaktors γ . Beispielsweise ist in [XZL]16, Theorem 4] durch $\gamma \leq 2 \left\| (\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^\top)^{\frac{1}{2}} \right\|_2$ eine obere Schranke für global diskontierte LQ-Regulierungsprobleme gegeben, sodass eine stabilisierende Lösung existiert. Analog liefert [YWM⁺ 19, Theorem 1] obere Schranken für Zwei-Spieler-LQ-Differenzialspiele. Zur Berechnung dieser Schranken müssten jedoch sowohl die Parameter der jeweiligen Gütefunktionen als auch die jeweiligen Eingangsmatrizen \mathbf{B} bekannt sein. Zu beachten ist weiterhin, dass $\mathbf{x}_r(t)$ nicht zu unendlichen

⁸⁰ Analog hierzu stellt Satz 3.3 im zeitdiskreten Fall Stabilität sicher, sofern dort die Diskontierung nicht zu stark gewählt wird, d. h. γ nicht zu klein gewählt wird.

Gesamtkosten führen darf. In [XZLJ16] wird daher beispielsweise zusätzlich zur oberen Schranke von γ gefordert, dass $\mathbf{x}_r(t)$ beschränkt ist und $\gamma > 0$ gilt⁸¹. Die Verwendung einer globalen Diskontierung ist somit durchaus möglich, sofern die Diskontierung mit Bedacht gewählt wird. Dies ist insbesondere für Algorithmen relevant, deren Konvergenzaussagen von der Verwendung einer globalen Diskontierung profitieren⁸².

Nachdem in diesem Abschnitt eine ADP-kompatible lokale Referenztrajektorienrepräsentation vorgestellt wurde, wird diese Solltrajektorienrepräsentation im nachfolgenden Abschnitt für den Entwurf eines ADP-basierten Solltrajektorienreglers verwendet.

4.2.4 ADP-Umsetzung

In diesem Abschnitt werden die Struktur des ADP-basierten Solltrajektorienfolgereglers, der die in Abschnitt 4.2.1 präsentierte Trajektorienrepräsentation verwendet, und der zugrunde liegende ADP-Algorithmus, der genutzt wird, vorgestellt. Anschließend werden in Abschnitt 4.2.5 Simulationsergebnisse präsentiert.

Um Problem 4.3 basierend auf Messdaten und ohne Kenntnis der Werte von \mathbf{A} und \mathbf{B} lösen zu können, wird eine zeitkontinuierliche Value Iteration (vgl. Abschnitt 2.1.4.2), konkret der Algorithmus von Bian und Jiang [BJ16a], verwendet. Dieser Ansatz benötigt keine explizite Kenntnis eines zulässigen initialen Regelgesetzes und gehört zudem zur Klasse der Off-Policy-Algorithmen inklusive der in Abschnitt 2.1.4.4 genannten Vorteile. Die zeitkontinuierliche Value Iteration [BJ16a] ist der Klasse der IRL-Algorithmen zuzuordnen und basiert auf einer Grenzwertbetrachtung der zeitdiskreten VI (vgl. Abschnitt 2.1.4.2). Nachfolgend stellt t die physikalische Zeit (und letztlich auch die Zeit der Datenaufzeichnung) dar, wohingegen s die Zeitvariable des Adaptionsvorgangs beschreibt (dies entspricht dem Iterationsindex l der zeitdiskreten VI). Das Ziel der Methode ist somit die Konvergenz

$$\lim_{s \rightarrow \infty} V^{\mu_s}(\tilde{\mathbf{x}}(t)) = V^*(\tilde{\mathbf{x}}(t)). \quad (4.70)$$

Aus (2.20b) folgt

$$V^{\mu_{s+T_{\text{IRL}}}}(\tilde{\mathbf{x}}(t)) = \min_{\boldsymbol{\mu}} \left(\int_t^{t+T_{\text{IRL}}} r(\tilde{\mathbf{x}}, \boldsymbol{\mu}) \, d\tau + V^{\mu_s}(\tilde{\mathbf{x}}(t + T_{\text{IRL}})) \right) \quad (4.71)$$

⁸¹ Für eine instabile Exosystemdynamik würde sich, aufgrund der Nicht-Steuerbarkeit von $\boldsymbol{\zeta}(t)$, auch bei Verwendung einer globalen Diskontierung zusätzlich zu einer oberen Schranke für γ (vgl. [XZLJ16, Theorem 4]) die untere Schranke γ_{\min} (4.56) ergeben.

⁸² Siehe beispielsweise [LP03, Theorem 7.1] für Konvergenzaussagen des zeitdiskreten LSPI-Ansatzes. In [HWL21, Theorem 2] wird zudem eine untere Schranke für die Wahl eines Diskontierungsfaktors eines global diskontierten zeitkontinuierlichen Optimierungsproblems gegeben, um Stabilität zu gewährleisten.

(vgl. (2.24)) und somit

$$\begin{aligned} \frac{\partial V^{\mu_s}(\tilde{\mathbf{x}}(t))}{\partial s} &= \min_{\boldsymbol{\mu}} \left(r(\tilde{\mathbf{x}}(t), \boldsymbol{\mu}) + \left(\nabla_{\tilde{\mathbf{x}}} V^{\mu_s}(\tilde{\mathbf{x}}(t)) \right)^\top \frac{d\tilde{\mathbf{x}}(t)}{dt} \right) \\ &\stackrel{(2.14)}{=} \min_{\boldsymbol{\mu}} H(\tilde{\mathbf{x}}, \nabla_{\tilde{\mathbf{x}}} V^{\mu_s}(\tilde{\mathbf{x}}), \boldsymbol{\mu}) \end{aligned} \quad (4.72)$$

(vgl. [BJ16a], [Büh20, Abschnitt 3.1.2]). Die Value Function $V^{\mu_s}(\tilde{\mathbf{x}})$ sowie $(\nabla_{\tilde{\mathbf{x}}} V^{\mu_s}(\tilde{\mathbf{x}}))^\top \frac{d\tilde{\mathbf{x}}(t)}{dt}$ werden entsprechend [BJ16a] durch lineare Funktionsapproximatoren (vgl. Abschnitt 2.1.3)

$$V^{\mu_s}(\tilde{\mathbf{x}}) = \boldsymbol{\phi}^\top(\tilde{\mathbf{x}}) \hat{\mathbf{w}}_\phi(s) \quad (4.73)$$

und

$$\left(\nabla_{\tilde{\mathbf{x}}} V^{\mu_s}(\tilde{\mathbf{x}}) \right)^\top \frac{d\tilde{\mathbf{x}}(t)}{dt} = \boldsymbol{\psi}^\top(\tilde{\mathbf{x}}, \mathbf{u}) \hat{\mathbf{w}}_\psi(s) \quad (4.74)$$

beschrieben. Da im betrachteten LQ-Fall nach Problem 4.3 die optimale Value Function durch $V^*(\tilde{\mathbf{x}}) = \tilde{\mathbf{x}}^\top \tilde{\mathbf{P}}^* \tilde{\mathbf{x}}$ gegeben ist (vgl. [Föll16, S. 342]), lassen sich wegen⁸³

$$V^*(\tilde{\mathbf{x}}) = \tilde{\mathbf{x}}^\top \tilde{\mathbf{P}}^* \tilde{\mathbf{x}} = \underbrace{(\tilde{\mathbf{x}} \otimes_r \tilde{\mathbf{x}})^\top}_{=: \boldsymbol{\phi}(\tilde{\mathbf{x}})} \underbrace{\text{vecr}(\tilde{\mathbf{P}}^*)}_{=: \mathbf{w}_\phi^*} \quad (4.75)$$

und

$$\begin{aligned} \left(\nabla_{\tilde{\mathbf{x}}} V^*(\tilde{\mathbf{x}}) \right)^\top \frac{d\tilde{\mathbf{x}}(t)}{dt} &= 2\tilde{\mathbf{x}}^\top \tilde{\mathbf{P}}^* \left(\tilde{\mathbf{A}}\tilde{\mathbf{x}} + \tilde{\mathbf{B}}\mathbf{u} \right) \\ &= (\tilde{\mathbf{x}} \otimes_r \tilde{\mathbf{x}})^\top \text{vecr} \left(\tilde{\mathbf{A}}^\top \tilde{\mathbf{P}}^* + \tilde{\mathbf{P}}^* \tilde{\mathbf{A}} \right) + (\tilde{\mathbf{x}} \otimes \mathbf{u})^\top 2\text{vec} \left(\tilde{\mathbf{B}}^\top \tilde{\mathbf{P}}^* \right) \\ &= \underbrace{\left[(\tilde{\mathbf{x}} \otimes_r \tilde{\mathbf{x}})^\top \quad (\tilde{\mathbf{x}} \otimes \mathbf{u})^\top \right]}_{=: \boldsymbol{\psi}^\top(\tilde{\mathbf{x}}, \mathbf{u})} \underbrace{\begin{bmatrix} \text{vecr} \left(\tilde{\mathbf{A}}^\top \tilde{\mathbf{P}}^* + \tilde{\mathbf{P}}^* \tilde{\mathbf{A}} \right) \\ 2\text{vec} \left(\tilde{\mathbf{B}}^\top \tilde{\mathbf{P}}^* \right) \end{bmatrix}}_{=: \mathbf{w}_\psi^*} \end{aligned} \quad (4.76)$$

geeignete Funktionsapproximatoren $\boldsymbol{\phi}(\tilde{\mathbf{x}})$ und $\boldsymbol{\psi}(\tilde{\mathbf{x}}, \mathbf{u})$ wählen. Aus der Minimierung der Hamilton-Funktion (vgl. (2.14), (2.18) und (4.72)) sowie der Linearität von $\boldsymbol{\psi}(\tilde{\mathbf{x}}, \mathbf{u})$ bezüglich \mathbf{u} folgt

$$\begin{aligned} \boldsymbol{\mu}^*(\tilde{\mathbf{x}}) &= \arg \min_{\mathbf{u}} \left(\left(\nabla_{\tilde{\mathbf{x}}} V^*(\tilde{\mathbf{x}}) \right)^\top \frac{d\tilde{\mathbf{x}}(t)}{dt} + r(\tilde{\mathbf{x}}, \mathbf{u}) \right) \\ &= \arg \min_{\mathbf{u}} \left(\boldsymbol{\psi}^\top(\tilde{\mathbf{x}}, \mathbf{u}) \mathbf{w}_\psi^* + r(\tilde{\mathbf{x}}, \mathbf{u}) \right) \\ &= -\frac{1}{2} \mathbf{R}^\top (\nabla_{\mathbf{u}} \boldsymbol{\psi}(\tilde{\mathbf{x}}, \mathbf{u}))^\top \mathbf{w}_\psi^* \\ &= -\tilde{\mathbf{K}}^* \tilde{\mathbf{x}} \end{aligned} \quad (4.77)$$

⁸³ Hierbei bezeichnet \otimes das Kronecker-Produkt, \otimes_r das reduzierte Kronecker-Produkt, welches nur nicht-redundante Elemente beinhaltet, $\text{vec}(\cdot)$ die Vektorisierung einer Matrix, wobei die Spalten vertikal konkateniert werden, und $\text{vecr}(\cdot)$ die Vektorisierung einer symmetrischen Matrix \mathbf{M} , sodass $\mathbf{x}^\top \mathbf{M} \mathbf{x} = (\mathbf{x} \otimes_r \mathbf{x})^\top \text{vecr}(\mathbf{M})$ gilt.

mit

$$\tilde{\mathbf{K}}^* = \frac{1}{2} \mathbf{R}^{-1} \text{mat} \left(\left[\mathbf{w}_\psi^* \right]_{\frac{\tilde{n}(\tilde{n}+1)}{2} : (\frac{\tilde{n}(\tilde{n}+1)}{2} + \tilde{n}p)}, p, \tilde{n} \right), \quad (4.78)$$

wobei $\tilde{n} = n + n_\zeta$ der Dimension des erweiterten Zustands $\tilde{\mathbf{x}}$ und p der Dimension der Eingangsgröße \mathbf{u} entspricht⁸⁴. Damit reduziert sich das ADP-Problem auf das datenbasierte Bestimmen insbesondere des Gewichts \mathbf{w}_ψ^* . Das Vorgehen der hierbei verwendeten Value Iteration [BJ16a] ist in Anhang B.8 skizziert.

Da eine ausreichende Systemanregung, d. h. Regularität der Matrizen \mathbf{D}_ϕ (B.34) und \mathbf{D}_ψ (B.35), während der Datenaufzeichnung für einen erfolgreichen Adaptionprozess benötigt wird [BJ16a, Assumption 3] (vgl. auch Kapitel 5), wird Gaußsches weißes Rauschen mit der Standardabweichung σ_{expl} zum Systemeingang $\mathbf{u}(t)$ addiert. Weil hierdurch jedoch nur der Systemzustand $\mathbf{x}(t)$ und nicht der erweiterte Zustand $\tilde{\mathbf{x}}(t) = [\mathbf{x}^\top(t) \quad \zeta^\top(t)]^\top$ angeregt wird, muss zudem der Zustand $\zeta(t)$, der den Solltrajektorienverlauf repräsentiert, angeregt werden. Um angeregte Daten $\zeta(t)$ für den Adaptionprozess zu generieren, ohne dabei den Systemzustand $\mathbf{x}(t)$ weiter zu beeinflussen, wird die in Abbildung 4.3 gezeigte Struktur verwendet.

Der Referenzzustand $\zeta_{\text{train}}(t)$ wird zu Beginn jedes Integrationsintervalls der Länge T_{IRL} zufällig initialisiert, indem für jedes Element ein zufälliger Wert aus der Normalverteilung $\mathcal{N}(0; \sigma_{\text{expl}}^2)$ gezogen wird. Hierdurch wird einerseits sichergestellt, dass unterschiedliche Referenzzustände in den Trainingsdaten vorhanden sind, andererseits folgt ζ_{train} innerhalb eines IRL-Integrationsintervalls der Dauer T_{IRL} der durch $\mathbf{D}' = \mathbf{D} - \frac{\gamma}{2} \mathbf{I}$ beschriebenen Dynamik, weshalb eine ADP-kompatible Solltrajektorienarstellung (vgl. Definition 4.1) zur

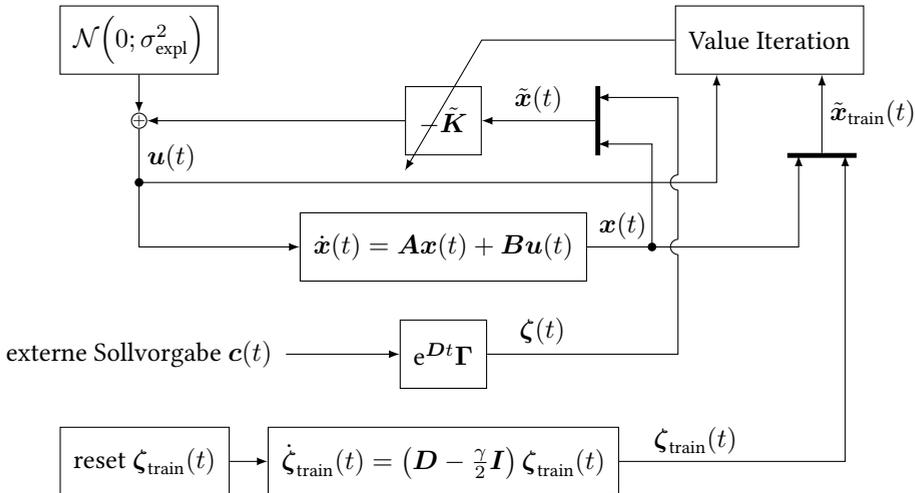


Abbildung 4.3: Struktur des vorgestellten zeitkontinuierlichen ADP-Trajektorienfolgeregelungsansatzes.

⁸⁴ Der Operator $\text{mat}(\cdot)$ bildet aus einem Vektor eine Matrix, wobei $\text{mat}(\text{vec}(\mathbf{M}), \tilde{n}, p) = \mathbf{M} \in \mathbb{R}^{\tilde{n} \times p}$ gilt, zudem bezeichne $[\mathbf{v}]_{n_1:n_2} = [v_{n_1} \quad \dots \quad v_{n_2}]^\top$ mit $n_1, n_2 \in \mathbb{N}_{>0}, n_1 < n_2$, einen Teilvektor eines Vektors \mathbf{v} .

Erzeugung der Lerndaten vorliegt. Unabhängig davon kann der gewünschte Verlauf der Trajektorie durch beliebige Vorgabe des Parameters $c(t)$ erfolgen.

4.2.5 Simulationsergebnisse

Die simulative Evaluation des zuvor vorgestellten ADP-basierten Solltrajektorienfolgereglers erfolgt anhand einer Implementierung in MATLAB R2020a (64 Bit) und SIMULINK 10.1. Zunächst werden das während der Simulation verwendete Beispielsystem und die Trainingsparameter vorgestellt. Danach werden durch Wahl der Hyperparameter λ_{ij} und ν_{ij} verschiedene Solltrajektorienparametrisierungen betrachtet (vgl. Bemerkung 4.1) und mit der aus der Literatur bekannten stationären Sollwertvorgabe verglichen. Abschließend wird der Einfluss von Messrauschen, das auf den Systemzustand $\mathbf{x}(t)$ addiert wird, simulativ untersucht.

4.2.5.1 Systemmodell und Parametrierung des ADP-Algorithmus

Betrachtet wird im Folgenden ein Feder-Masse-Dämpfer-System mit der Masse $m_{\text{sys}} = 1$ kg, der Federsteifigkeit $k_{\text{sys}} = 1$ N m⁻¹ und der Dämpfung $d_{\text{sys}} = 1$ N s m⁻¹, auf das eine Kraft als Stellgröße wirkt. Das System ist demnach durch die Zustandsdifferenzialgleichung

$$\dot{\mathbf{x}}(t) = \begin{bmatrix} 0 & 1 \\ -\frac{k_{\text{sys}}}{m_{\text{sys}}} & -\frac{d_{\text{sys}}}{m_{\text{sys}}} \end{bmatrix} \mathbf{x}(t) + \begin{bmatrix} 0 \\ \frac{1}{m_{\text{sys}}} \end{bmatrix} u(t) = \begin{bmatrix} 0 & 1 \\ -1 & -1 \end{bmatrix} \mathbf{x}(t) + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u(t) \quad (4.79)$$

beschrieben. Zudem sei $\mathbf{R} = 1$, $\mathbf{Q} = \begin{bmatrix} 100 & 0 \\ 0 & 0 \end{bmatrix}$ und $\gamma = 0,01$ (vgl. Problem 4.3). Somit soll der Zustand $x_1(t)$ einer noch näher zu spezifizierenden Solltrajektorie folgen. Der Initialzustand werde zu $\mathbf{x}(0) = [1 \quad 1]^T$ gesetzt.

Für den Trainingsvorgang wird das Integrationsintervall zu $T_{\text{IRL}} = 0,01$ s, das Anregungsrauschen zu $\sigma_{\text{expl}} = 1$ und die Anzahl der verwendeten Tupel zu $M = 200$ gesetzt. Die zu lernenden Gewichte werden mit $\hat{\mathbf{w}}_{\psi}(0) = \mathbf{0}$ und $\hat{\mathbf{w}}_{\phi}(0) = \text{vecr}(\mathbf{Z}^T \mathbf{Z})$ initialisiert, wobei $\mathbf{Z} \in \mathbb{R}^{\tilde{n} \times \tilde{n}}$ zufällige Elemente⁸⁵ enthält, wodurch die initiale Value Function positiv definit ist (vgl. [BJ16a, Theorem 2]).

4.2.5.2 Stationäre Sollwertvorgabe (Vergleichsmethode)

Der vorgestellte Ansatz erlaubt als Spezialfall auch die in der Literatur häufig verwendete stationäre Sollwertvorgabe (vgl. Abschnitt 2.2). Dieser Fall, welcher der Vorgabe eines Polynoms nullter Ordnung entspricht, soll zu Vergleichszwecken hinzugezogen werden. Mit $J_1 = 1$, $\lambda_{11} = 0$ und $\nu_{11} = 1$ folgt $\mathbf{D} = 0$, $\mathbf{S} = [1 \quad 0]^T$ und

$$\tilde{\mathbf{K}}^* = [9,0499 \quad 3,3703 \quad -9,9288]. \quad (4.80)$$

⁸⁵ Konkret wird eine Gleichverteilung im Intervall $(0; 0,1]$ verwendet.

Der Parameter $c(t)$ entspricht in diesem Fall dem Solltrajektorienverlauf $x_{r,1}(t)$. Abbildung 4.4 zeigt den Verlauf der Reglerparameter $\tilde{\mathbf{K}}(t)$ und die Gewichtsfehlnorm $\|\tilde{\mathbf{K}}(t) - \tilde{\mathbf{K}}^*\|_2$. Die ohne Verwendung des Systemmodells gelernten Reglergewichte konvergieren bis auf eine Fehlernorm von $6,28 \cdot 10^{-14}$.

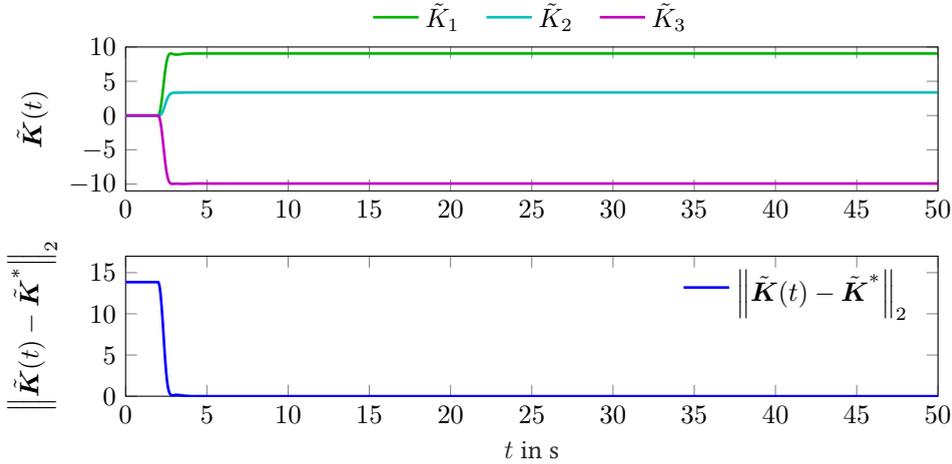


Abbildung 4.4: Reglergewichte $\tilde{\mathbf{K}}(t)$ (oben) und Gewichtsfehlnorm $\|\tilde{\mathbf{K}}(t) - \tilde{\mathbf{K}}^*\|_2$ (unten) für das Beispiel der stationären Sollvorgabe. Während der ersten 2 s werden $M = 200$ Datentupel aufgezeichnet.

4.2.5.3 Harmonischer Oszillator als Solltrajektorienparametrierung

In diesem Beispiel werden zur Beschreibung von $x_{r,1}$ zwei konjugiert-komplexe Polpaare $\lambda_{11} = \lambda_{12}^* = j\omega_1$ und $\lambda_{13} = \lambda_{14}^* = j\omega_2$ mit $\nu_{1j} = 1, j = 1, \dots, 4$, verwendet (vgl. Annahme 4.1). Während die Hyperparameter durch die beispielhaft gewählten Werte $\omega_1 = 0,5$ und $\omega_2 = 1,3$ die Frequenzen der darstellbaren harmonischen Schwingungen festlegen, lassen sich durch $c_{110} = c_{120}^* \in \mathbb{C}$ und $c_{130} = c_{140}^* \in \mathbb{C}$ nach Bemerkung 4.1 deren Phasenverschiebungen und Amplituden parametrieren. Da die Konstruktion nach (4.14) zu einem komplexwertigen \mathbf{D} führt, wird die in Anhang B.3 gezeigte Zustandstransformation angewandt. Dies führt auf

$$\mathbf{D} = \begin{bmatrix} 0 & \omega_1 & 0 & 0 \\ -\omega_1 & 0 & 0 & 0 \\ 0 & 0 & 0 & \omega_2 \\ 0 & 0 & -\omega_2 & 0 \end{bmatrix} \quad \text{und} \quad \mathbf{S} = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}. \quad (4.81)$$

Das optimale Regelgesetz berechnet sich nach Satz 4.1 zu

$$\tilde{\mathbf{K}}^* = [9,0499 \quad 3,3703 \quad -9,6992 \quad -2,1628 \quad -8,1630 \quad -5,5458]. \quad (4.82)$$

Die Sollzustandstrajektorie wird durch den Parameter c beschrieben. Dieser wird zu den in Tabelle 4.1 gezeigten Zeitpunkten auf die dort aufgeführten Werte gesetzt und anschließend konstant gehalten.

Zeit t	$c_{110} = c_{120}^*$	$2 c_{110} $	$c_{130} = c_{140}^*$	$2 c_{130} $
0 s	$0,1768 + j0,1768$	0,5	$0,2500 + j0$	0,5
9,7 s	$0 + j0$	0	$0,1294 + j0,4830$	1
18,7 s	$0,2000 + j0$	0,4	$0 + j0,1000$	0,2

Tabelle 4.1: Wahl des Parameters c für das Beispiel der harmonischen Referenzzustandsdarstellung. Nach (4.11) stellt $2|c_{110}|$ die Amplitude der Schwingung mit Frequenz $\omega_1 = 0,5$ und $2|c_{130}|$ die Amplitude zur Frequenz $\omega_2 = 1,3$ dar.

Der resultierende Sollzustandsverlauf $x_{r,1}(t)$, Zustandsverlauf $x(t)$ und Stellgrößenverlauf $u(t)$ ist in Abbildung 4.5 gezeigt. Während der ersten zwei Sekunden wird der Datenspeicher mit $M = 200$ Tupeln gefüllt. Danach werden die Gewichte \hat{w}_ϕ und \hat{w}_ψ adaptiert und das Anregungsrauschen zu null gesetzt. Der Zustand x_1 folgt dem Sollzustandsverlauf $x_{r,1}$ optimal im Sinne des Gütemaßes (4.46). Dies bestätigt Abbildung 4.6, in welcher der zeitliche Verlauf der gelernten Reglerparameter $\tilde{K}(t)$ sowie deren Gewichtsfehlnorm $\left\| \tilde{K}(t) - \tilde{K}^* \right\|_2$ zu sehen sind. Die Reglergewichte konvergieren hierbei bis auf eine Fehlernorm von $2,83 \cdot 10^{-13}$. Der Rang der Datenmatrizen D_ϕ und D_ψ (vgl. Anhang B.8) ist in Abbildung 4.7 gezeigt, welche eine ausreichende Anregung der Daten bestätigt.

In Abbildung 4.5 sind des Weiteren die Trajektorien $x_{1,s}(t)$, $x_{2,s}(t)$ und $u_s(t)$ der Vergleichslösung mit stationärer Sollzustandsvorgabe nach Abschnitt 4.2.5.2 gezeigt. Da die Vorgabe des aktuellen Sollzustands $x_{r,1}(t)$ keinerlei Information über den weiteren Verlauf der Solltrajektorie beinhaltet, nimmt der Vergleichsregler mit stationärer Sollwertvorgabe stets an, der aktuelle Sollzustand solle bis ins Unendliche gehalten werden. Dies führt zu einem merklichen Zeitversatz der Trajektorie $x_{1,s}(t)$. Wie Abbildung 4.8 verdeutlicht, korreliert dieser Zeitversatz der Trajektorie mit deutlich höheren Kosten. Die in dieser Arbeit präsentierte Methode berücksichtigt hingegen den zukünftigen Verlauf der Solltrajektorie und führt aufgrund der damit verbundenen Vorausschaufähigkeit zu signifikant geringeren Kosten im Sinne des Gütefunktional. So werden die kumulierten Kosten am Ende des Simulationsbeispiels, wie in Abbildung 4.8 ersichtlich, um etwa 69 % reduziert.

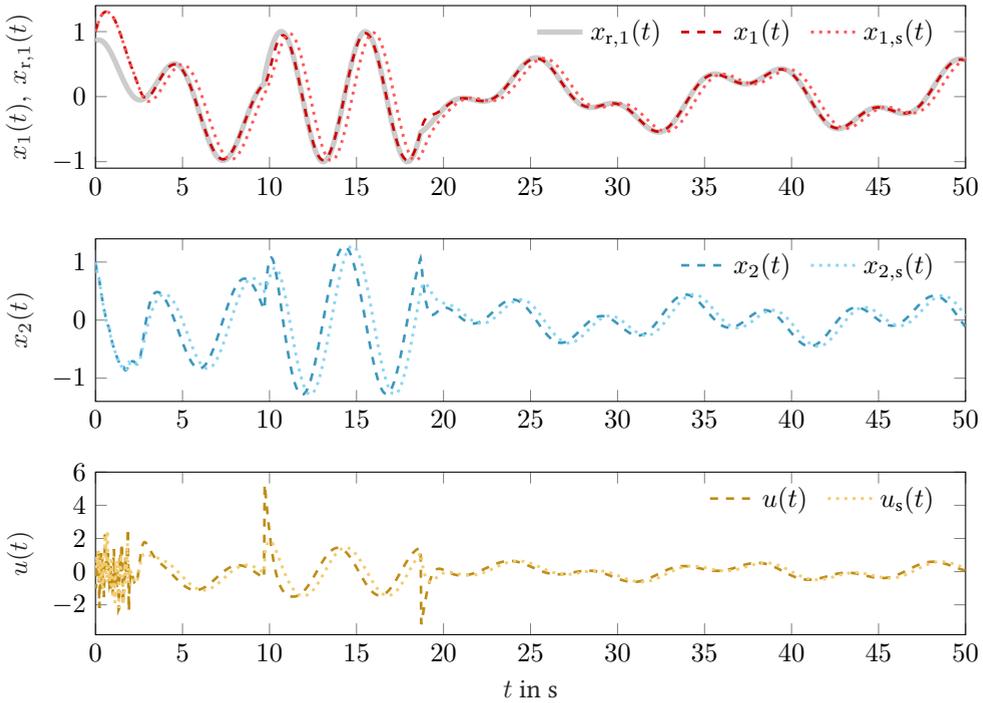


Abbildung 4.5: Sollzustandsverlauf $x_{r,1}(t)$, Zustandsverlauf $\mathbf{x}(t)$ und Stellgröße $u(t)$ für das Beispiel des harmonischen Oszillators als Solltrajektorienparametrierung im Vergleich zur stationären Sollzustandsvorgabe (gekennzeichnet durch den Index s).

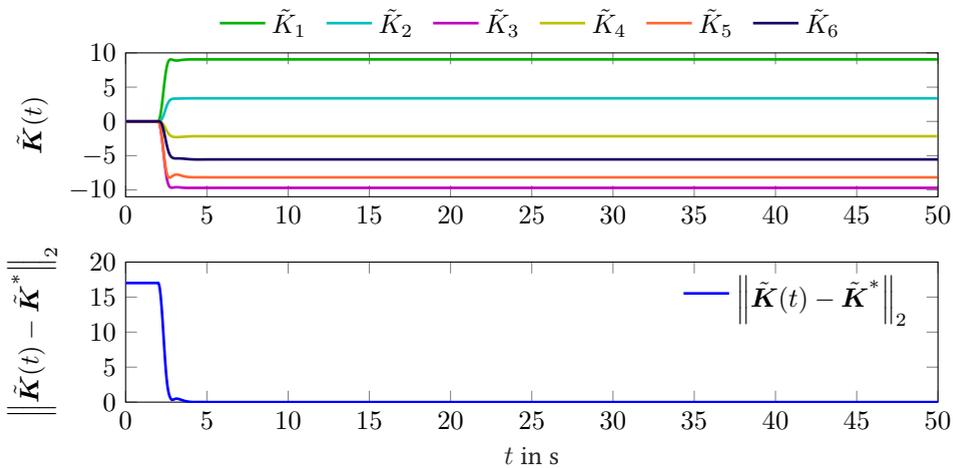


Abbildung 4.6: Reglergewichte $\tilde{K}(t)$ (oben) und Gewichtsfehlnorm $\|\tilde{K}(t) - \tilde{K}^*\|_2$ (unten) für das Beispiel des harmonischen Oszillators als Solltrajektorienparametrierung. Während der ersten 2 s werden $M = 200$ Datentupel aufgezeichnet.

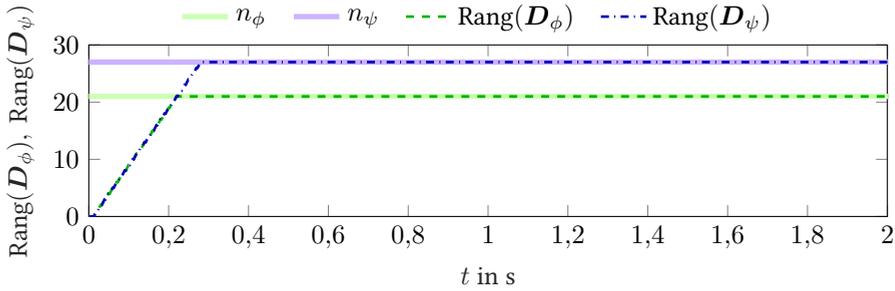


Abbildung 4.7: Rang der Datenmatrizen $D_\phi \in \mathbb{R}^{n_\phi \times n_\phi}$ und $D_\psi \in \mathbb{R}^{n_\psi \times n_\psi}$ während der ersten 2 s. Für das gewählte Beispiel des harmonischen Oszillators als Solltrajektorienparametrierung gilt $n_\phi = \frac{\tilde{n}(\tilde{n}+1)}{2} = 21$ und $n_\psi = \frac{\tilde{n}(\tilde{n}+1)}{2} + \tilde{n}p = 27$.

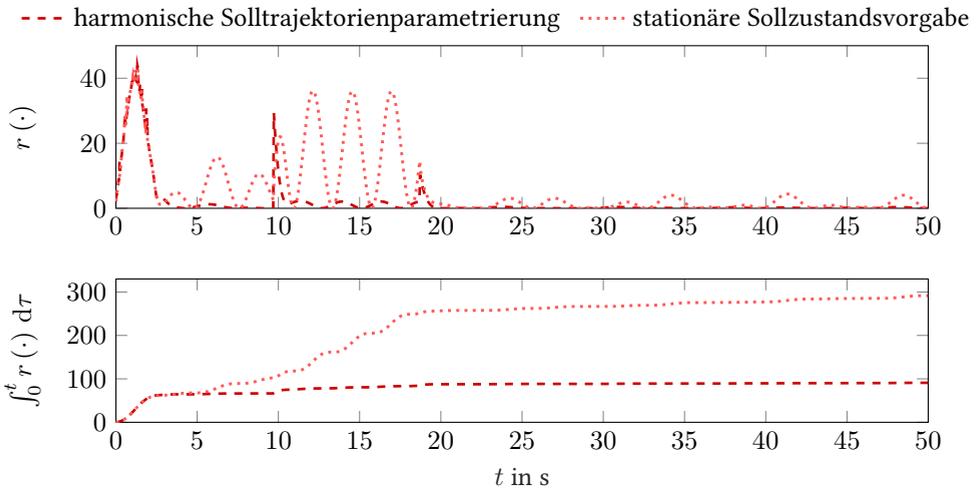


Abbildung 4.8: Kosten $r(\mathbf{x}(t), \mathbf{x}_r(t), \mathbf{u}(t))$ und kumulierte Kosten $\int_0^t r(\mathbf{x}(\tau), \mathbf{x}_r(\tau), \mathbf{u}(\tau)) d\tau$ für das Beispiel des harmonischen Oszillators als Solltrajektorienparametrierung im Vergleich zur stationären Sollzustandsvorgabe.

4.2.5.4 Polynomvorgabe

Im Folgenden wird zur Beschreibung des Sollverlaufs des ersten Systemzustands ein Polynom zweiten Grades verwendet, d. h. es gilt

$$x_{r,1}(t) = c_{110} + c_{111}t + c_{112}t^2. \quad (4.83)$$

Dies entspricht nach Bemerkung 4.1 einem dreifachen Eigenwert bei null ($\lambda_{11} = 0, \nu_{11} = 3$). Mit

$$D = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \quad \text{und} \quad S = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (4.84)$$

berechnet sich die optimale Zustandsrückführungsmatrix zu

$$\tilde{K}^* = [9,0499 \quad 3,3703 \quad -9,9288 \quad -4,3181 \quad -0,8922]. \quad (4.85)$$

Der Parameter c , der den Solltrajektorienverlauf beschreibt, wird, wie in Tabelle 4.2 gezeigt, gesetzt und zwischen den dort aufgeführten Zeitpunkten konstant gehalten.

Zeit t	$\mathbf{x}_{r,1}(t)$	$\dot{\mathbf{x}}_{r,1}(t)$	$\ddot{\mathbf{x}}_{r,1}(t)$	c_{110}	c_{111}	c_{112}
0 s	0	0	0	0	0	0
5 s	0	1	0	-5	1	0
10 s	5	0	0	5	0	0
15 s	5	0	-0,6	-62,5	9	-0,3
18 s	2,3	-1,4	1	189,5	-19,4	0,5
20 s	1,5	1	-0,5	-118,5	11	-0,25
25 s	0,25	-2	1,4	487,75	-37	0,7
26,4 s	-1,1780	0	0	-1,1780	0	0
28 s	-1,1780	0	0,07	26,2620	-1,96	0,035
35 s	0,5370	0,5	-0,1	-78,2130	4	-0,05

Tabelle 4.2: Wahl des Parameters c für das Beispiel der Polynomvorgabe. Hierbei entspricht $\zeta(t) = [\mathbf{x}_{r,1}(t) \quad \dot{\mathbf{x}}_{r,1}(t) \quad \ddot{\mathbf{x}}_{r,1}(t)]^T$ (vgl. (4.12)).

Der Sollzustandsverlauf $x_{r,1}(t)$, Zustandsverlauf $\mathbf{x}(t)$ und Stellgrößenverlauf $u(t)$ ist in Abbildung 4.9 gegeben. Abbildung 4.10 zeigt den zeitlichen Verlauf der geschätzten Reglerparameter $\tilde{K}(t)$ und die Gewichtsfehlernorm $\left\| \tilde{K}(t) - \tilde{K}^* \right\|_2$. Die Reglergewichte konvergieren bis auf eine Fehlernorm von $2,63 \cdot 10^{-8}$. Der Rang der Datenmatrizen D_ϕ und D_ψ ist Abbildung 4.11 zu entnehmen.

Auch in diesem Fall weist der in dieser Arbeit vorgestellte Ansatz, eine ADP-kompatible Beschreibungsform des Solltrajektorienverlaufs zu verwenden, Vorteile gegenüber der Vergleichsmethode nach Abschnitt 4.2.5.2 auf. So folgt gemäß Abbildung 4.9 der Zustand $x_1(t)$ dem Sollzustand $x_{r,1}(t)$ präziser im Vergleich zu $x_{1,s}(t)$. Letztlich kann, wie Abbildung 4.12 offenbart, der neuartige Ansatz auch in diesem Beispiel die Kosten im Vergleich zur stationären Sollzustandsvorgabe signifikant reduzieren. Die kumulierten Kosten werden in dieser beispielhaften Simulation um 38 % verringert.

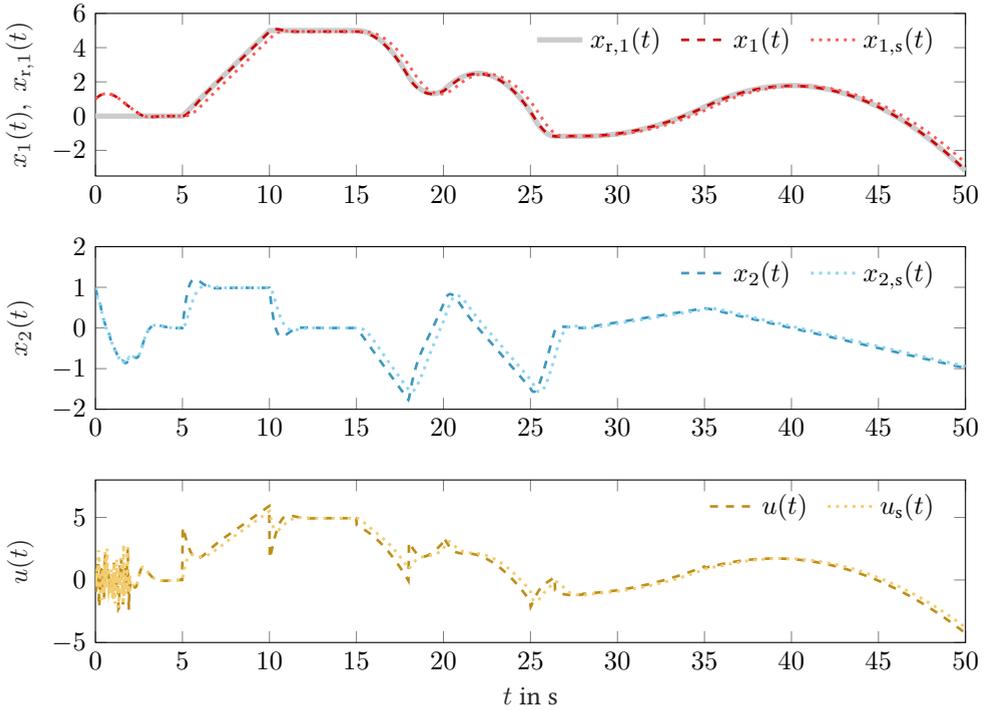


Abbildung 4.9: Sollzustandsverlauf $x_{r,1}(t)$, Zustandsverlauf $\boldsymbol{x}(t)$ und Stellgröße $u(t)$ für das Beispiel der Polynomvorgabe im Vergleich zur stationären Sollzustandsvorgabe (gekennzeichnet durch den Index s).

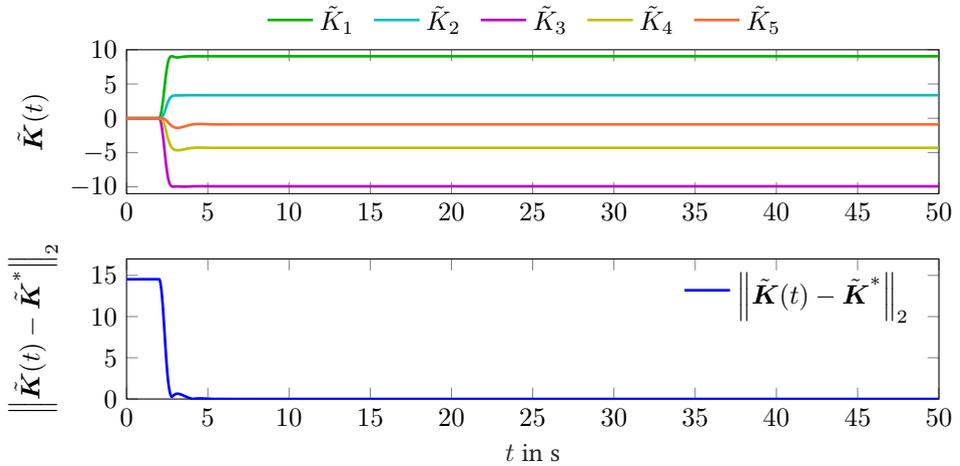


Abbildung 4.10: Reglergewichte $\tilde{\mathbf{K}}(t)$ (oben) und Gewichtsfehlernorm $\|\tilde{\mathbf{K}}(t) - \tilde{\mathbf{K}}^*\|_2$ (unten) für das Beispiel der Polynomvorgabe. Während der ersten 2 s werden $M = 200$ Datentupel aufgezeichnet.

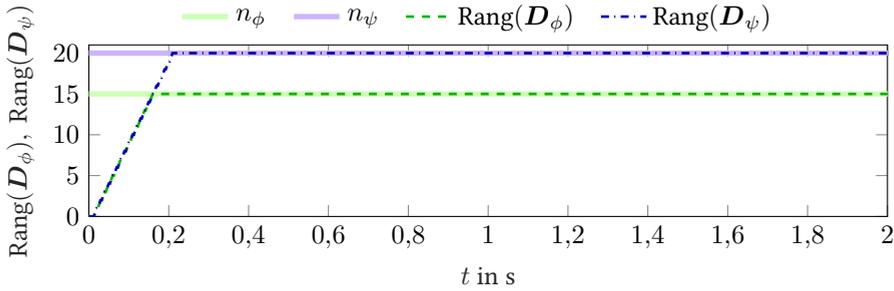


Abbildung 4.11: Rang der Datenmatrizen $D_\phi \in \mathbb{R}^{n_\phi \times n_\phi}$ und $D_\psi \in \mathbb{R}^{n_\psi \times n_\psi}$ während der ersten 2 s. Für das gewählte Beispiel der Polynomvorgabe gilt $n_\phi = \frac{\tilde{n}(\tilde{n}+1)}{2} = 15$ und $n_\psi = \frac{\tilde{n}(\tilde{n}+1)}{2} + \tilde{n}p = 20$.

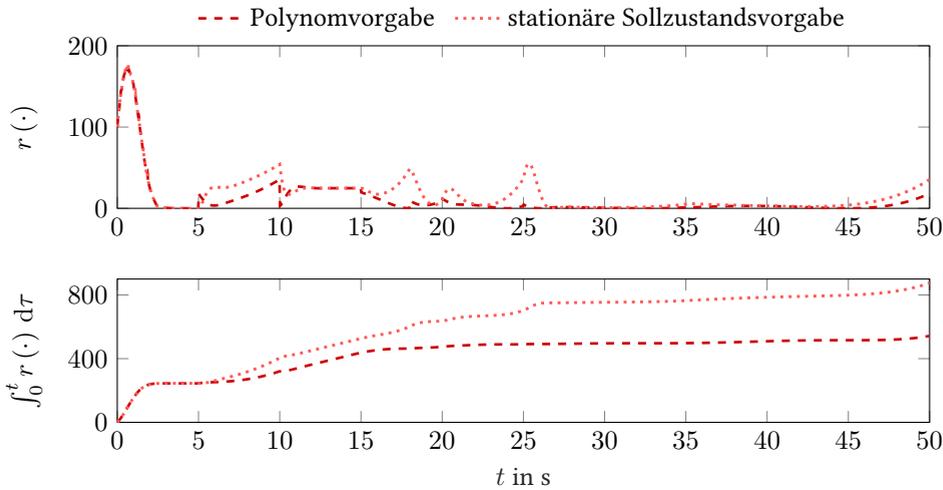


Abbildung 4.12: Kosten $r(\mathbf{x}(t), \mathbf{x}_r(t), \mathbf{u}(t))$ und kumulierte Kosten $\int_0^t r(\mathbf{x}(\tau), \mathbf{x}_r(\tau), \mathbf{u}(\tau)) d\tau$ für das Beispiel der Polynomvorgabe im Vergleich zur stationären Sollzustandvorgabe.

4.2.5.5 Einfluss von Messrauschen

Abschließend wird für das gegebene Simulationsbeispiel noch der Einfluss von weißem Gaußschem Messrauschen mit unterschiedlicher Standardabweichung σ_m , das zu $\mathbf{x}(t)$ addiert wird, untersucht. Hierbei wird die zuvor verwendete Solltrajektorienarstellung der Polynomvorgabe wie in Abschnitt 4.2.5.4 verwendet. Abbildung 4.13 zeigt den Zustandsverlauf $x_1(t)$ sowie die Stellgröße $u(t)$ für Standardabweichungen des Messrauschens von $\sigma_m = 0$ bis $\sigma_m = 2$. Ein Messrauschen bis $\sigma_m = 0,2$ führt zu kaum merklichen Abweichungen. Bis zu einem Messrauschen von $\sigma_m = 1$ bleibt der geschlossene Regelkreis stabil und kann der vorgegebenen Solltrajektorie folgen. Für $\sigma_m = 2$ resultiert ein instabiles Regelgesetz. Die Fehlernorm im Vergleich zum optimalen Regelgesetz ist für unterschiedliches Messrauschen in Abbildung 4.14

gegeben. Dabei ist zu erkennen, dass der gelernte Regler mit zunehmendem Messrauschen stärker vom Optimalregler abweicht. Für $\sigma_m = 2$ divergiert \tilde{K} für das Beispielszenario.

Zwar weisen in allen betrachteten Fällen die Datenmatrizen D_ϕ und D_ψ vollen Rang auf und das Messrauschen begünstigt diese Tatsache sogar, jedoch ist das Systemverhalten bei zu starkem Messrauschen im Vergleich zur Systemanregung nicht mehr adäquat in den Messdaten repräsentiert⁸⁶. Somit treten auch hier vergleichbare Effekte wie im zeitdiskreten Fall in Abschnitt 3.3.4.4 auf und die Notwendigkeit einer geeigneten Systemanregung verstärkt sich insbesondere im Fall vorhandenen Messrauschens.

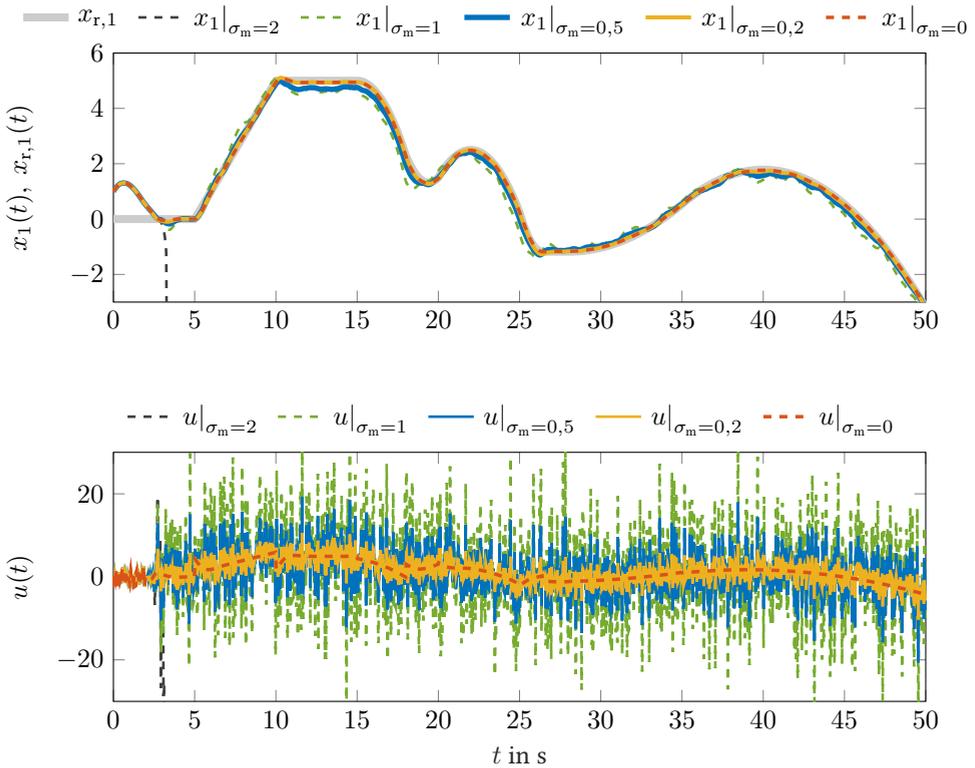


Abbildung 4.13: Sollzustandsverlauf $x_{r,1}(t)$, Zustandsverlauf $x_1(t)$ und Stellgröße $u(t)$ für das Beispiel der Polynomvorgabe unter Einfluss von Messrauschen.

⁸⁶ Um stärkeres additives Messrauschen handhabbar zu machen, kann hingegen (im Rahmen der numerischen Grenzen) mit einer größeren Systemauslenkung während der Datenaufzeichnung gearbeitet werden, d. h. beispielsweise mit einem stärkeren Anregungsrauschen. Für $\sigma_{\text{expl}} = 30$ und ein durch $\sigma_m = 2$ beschriebenes Messrauschen ergibt sich die Fehlernorm des gelernten Regelgesetzes beispielsweise zu 4,2 und der geschlossene Regelkreis ist stabil. Diese Wahl des Anregungs- und Messrauschens führt zu einer ähnlichen Abweichung vom optimalen Regelgesetz wie zuvor mit $\sigma_{\text{expl}} = 1$ und $\sigma_m = 1$ mit einer Fehlernorm von 4,7.

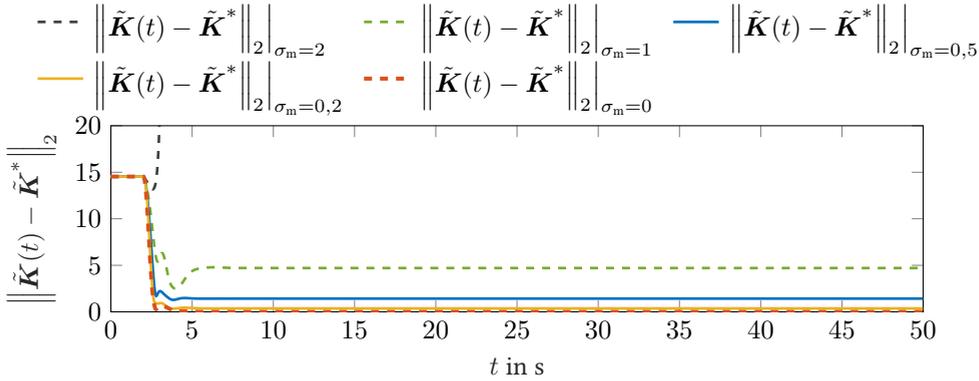


Abbildung 4.14: Gewichtsfehlnorm $\|\tilde{\mathbf{K}}(t) - \tilde{\mathbf{K}}^*\|_2$ für das Beispiel der Polynomvorgabe unter Einfluss von Messrauschen. Während der ersten 2 s werden $M = 200$ Datentupel aufgezeichnet.

4.3 Zusammenfassung

In diesem Kapitel wurde erstmals eine allgemeine Definition zeitkontinuierlicher ADP-kompatibler Solltrajektorien-
darstellungen vorgestellt. Des Weiteren wurde eine neuartige Klasse ADP-kompatibler Solltrajektorien in Form einer Superposition von Lösungen homogener linearer Differenzialgleichungen präsentiert und analysiert. Diese verallgemeinerte Darstellung inkludiert insbesondere konstante Sollzustandsvorgaben, Polynomverläufe, Exponentialfunktionen und harmonische Funktionen. Für LQ-Solltrajektorienfolge-
regelungsprobleme mit gegebenenfalls gedämpfter Solltrajektorien-
dynamik, die diese neue, ADP-kompatible Solltrajektorien-
darstellung nutzen, wurde unter milden Annahmen die Existenz und Eindeutigkeit der stabilisierenden optimalen Lösung gezeigt.

Die Einbeziehung des durch den Parameter \mathbf{c} repräsentierten Solltrajektorienverlaufs in die Value Function V ermöglicht schließlich ein ADP-basiertes, modellfreies Erlernen sowohl der optimalen Value Function als auch des zugehörigen Regelgesetzes. Der hierzu verwendete ADP-Algorithmus ist dabei grundsätzlich austauschbar, beispielhaft wurde die IRL-basierte Value Iteration nach Bian und Jiang [BJ16a] genutzt. Der für den Trainingsvorgang verwendete Solltrajektorienparameter ζ_{train} folgt hierbei während eines Integrationsintervalls der Dauer T_{IRL} der durch $\mathbf{D} - \frac{\gamma}{2}\mathbf{I}$ beschriebenen Dynamik und wird nur zu Beginn der neuen Integrationsintervalle zugunsten einer ausreichenden Anregung neu initialisiert. Insbesondere gehören die während des Trainings verwendeten Datentupel $(\tilde{\mathbf{x}}(t_j), \Theta_j, \tilde{\mathbf{x}}(t_j + T_{\text{IRL}}))$, $j = 1, \dots, M$, (vgl. Anhang B.8) daher zu ADP-kompatiblen Solltrajektorien (vgl. auch Bemerkung 3.1). Aufgrund der in Abbildung 4.3 gezeigten parallelen Struktur von $\zeta(t)$ und $\zeta_{\text{train}}(t)$ kann der Parameter \mathbf{c} , der den Solltrajektorienverlauf aus lokaler Perspektive zum Zeitpunkt t beschreibt, jederzeit beliebig vorgegeben werden. Der im vorliegenden Kapitel vorgestellte Ansatz kann letztlich als Übertragung der ADP-kompatiblen zeitdiskreten Methode aus Abschnitt 3.2 auf

zeitkontinuierliche Systeme interpretiert werden⁸⁷. Als Grenze des Ansatzes lässt sich nennen, dass in der Zukunft liegende Sprünge durch die Solltrajektorien­darstellung nach (4.7) nicht exakt abgebildet, sondern lediglich approximiert werden können. Zudem muss der Einfluss der Systemanregung auf die Trainingsdatentupel signifikant stärker sein als der Einfluss des Messrauschens (vgl. Abschnitt 4.2.5.5). Letzteres ist jedoch kein spezifisches Problem der ADP-kompatiblen Solltrajektorien­folgeregelung, sondern Systemidentifikationsproblemen und adaptiven Methoden inhärent (vgl. [Bit84], [KT04, S. 122]).

Die vorgestellten Simulationsergebnisse zeigen die Flexibilität der vorgestellten Methode auf. So lassen sich im Rahmen ADP-kompatibler Solltrajektorien­darstellungen vielfältige Referenztrajektorienverläufe in modellfreie zeitkontinuierliche ADP-Algorithmen integrieren. Die explizite Abhängigkeit der gelernten Value Function und des zugehörigen Optimalreglers vom Parameter c erlaubt im Gegensatz zu bestehenden Ansätzen einerseits, vielfältige Solltrajektorien von außen vorzugeben, andererseits kann hierdurch der zukünftige Verlauf der Solltrajektorie berücksichtigt werden. Insbesondere führt die Einbeziehung des Solltrajektorienverlaufs zu einem vorausschauenden Verhalten, das, im Vergleich zu einer aus der Literatur bekannten konstanten Sollwertvorgabe, mit deutlich reduzierten Gesamtkosten korreliert. Somit beantworten die neuartigen, flexibel einsetzbaren Mechanismen für ADP-basierte Solltrajektorien­folgeregler die in Abschnitt 2.4.1 formulierte Forschungsfrage 1 nach ADP-kompatiblen Solltrajektorien­darstellungen und deren Integration in modellfreie ADP-Ansätze für den zeitkontinuierlichen Fall.

⁸⁷ Aufgrund der Stetigkeitsanforderungen nach Definition 4.1 sei an dieser Stelle angemerkt, dass eine direkte Übertragung der in Abschnitt 3.3 vorgestellten zeitdiskreten Repräsentation, welche beliebige Sprünge innerhalb des Vorausschauhorizonts exakt berücksichtigt, auf den zeitkontinuierlichen Fall nicht möglich ist.

5 Konvergenzbedingungen zeitkontinuierlicher adaptiver Optimalregler

Ein zentrales Element ADP-basierter Optimalregler stellt eine geeignete Anregung des Systems dar. Diese ermöglicht adaptiven Optimalreglern, Informationen über das Systemverhalten sowie die damit verbundenen Kosten zu gewinnen und sich kostenoptimal anzupassen (vgl. Abschnitt 2.3, (3.29), (3.115), (B.37) und (B.40)).

Den Ausgangspunkt des vorliegenden Kapitels⁸⁸ bildet daher die Analyse der zeitkontinuierlichen PE-Bedingung (2.27) und deren Bedeutung für die Konvergenz eines gradientenbasierten Policy-Evaluation-Schrittes (vgl. Abschnitt 2.1.4). Betrachtet wird hierbei ein zeitkontinuierliches, eingangsaффines Differenzialspiel, das eine Generalisierung des durch (2.10) und (2.11) gegebenen Optimierungsproblems darstellt. Auf dieser Basis wird der zentrale Beitrag dieses Kapitels in Form von hinreichenden Bedingungen an den Systemzustand, die garantieren, dass das zur Adaption benötigte Signal die PE-Bedingung erfüllt, hergeleitet. Dieses Signal folgt bei Verwendung polynomieller Basisfunktionen zur Critic-Funktionsapproximation aus einer nichtlinearen Transformation des Systemzustands. Aufbauend auf diesen theoretischen Erkenntnissen wird schließlich am Beispiel eines exakt zustandslinearisierbaren Systems ein geeignetes Anregungssignal mithilfe eines Vorsteuerungsentwurfs konstruiert, um Simulationen zu ermöglichen. Anhand dieses Simulationsbeispiels werden die vorgestellten Anregungssignale mit einer Anregung durch weißes Rauschen verglichen und Unterschiede diskutiert.

In den Abschnitten 5.1–5.3 erfolgt zunächst eine Einführung und Analyse in die betrachtete Problemstellung, bevor in den Abschnitten 5.4–5.7 die Hauptergebnisse dieses Kapitels präsentiert werden.

5.1 Eingangsaффines Differenzialspiel mit unbekanntem Gegenspielern

In diesem Kapitel werden eingangsaффine Differenzialspiele, bei denen die Spieler a priori kein genaues Modell über das Verhalten der jeweils anderen Spieler besitzen, gemäß der nachfolgenden Definition betrachtet. Um die Ziele der N auf ein System einwirkenden Regler

⁸⁸ Dieses Kapitel basiert auf einem im Rahmen der vorliegenden Arbeit entstandenen eigenen Beitrag [KKBH23].

formal definieren zu können, wird zudem das Lösungskonzept des Nash-Gleichgewichts gemäß Definition 5.2 herangezogen⁸⁹.

Definition 5.1 (Eingangsaффines Nicht-Nullsummen-Differenzialspiel mit unbekanntem Gegenspielern)

Ein eingangsaффines Nicht-Nullsummen-Differenzialspiel mit unbekanntem Gegenspielern sei charakterisiert durch:

1. Eine eingangsaффine Systemdynamik

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \sum_{i=1}^N \mathbf{g}_i(\mathbf{x})\mathbf{u}_i =: \mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\mathbf{u} =: \mathbf{f} + \mathbf{g}\mathbf{u}, \quad (5.1)$$

$$\mathbf{x}(0) = \mathbf{x}_0,$$

mit dem Systemzustand $\mathbf{x} \in \mathbb{R}^n$, dem Anfangszustand \mathbf{x}_0 und den Stellgrößen $\mathbf{u}_i \in \mathbb{R}^{p_i}$ der N Regler (auch Spieler genannt), $i \in \mathcal{N} = \{1, 2, \dots, N\}$, $N \in \mathbb{N}_{\geq 1}$, $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ und $\mathbf{g}_i : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times p_i}$. Die Kurzschreibweisen $\mathbf{g}(\mathbf{x}) := [\mathbf{g}_1(\mathbf{x}) \quad \mathbf{g}_2(\mathbf{x}) \quad \dots \quad \mathbf{g}_N(\mathbf{x})]$, $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times p}$, $p := \sum_{i=1}^N p_i$ und $\mathbf{u} := [\mathbf{u}_1^\top \quad \mathbf{u}_2^\top \quad \dots \quad \mathbf{u}_N^\top]^\top \in \mathbb{R}^p$ werden zur Vereinfachung der Notation verwendet. Der Initialzeitpunkt t_0 wird o. B. d. A. zu $t_0 = 0$ gesetzt. Zudem seien $\mathbf{f}(\mathbf{x})$ und $\mathbf{g}(\mathbf{x})$ Lipschitz-stetig auf der kompakten Menge $\mathcal{X} \subset \mathbb{R}^n$, die den Ursprung enthält. Außerdem gelte $\mathbf{f}(\mathbf{0}) = \mathbf{0}$. Das System sei stabilisierbar auf \mathcal{X} .

2. Die Gütemaße

$$\begin{aligned} J_i(\mathbf{x}_0, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_N) &= \int_0^\infty q_i(\mathbf{x}(\tau)) + \sum_{j=1}^N \boldsymbol{\mu}_j^\top(\mathbf{x}(\tau)) \mathbf{R}_{ij} \boldsymbol{\mu}_j(\mathbf{x}(\tau)) \, d\tau \\ &=: \int_0^\infty r_i(\mathbf{x}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_N) \, d\tau, \end{aligned} \quad (5.2)$$

wobei $q_i : \mathbb{R}^n \rightarrow \mathbb{R}$ eine positiv definite Funktion (im Sinne von [NA05, S. 53]) darstellt, und $\mathbf{R}_{ij} \succeq \mathbf{0}$, $\forall j \in \mathcal{N}$, $i \neq j$, und $\mathbf{R}_{ii} = \mathbf{R}_{ii}^\top \succ \mathbf{0}$ gilt.

3. Das Ziel jedes Reglers $i \in \mathcal{N}$, ein Regelgesetz $\boldsymbol{\mu}_i^*(\mathbf{x})$ zu finden, wobei $\{\boldsymbol{\mu}_1^*, \dots, \boldsymbol{\mu}_N^*\}$ eine Feedback-Nash-Lösung (vgl. Definition 5.2) darstellt und aus Sicht von Regler i

- die Systemdynamik bekannt ist,
- Messungen von $\mathbf{x}(t)$ und $\boldsymbol{\mu}_j(t)$, $j \in \mathcal{N}$, verfügbar sind⁹⁰
- und Gütemaße J_j und Regelgesetze $\boldsymbol{\mu}_j(\mathbf{x})$ der anderen Spieler, d. h. $\forall j \neq i$, unbekannt sind.

⁸⁹ Da durch $\boldsymbol{\mu}_i$, $i \in \{1, 2, \dots, N\}$, im Folgenden Zustandsrückführungen bezeichnet werden, handelt es sich hierbei um ein sogenanntes Feedback-Nash-Gleichgewicht (vgl. [BO99, Abschnitt 6.5.2]).

⁹⁰ Die Notation $\boldsymbol{\mu}_j(t)$ stellt Zeitsignale dar, im Gegensatz zu Regelgesetzen $\boldsymbol{\mu}_j(\mathbf{x})$.

Definition 5.2 (Nash-Gleichgewicht [BO99, Definition 4.1])

Ein N -Tupel $\{\boldsymbol{\mu}_1^*, \dots, \boldsymbol{\mu}_N^*\}$ stellt ein Nash-Gleichgewicht eines N -Spieler-Nicht-Nullsummen-Differenzialspiels mit unendlichem Optimierungshorizont dar, wenn $\forall i \in \mathcal{N}, \mathcal{N} := \{1, 2, \dots, N\}$,

$$J_i(\mathbf{x}_0, \boldsymbol{\mu}_1^*, \dots, \boldsymbol{\mu}_i^*, \dots, \boldsymbol{\mu}_N^*) \leq J_i(\mathbf{x}_0, \boldsymbol{\mu}_1^*, \dots, \boldsymbol{\mu}_i, \dots, \boldsymbol{\mu}_N^*) \quad (5.3)$$

gilt.

Da in diesem Kapitel der Fokus auf der Analyse und Erfüllung der PE-Bedingung liegt, wird im Folgenden der Regulationsfall behandelt. Die Erweiterbarkeit des Differenzialspiels nach Definition 5.1 auf den Solltrajektorienfolgeregelungsfall wird jedoch in der nachfolgenden Bemerkung skizziert.

Bemerkung 5.1

ADP-kompatible Solltrajektorien nach Definition 4.1, wie beispielsweise die in Abschnitt 4.2.1 vorgestellte zeitkontinuierliche Solltrajektorienendarstellung, können in diesem Differenzialspiel nach Definition 5.1 prinzipiell ebenfalls verwendet werden. Hierzu wird anstelle des Systemzustands $\mathbf{x}(t)$ der erweiterte Systemzustand $\tilde{\mathbf{x}}(t) := [\mathbf{x}^\top(t) \quad \boldsymbol{\zeta}^\top(t)]^\top$ mit $\dot{\boldsymbol{\zeta}}(t) = \mathbf{D}\boldsymbol{\zeta}(t)$ (vgl. (4.19)) bzw. $\dot{\boldsymbol{\zeta}}(t) = \mathbf{D}'\boldsymbol{\zeta}(t)$ (vgl. (4.48)) betrachtet. Neben einer geeigneten Anregung von $\mathbf{x}(t)$, die grundsätzlich, ähnlich wie im vorliegenden Kapitel, über die Stellgröße $\mathbf{u}(t)$ erfolgen kann, muss dabei zusätzlich $\boldsymbol{\zeta}(t)$ angeregt werden. Dies kann, insbesondere bei Off-Policy-Algorithmen, prinzipiell über die Wahl von $\boldsymbol{\zeta}(t)$ während des Trainingsvorgangs erfolgen (vgl. auch Abbildung 4.3, Abbildung 6.4 und (6.33) für die Anregung von Solltrajektorienparametern $\boldsymbol{\zeta}$).

5.2 Policy Iteration für Nicht-Nullsummen-Differenzialspiele

Für die Lösung des Nicht-Nullsummen-Differenzialspiels nach Definition 5.1 soll im Folgenden der Policy-Iteration-Algorithmus für Differenzialspiele (vgl. [SLW17, Algorithm 1], [LLW14, Algorithm 1], [VL11, Algorithm 1]) genutzt werden. Ausgehend von den Gütefunktionalen (5.2) wird zunächst ein Maß definiert, das die akkumulierten Kosten beschreibt, die dem Regler i entstehen, wenn sich das System zum Zeitpunkt t im Zustand $\mathbf{x}(t)$ befindet und die Regelgesetze $\boldsymbol{\mu}(\mathbf{x}) = [\boldsymbol{\mu}_1^\top(\mathbf{x}) \quad \dots \quad \boldsymbol{\mu}_N^\top(\mathbf{x})]^\top$ angewandt werden. Damit die Gesamtkosten endlich sind und das gesuchte Maß existiert, ist die nachfolgende Definition erforderlich.

Definition 5.3 (Zulässige Regelgesetze⁹¹ [VL11, Definition 1])

Die Regelgesetze $\mu_i(\mathbf{x}), \forall i \in \mathcal{N}$, werden als zulässig bezüglich (5.2) auf der Menge \mathcal{X} bezeichnet ($\boldsymbol{\mu}(\mathbf{x}) = [\boldsymbol{\mu}_1^\top(\mathbf{x}) \ \cdots \ \boldsymbol{\mu}_N^\top(\mathbf{x})]^\top \in \Psi(\mathcal{X})$), wenn gilt:

1. $\mu_i(\mathbf{x}) \in \mathcal{C}^0(\mathcal{X})$,
2. $\mu_i(\mathbf{0}) = \mathbf{0}$,
3. $\mu_i(\mathbf{x})$ das System (5.1) auf \mathcal{X} stabilisiert und
4. die Gütemaße J_i nach (5.2) $\forall \mathbf{x}_0 \in \mathcal{X}$ endlich sind.

Für $\boldsymbol{\mu}(\mathbf{x}) \in \Psi(\mathcal{X})$ sind die Value Functions $V_i^\mu, i \in \mathcal{N}$, dann analog zu (2.12) durch

$$V_i^\mu(\mathbf{x}) := V_i^\mu(\mathbf{x}(t)) = \int_t^\infty r_i(\mathbf{x}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_N) d\tau \quad (5.4)$$

definiert. Zudem gelte Annahme 5.1.

Annahme 5.1

Sei $V_i^\mu(\mathbf{x}) \in \mathcal{C}^1(\mathcal{X})$.

Aus (5.4) ergeben sich die Lyapunov-Gleichungen

$$0 = r_i(\mathbf{x}, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_N) + (\nabla_{\mathbf{x}} V_i^\mu(\mathbf{x}))^\top (\mathbf{f} + \mathbf{g}\boldsymbol{\mu}) \quad (5.5)$$

(vgl. [VL11, LLW14]), bzw. für $\boldsymbol{\mu} = \boldsymbol{\mu}^*$ und $V_i^\mu(\mathbf{x}) = V_i^{\mu^*}(\mathbf{x}) =: V_i^*(\mathbf{x})$ die HJB-Gleichungen

$$0 = r_i(\mathbf{x}, \boldsymbol{\mu}_1^*, \dots, \boldsymbol{\mu}_N^*) + (\nabla_{\mathbf{x}} V_i^*(\mathbf{x}))^\top (\mathbf{f} + \mathbf{g}\boldsymbol{\mu}^*). \quad (5.6)$$

Mithilfe der Definition der Hamilton-Funktion⁹²

$$H_i(\mathbf{x}, \nabla_{\mathbf{x}} V_i, \boldsymbol{\mu}) = r_i(\mathbf{x}, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_N) + (\nabla_{\mathbf{x}} V_i(\mathbf{x}))^\top (\mathbf{f} + \mathbf{g}\boldsymbol{\mu}) \quad (5.7)$$

lässt sich mit Algorithmus 5.1 ein Lösungsalgorithmus in Form einer Policy Iteration für das Differenzialspiel gemäß Definition 5.1 angeben. Dies folgt aus Satz 5.1 (vgl. [SLW17, Theorem 2]), wobei die Konvergenz gegen $V_i^*, \forall i \in \mathcal{N}$, gerade der Nash-Lösung entspricht (vgl. [KKD14, (3)]). Die Policy Iteration kann entweder fortgeführt werden oder durch eine Abbruchbedingung, beispielsweise bei Konvergenz von $V_i^{[l]}, \forall i \in \mathcal{N}$, beendet werden.

Satz 5.1 (Konvergenz der Policy Iteration)

Es seien $\mu_i^{[l]}(\mathbf{x})$ und $V_i^{[l]}(\mathbf{x}), i \in \mathcal{N}$, wie in Algorithmus 5.1 gegeben und \mathcal{X} kompakt. Dann konvergiert die Folge $V_i^{[l]}$ für $l \rightarrow \infty$ gegen V_i^* .

⁹¹ (engl.): *admissible policies*.

⁹² Gleichung (5.7) ist als Funktion mit beliebigem V_i und $\boldsymbol{\mu}$ zu verstehen, es muss also nicht notwendigerweise $V_i = V_i^\mu$ gelten.

Beweis:Siehe [SLW17, Theorem 2]. □

Algorithmus 5.1 Policy Iteration für Differenzialspiele ([SLW17, Algorithm 1], [LLW14, Algorithm 1], [VL11, Algorithm 1])

1: **Initialisiere** Iterationsindex $l := 0$, Initialregler $\boldsymbol{\mu}^{[0]} \in \Psi(\mathcal{X})$.**Schritt 1** (Policy Evaluation):2: $\forall i \in \mathcal{N}$, finde $V_i^{[l+1]}$, sodass gilt:

$$0 = r_i(\mathbf{x}, \boldsymbol{\mu}_1^{[l]}, \dots, \boldsymbol{\mu}_N^{[l]}) + \left(\nabla_{\mathbf{x}} V_i^{[l+1]}(\mathbf{x}) \right)^\top (\mathbf{f} + \mathbf{g}\boldsymbol{\mu}^{[l]}), \quad (5.8a)$$

$$V_i^{[l+1]}(\mathbf{0}) = 0. \quad (5.8b)$$

Schritt 2 (Policy Improvement):3: $\forall i \in \mathcal{N}$, aktualisiere das Regelgesetz

$$\begin{aligned} \boldsymbol{\mu}_i^{[l+1]}(\mathbf{x}) &= \arg \min_{\boldsymbol{\mu}_i(\mathbf{x}) \in \Psi(\mathcal{X})} H_i(\mathbf{x}, \nabla_{\mathbf{x}} V_i^{[l+1]}(\mathbf{x}), \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_N) \\ &= -\frac{1}{2} \mathbf{R}_{ii}^{-1} \mathbf{g}_i^\top(\mathbf{x}) \nabla_{\mathbf{x}} V_i^{[l+1]}(\mathbf{x}) \end{aligned} \quad (5.9)$$

4: und setze $l := l + 1$. Gehe zu Schritt 1.

5.3 Funktionsapproximation und Anregungsbedingung

Das Lösen der nichtlinearen Lyapunov-Gleichungen (5.8) im Policy-Evaluation-Schritt in Algorithmus 5.1 stellt im Allgemeinen eine nicht-triviale Herausforderung dar (vgl. [VL11], [LLW14]). Die Value Function wird daher durch lineare Funktionsapproximation (siehe Abschnitt 2.1.3 sowie [BBdE10], [VL11] und [LLW14]) approximiert, d. h.

$$V_i^{[l]}(\mathbf{x}) = \mathbf{w}_i^{[l]\top} \boldsymbol{\phi}_i(\mathbf{x}) + \epsilon_i^{[l]}(\mathbf{x}) \quad (5.10)$$

mit den Basisfunktionen $\boldsymbol{\phi}_i : \mathbb{R}^n \rightarrow \mathbb{R}^{h_i}$, beschränkten Gewichtungsvektoren $\mathbf{w}_i^{[l]} \in \mathbb{R}^{h_i}$ und beschränkten Approximationsfehlern $\epsilon_i^{[l]}(\mathbf{x}) \in \mathcal{C}^1(\mathcal{X})$, die im Allgemeinen durch die endliche Dimension h_i der Funktionsapproximatoren $\boldsymbol{\phi}_i$ resultieren können (vgl. [VL10]). Zudem sei

$$V_i^*(\mathbf{x}) = \mathbf{w}_i^{*\top} \boldsymbol{\phi}_i(\mathbf{x}) + \epsilon_i(\mathbf{x}), \quad (5.11)$$

mit konstanten und beschränkten optimalen⁹³ Gewichtungsvektoren $\mathbf{w}_i^* \in \mathbb{R}^{h_i}$ und beschränkten Approximationsfehlern $\epsilon_i(\mathbf{x})$. Um die Notation übersichtlicher zu gestalten, wird zunächst die folgende Annahme getroffen⁹⁴.

Annahme 5.2

Seien $\epsilon_i^{[l]}(\mathbf{x}) = 0, \forall l, \forall i \in \mathcal{N}$, und somit auch $\epsilon_i(\mathbf{x}) = 0$.

Da die optimalen Gewichte \mathbf{w}_i^* des Differenzialspiels und auch die optimalen Gewichte $\mathbf{w}_i^{[l+1]}$ des Policy-Evaluation-Schrittes (5.8) a priori unbekannt sind, werden die geschätzten Value Functions zu

$$\hat{V}_i^{[l]}(\mathbf{x}) = \hat{\mathbf{w}}_i^{[l]\top} \phi_i(\mathbf{x}) \quad (5.12)$$

mit den geschätzten Gewichtungsvektoren $\hat{\mathbf{w}}_i^{[l]} \in \mathbb{R}^{h_i}$ definiert. Ziel der ADP-basierten Regler ist somit, die Gewichte $\hat{\mathbf{w}}_i^{[l]}$ anzupassen, um im Policy-Evaluation-Schritt (5.8) die Lösungen $\mathbf{w}_i^{[l+1]}$ und letztlich die Nash-Lösungen

$$\lim_{l \rightarrow \infty} \mathbf{w}_i^{[l]} = \mathbf{w}_i^* \quad (5.13)$$

(vgl. Satz 5.1) zu bestimmen. Solange jedoch $\hat{\mathbf{w}}_i^{[l]} \neq \mathbf{w}_i^{[l]}$ gilt, ist im Allgemeinen $\hat{V}_i^{[l]} \neq V_i^{[l]}$. Daher wird im Folgenden ein Gradientenabstiegsverfahren zur Adaption von $\hat{\mathbf{w}}_i^{[l]}$ im Policy-Evaluation-Schritt analysiert. Nachfolgend bezeichnet $\hat{\boldsymbol{\mu}}^{[l]}(\mathbf{x}) = \left[\hat{\boldsymbol{\mu}}_1^{[l]\top}(\mathbf{x}) \quad \dots \quad \hat{\boldsymbol{\mu}}_N^{[l]\top}(\mathbf{x}) \right]^\top$ die geschätzten Regelgesetze in der l -ten Iteration der Policy Iteration.

Die Hamilton-Funktion (5.7) kann als Maß zur Fehlerbeschreibung des Policy-Evaluation-Schrittes (5.8) betrachtet werden. Es gilt

$$\begin{aligned} H_i(\mathbf{x}, \nabla_{\mathbf{x}} \hat{V}_i^{[l+1]}, \hat{\boldsymbol{\mu}}^{[l]}) &= r_i(\mathbf{x}, \hat{\boldsymbol{\mu}}_1^{[l]}, \dots, \hat{\boldsymbol{\mu}}_N^{[l]}) + \left(\nabla_{\mathbf{x}} \hat{V}_i^{[l+1]} \right)^\top \left(\mathbf{f} + \mathbf{g} \hat{\boldsymbol{\mu}}^{[l]} \right) \\ &= r_i(\mathbf{x}, \hat{\boldsymbol{\mu}}_1^{[l]}, \dots, \hat{\boldsymbol{\mu}}_N^{[l]}) \\ &\quad + \hat{\mathbf{w}}_i^{[l+1]\top} \frac{\partial \phi_i(\mathbf{x})}{\partial \mathbf{x}} \left(\mathbf{f}(\mathbf{x}) + \sum_{j=1}^N \mathbf{g}_j(\mathbf{x}) \hat{\boldsymbol{\mu}}_j^{[l]} \right) \\ &=: r_i + \hat{\mathbf{w}}_i^{[l+1]\top} \boldsymbol{\sigma}_i, \end{aligned} \quad (5.14)$$

⁹³ Als optimale Gewichte werden im Rahmen dieses Kapitels Gewichte bezeichnet, die eine Lösung im Sinne des hier gewählten Lösungskonzepts (Nash-Gleichgewicht nach Definition 5.2) für das durch (5.1) und (5.2) beschriebene Differenzialspiel darstellen.

⁹⁴ An entsprechenden Stellen werden Hinweise auf den Einfluss von $\epsilon_i^{[l]}(\mathbf{x}) \neq 0$ gegeben (vgl. beispielsweise Bemerkung 5.2). Zudem sei für eine Diskussion des Einflusses nicht-exakter Funktionsapproximation auf [AKL05], [VL11] und [LLW14] verwiesen. So wird beispielsweise in [VL11, Proposition 2] gezeigt, dass für $\boldsymbol{\mu}(\mathbf{x}) \in \Psi(\mathcal{X})$ bei geeigneter Wahl der Basisfunktionen ϕ_i für $h_i \rightarrow \infty$ der Einfluss des Approximationsfehlers $\epsilon_i^{[l]}(\mathbf{x})$ verschwindet.

wobei die Kurzschreibweise $\sigma_i = \left. \frac{d\phi_i(\mathbf{x}(t))}{dt} \right|_{\mathbf{u}=\hat{\boldsymbol{\mu}}^{[l]}}$ eingeführt wurde. Die Minimierung des Quadrats der Hamilton-Funktion $\frac{1}{2} \left(H_i(\mathbf{x}, \nabla_{\mathbf{x}} \hat{V}_i^{[l+1]}, \hat{\boldsymbol{\mu}}^{[l]}) \right)^2$ bezüglich $\hat{\mathbf{w}}_i^{[l+1]}$ mittels Gradientenabstieg liefert die Adaptionsgesetze

$$\begin{aligned} \frac{d\hat{\mathbf{w}}_i^{[l+1]}}{dt} &= -\frac{\eta_i}{2} \left(\frac{\partial \left(H_i(\mathbf{x}, \nabla_{\mathbf{x}} \hat{V}_i^{[l+1]}, \hat{\boldsymbol{\mu}}^{[l]}) \right)^2}{\partial \hat{\mathbf{w}}_i^{[l+1]}} \right)^{\top} \\ &= -\eta_i H_i(\mathbf{x}, \nabla_{\mathbf{x}} \hat{V}_i^{[l+1]}, \hat{\boldsymbol{\mu}}^{[l]}) \boldsymbol{\sigma}_i \\ &= -\eta_i \left(r_i + \hat{\mathbf{w}}_i^{[l+1]\top} \boldsymbol{\sigma}_i \right) \boldsymbol{\sigma}_i \end{aligned} \quad (5.15)$$

mit den Lernraten $\eta_i \in \mathbb{R}_{>0}$.

Mit $V_i^{[l+1]} = V_i^{\hat{\boldsymbol{\mu}}^{[l]}}(\mathbf{x})$, d. h., wenn $V_i^{[l+1]}$ die Lösung des Policy-Evaluation-Schritts zum aktuellen Regelgesetz $\hat{\boldsymbol{\mu}}^{[l]}$ darstellt, gilt nach (5.5) und (5.10)

$$\begin{aligned} 0 &= r_i(\mathbf{x}, \hat{\boldsymbol{\mu}}_1^{[l]}, \dots, \hat{\boldsymbol{\mu}}_N^{[l]}) + \left(\nabla_{\mathbf{x}} V_i^{\hat{\boldsymbol{\mu}}^{[l]}}(\mathbf{x}) \right)^{\top} (\mathbf{f} + \mathbf{g}\hat{\boldsymbol{\mu}}^{[l]}) \\ &= r_i + \left(\nabla_{\mathbf{x}} \left(\mathbf{w}_i^{[l+1]\top} \phi_i(\mathbf{x}) + \epsilon_i^{[l+1]}(\mathbf{x}) \right) \right)^{\top} (\mathbf{f} + \mathbf{g}\hat{\boldsymbol{\mu}}^{[l]}) \\ &= r_i + \boldsymbol{\sigma}_i^{\top} \mathbf{w}_i^{[l+1]} + \left(\nabla_{\mathbf{x}} \epsilon_i^{[l+1]}(\mathbf{x}) \right)^{\top} (\mathbf{f} + \mathbf{g}\hat{\boldsymbol{\mu}}^{[l]}) \\ &=: r_i + \boldsymbol{\sigma}_i^{\top} \mathbf{w}_i^{[l+1]} + \epsilon_{\text{H},i}^{[l+1]}, \end{aligned} \quad (5.16)$$

wobei der Fehler $\epsilon_{\text{H},i}^{[l+1]} := \left(\nabla_{\mathbf{x}} \epsilon_i^{[l+1]}(\mathbf{x}) \right)^{\top} (\mathbf{f} + \mathbf{g}\hat{\boldsymbol{\mu}}^{[l]})$ verschwindet, wenn eine exakte Approximation möglich ist, also $\epsilon_i^{[l+1]}(\mathbf{x}) = 0$ nach Annahme 5.2 gilt.

Aus der Definition des Gewichtsfehlers $\tilde{\mathbf{w}}_i^{[l+1]} := \mathbf{w}_i^{[l+1]} - \hat{\mathbf{w}}_i^{[l+1]}$ ergibt sich die Fehlerdynamik bei Verwendung des Gradientenabstiegs (5.15) zu

$$\begin{aligned} \frac{d\tilde{\mathbf{w}}_i^{[l+1]}}{dt} &= -\frac{d\hat{\mathbf{w}}_i^{[l+1]}}{dt} = \eta_i \left(r_i + \hat{\mathbf{w}}_i^{[l+1]\top} \boldsymbol{\sigma}_i \right) \boldsymbol{\sigma}_i \\ &= \eta_i \boldsymbol{\sigma}_i \left(r_i + \boldsymbol{\sigma}_i^{\top} \left(\mathbf{w}_i^{[l+1]} - \tilde{\mathbf{w}}_i^{[l+1]} \right) \right) \\ &= -\eta_i \boldsymbol{\sigma}_i \boldsymbol{\sigma}_i^{\top} \tilde{\mathbf{w}}_i^{[l+1]} + \eta_i \boldsymbol{\sigma}_i \left(r_i + \boldsymbol{\sigma}_i^{\top} \mathbf{w}_i^{[l+1]} \right) \\ &= -\eta_i \boldsymbol{\sigma}_i \boldsymbol{\sigma}_i^{\top} \tilde{\mathbf{w}}_i^{[l+1]} - \eta_i \boldsymbol{\sigma}_i \epsilon_{\text{H},i}^{[l+1]}, \end{aligned} \quad (5.17)$$

wobei der letzte Schritt unter Verwendung von (5.16) resultiert.

Um die Konvergenzeigenschaften der zentralen Fehlerdynamik (5.17) zu analysieren, wird zunächst die folgende Definition, welche die Anregung eines Signals mathematisch formalisiert, benötigt.

Definition 5.4 (PE-Funktionen nach [NA05, Definition 6.2])

Eine stückweise stetig differenzierbare, beschränkte Funktion (vgl. [NA05, Definition 6.1]) $\sigma : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^h$ wird als *persistently excited (PE)* für alle $t \geq t_0$ ⁹⁵ bezeichnet, wenn es positive Konstanten $\alpha, \alpha_1, T \in \mathbb{R}_{>0}$ gibt, sodass die äquivalenten Beziehungen⁹⁶

$$\int_t^{t+T} \sigma(\tau) \sigma^\top(\tau) d\tau \succeq \alpha \mathbf{I} \quad (5.18)$$

und

$$\frac{1}{T} \int_t^{t+T} |(\sigma(\tau))^\top \mathbf{e}| d\tau \geq \alpha_1 \quad (5.19)$$

für alle $t \geq t_0$ und jeden Einheitsvektor $\mathbf{e} \in \mathbb{R}^h$ gelten. Die Konstanten α bzw. α_1 werden als Grad der PE-Bedingung bezeichnet.

Basierend auf den Ergebnissen aus [NA05], [VL10] und [MN77] lassen sich damit die folgenden Konvergenzeigenschaften formulieren.

Lemma 5.1 (Konvergenz der Critic-Gewichte)

Seien $\{\hat{\mu}_1^{[l]}, \dots, \hat{\mu}_N^{[l]}\}$ zeitinvariante, nach Definition 5.3 zulässige Regelgesetze und Annahme 5.2 erfüllt. Die Critic-Gewichte $\hat{w}_i^{[l+1]}$ zur Approximation der Lösung $V_i^{[l+1]}$ der Lyapunov-Gleichungen (5.8) (Policy-Evaluation-Schritt) gemäß (5.12) werden durch den Gradientenabstieg (5.15) angepasst. Dann folgt, dass die Gewichtsfehler

$$\tilde{w}_i^{[l+1]} = w_i^{[l+1]} - \hat{w}_i^{[l+1]} \quad (5.20)$$

genau dann exponentiell gegen null konvergieren, wenn σ_i PE nach Definition 5.4 ist. Dabei gilt $\forall k \in \mathbb{N}_{\geq 0}$

$$\left\| \tilde{w}_i^{[l+1]}(kT_i) \right\|_2 \leq \exp\left(\frac{k}{2} \ln \rho_i\right) \left\| \tilde{w}_i^{[l+1]}(0) \right\|_2 \quad (5.21)$$

mit

$$\rho_i := 1 - \frac{2T_i \eta_i \alpha_{1,i}^2}{(1 + T_i \eta_i \sigma_{i,\max}^2)^2}, \quad (5.22)$$

wobei $\sigma_{i,\max} \geq \|\sigma_i\|_2$.

⁹⁵ Der Zeitpunkt t_0 , ab dem die Funktion σ PE ist, ergibt sich üblicherweise aus dem Kontext. Im Folgenden sei $t_0 = 0$, sofern nicht anders definiert.

⁹⁶ Die Äquivalenz ergibt sich direkt aus [NA05, Theorem 2.16].

Beweis:

Unter Annahme 5.2 ist im Fall $\eta_i = 1$ nach [MN77, Theorem 1] (5.19) notwendig und hinreichend für die gleichmäßige asymptotische Stabilität der Ruhelage von (5.17). Nun ist laut (5.18) σ_i mit α_i und T_i genau dann PE, wenn $\bar{\sigma}_i = \sqrt{\eta_i} \sigma_i$ PE mit $\bar{\alpha}_i = \eta_i \alpha_i$ und $T_i, \forall \eta_i > 0$, ist. Zudem ist $\bar{\sigma}_i$ genau dann PE, wenn

$$\frac{d\tilde{w}_i^{[l+1]}}{dt} = -\eta_i \sigma_i \sigma_i^\top \tilde{w}_i^{[l+1]} \quad (5.23)$$

gleichmäßig asymptotisch stabil ist (vgl. [MN77, Theorem 1]). Aufgrund der Linearität von (5.23) ist gleichmäßige asymptotische Stabilität identisch zu exponentieller Stabilität, weshalb σ_i genau dann PE ist, falls (5.23) exponentiell stabil ist. Die Abschätzung der Konvergenzgeschwindigkeit ρ_i in (5.22), welche die exponentielle Abnahme der Gewichtsfehler nach (5.21) charakterisiert, wird analog zu [NA05, Theorem 2.16] bestimmt (vgl. [VL10, Technical Lemma 2]). \square

Bemerkung 5.2

Im Fall einer nicht exakten Approximierbarkeit der gesuchten Value Function $V_i^{[l+1]}$ durch $w_i^{[l+1]}$, d. h., falls Annahme 5.2 nicht erfüllt ist, sondern beschränkte Approximationsfehler $\epsilon_i^{[l+1]} \leq \bar{\epsilon}_i^{[l+1]}$ bestehen bleiben, resultiert exponentielle Konvergenz von $\tilde{w}_i^{[l+1]}$ gegen die Nachbarschaft von null [VL10, Technical Lemma 2].

Die Policy Evaluation (5.8) kann somit durch die Adaption gemäß (5.15) bis zur Konvergenz von $\hat{w}_i^{[l+1]}$, $\forall i \in \mathcal{N}$, erfolgen. Aus den approximierten Value Functions $\hat{V}_i^{[l+1]}$ nach Konvergenz des Gradientenabstiegsverfahrens gemäß (5.15) folgen nach (5.9) im Policy-Improvement-Schritt die approximierten Regelgesetze

$$\hat{\mu}_i^{[l+1]}(\mathbf{x}) = -\frac{1}{2} \mathbf{R}_{ii}^{-1} \mathbf{g}_i^\top(\mathbf{x}) \left(\frac{\partial \phi_i(\mathbf{x})}{\partial \mathbf{x}} \right)^\top \hat{w}_i^{[l+1]}. \quad (5.24)$$

Aufgrund des analytischen Zusammenhangs der approximierten Regelgesetze zu den Critic-Gewichten w_i sind nach Liu et al. [LLW14] somit keine gesonderten Funktionsapproximatoren zur Beschreibung der Regelgesetze (sogenannte *Actors*, vgl. Abschnitt 2.1.3) notwendig.

Damit stellt die Policy Iteration nach Algorithmus 5.1 ein Lösungsverfahren gemäß Definition 5.1 dar⁹⁷. Da Lemma 5.1 belegt, dass die Erfüllung der PE-Bedingung nach Definition 5.4 notwendig und hinreichend für die exponentielle Konvergenz der durch $\tilde{w}_i^{[l+1]} = w_i^{[l+1]} - \hat{w}_i^{[l+1]}$, $i \in \mathcal{N}$, gegebenen Critic-Gewichtsfehler im Policy-Evaluation-Schritt unter Verwendung eines Gradientenabstiegs nach (5.15) ist, untersucht der Rest dieses Kapitels Bedingungen zur Erfüllung der PE-Eigenschaft von σ_i .

⁹⁷ Im Fall einer nicht exakten Approximation der Value Function folgt neben der Konvergenz der Policy Evaluation gegen die Nachbarschaft von $w_i^{[l+1]}$ (vgl. Bemerkung 5.2) auch Konvergenz der Policy Iteration gegen die Nachbarschaft der optimalen Lösung (vgl. [AKL05, Theorem 4], [VL11]).

5.4 Hinreichende Bedingungen zur Erfüllung der PE-Eigenschaft in ADP-basierten Differenzialspielen

In diesem Abschnitt werden für den im ADP-Kontext häufig verwendeten Fall polynomieller Funktionsapproximatoren Bedingungen an den Systemzustand $\mathbf{x}(t)$ hergeleitet, die garantieren, dass σ_i die PE-Eigenschaft erfüllt. Diese neuen hinreichenden Bedingungen können in ADP-basierten Problemen genutzt werden und bilden die theoretische Basis, um in Abschnitt 5.6 beispielhafte Anregungssignale zu entwerfen. Formal seien die Basisfunktionen $\phi_i(\mathbf{x})$ wie folgt beschrieben.

Annahme 5.3 (Basisfunktionen $\phi_i(\mathbf{x})$)

Die Elemente der Basisfunktionen $\phi_i(\mathbf{x})$, $i \in \mathcal{N}$, seien $\forall \bar{h} \in \{1, \dots, h\}$ gemäß

$$\begin{aligned} \phi_{i,\bar{h}}(\mathbf{x}) &= \bar{f}_{i,\bar{h}} \prod_{j=1}^n x_j^{f_{i,j,\bar{h}}}, \\ \bar{f}_{i,\bar{h}} &\in \mathbb{R}_{\neq 0}, \\ f_{i,j,\bar{h}} &\in \mathbb{N}_{\geq 0} \end{aligned} \tag{5.25}$$

gewählt, sodass $\exists j \in \{1, \dots, n\} : f_{i,j,\bar{h}} \neq 0$ und $\bar{f}_{i,\bar{h}} < \infty$.

Um das Signal \mathbf{x} formal mit der PE-Eigenschaft der transformierten Signale $\sigma_i = \dot{\phi}_i(\mathbf{x})$ in Zusammenhang zu bringen, wird der Begriff *sufficiently rich* (SR)⁹⁸ definiert.

Definition 5.5 (SR-Signal)

Sei $\mathcal{S}\{\cdot\}$ die Transformation eines Signals, die aus der Zusammenschaltung von Integratoren, Differenzialatoren und nichtlinearen, Lipschitz-stetigen Funktionen bestehen kann. Das Signal $\mathbf{x}(t) : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^n$ wird als *sufficiently rich* (SR) bezüglich \mathcal{S} bezeichnet, wenn $\mathcal{S}\{\mathbf{x}(t)\}$ PE ist.

Somit ist σ_i PE, wenn $\mathbf{x}(t)$ SR bezüglich $\dot{\phi}_i(\mathbf{x})$ ist. Daraus resultiert die folgende Problemstellung für dieses Unterkapitel.

Problem 5.1

Sei $\phi_i(\mathbf{x}) \forall i \in \mathcal{N}$ wie in Annahme 5.3. Gesucht ist $\mathbf{x}(t)$, sodass $\forall i \in \mathcal{N}$ das durch $\sigma_i = \dot{\phi}_i(\mathbf{x}(t))$ gegebene Signal die PE-Eigenschaft nach Definition 5.4 erfüllt, d. h. $\mathbf{x}(t)$ SR bezüglich $\dot{\phi}_i(\mathbf{x})$ ist.

⁹⁸ Boyd und Sastry [BS86] definieren den Begriff der SR-Signale für lineare Systeme (im Zusammenhang des *Model Reference Adaptive Control*) mit skalarem Eingang. Die in der vorliegenden Arbeit vorgestellte Definition kann daher als Verallgemeinerung auf den mehrdimensionalen und nichtlinearen Fall verstanden werden.

Da das zur Adaption verwendete Signal $\sigma_i = \dot{\phi}_i(\mathbf{x})$ aus nichtlinearen Basisfunktionen $\phi_i(\mathbf{x})$ resultiert⁹⁹, kann diese Transformation von $\mathbf{x}(t)$ nicht mithilfe klassischer Eigenschaften linearer Transformationen aus der PE-Literatur [NA05], [NA87] analysiert werden. So kann, abhängig von $\phi_i(\mathbf{x})$, das Signal σ_i PE sein, obwohl \mathbf{x} nicht PE ist, oder umgekehrt [LK98], [LK99]. Dies ist in Beispiel 5.1 illustriert.

Beispiel 5.1

1. \mathbf{x} ist PE $\not\Rightarrow$ σ_i ist PE

Mit der PE-Funktion

$$\mathbf{x}(t) = \begin{bmatrix} \frac{1}{2} \cos(2\omega t) \\ \cos(\omega t) \end{bmatrix}, \quad \omega \neq 0, \quad \text{und} \quad \phi_i(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2^2 \end{bmatrix} \quad (5.26)$$

folgt, dass

$$\sigma_i = \frac{d}{dt} \begin{bmatrix} x_1 \\ x_2^2 \end{bmatrix} = \begin{bmatrix} -\omega \sin(2\omega t) \\ -\omega \sin(2\omega t) \end{bmatrix} \quad (5.27)$$

keine PE-Funktion ist.

2. \mathbf{x} ist PE \Leftarrow σ_i ist PE

Zwar stellt

$$\mathbf{x}(t) = \begin{bmatrix} \sin(\omega t) \\ 0 \end{bmatrix}, \quad \omega \neq 0, \quad (5.28)$$

keine PE-Funktion dar, aus

$$\phi_i(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2^2 \end{bmatrix} \quad \text{und} \quad \sigma_i = \frac{d}{dt} \begin{bmatrix} x_1 \\ x_2^2 \end{bmatrix} = \begin{bmatrix} \omega \cos(\omega t) \\ \omega \sin(2\omega t) \end{bmatrix} \quad (5.29)$$

folgt jedoch, dass σ_i eine PE-Funktion ist.

Somit ist die Erfüllung der PE-Eigenschaft von \mathbf{x} weder notwendig noch hinreichend dafür, dass σ_i PE ist. Beispiel 5.1 veranschaulicht, dass durch polynomielle Basisfunktionen $\phi_i(\mathbf{x})$ aufgrund deren nichtlinearen Charakters im Vergleich zu \mathbf{x} sowohl Frequenzen ausgelöscht, als auch zusätzliche Frequenzen erzeugt werden können. Motiviert durch diese Erkenntnis ist die zentrale Idee im Folgenden, diese Frequenzen zu analysieren. Um Trajektorien $\mathbf{x}(t)$ zu finden, die eine Lösung von Problem 5.1 darstellen, wird zunächst die folgende allgemeine Struktur einer Anregungstrajektorie $\mathbf{x}_{PE}(t)$ angesetzt.

⁹⁹ Selbst für den Spezialfall eines linear-quadratischen Differenzialspiels mit linearer Systemdynamik und quadratischen Gütefunktionalen ist $\phi_i(\mathbf{x})$ nichtlinear.

Annahme 5.4 (Anregungstrajektorie $x_{\text{PE}}(t)$)

Die Anregungstrajektorie $x_{\text{PE}}(t) = [x_{\text{PE},1}(t), \dots, x_{\text{PE},n}(t)]^T$ des Systemzustands $x(t)$ wird $\forall o \in \{1, \dots, n\}$ zu

$$x_{\text{PE},o} = \sum_{j=1}^m g_{j,o} \sin(\omega_j t) + \sum_{j=1}^m \bar{g}_{j,o} \cos(\omega_j t), \quad (5.30)$$

$g_{j,o}, \bar{g}_{j,o} \in \mathbb{R}$, gewählt, sodass

1. $\forall o \in \{1, \dots, n\} \exists j \in \{1, \dots, m\} : g_{j,o} \neq 0$ oder $\bar{g}_{j,o} \neq 0$,
2. $|g_{j,o}| < \infty$ und $|\bar{g}_{j,o}| < \infty$ und
3. $1 \leq m < \infty$.

Hierbei sind $\omega_1, \dots, \omega_m$ Frequenzvariablen, die den Frequenzvektor $\omega = [\omega_1 \ \dots \ \omega_m]^T \in \mathbb{R}^m$ bilden.

In den nachfolgenden Schritten wird der Spielerindex i o. B. d. A. vernachlässigt, um die Notation übersichtlicher zu gestalten. Zum Ende des Abschnitts werden schließlich wieder alle N Spieler berücksichtigt. Der Frequenzvektor ω soll nun so gewählt werden, dass $x_{\text{PE}}(t)$ SR bezüglich $\dot{\phi}(x)$ ist. Unter Annahme 5.3 und Annahme 5.4 resultiert das folgende Lemma, das die in $\phi(x_{\text{PE}})$ auftretenden Frequenzen beschreibt.

Lemma 5.2

Werden $\phi(x)$ und $x_{\text{PE}}(t)$ gemäß Annahme 5.3 und Annahme 5.4 gewählt, so ergeben sich die Elemente von $\phi(x_{\text{PE}})$ zu

$$\phi_{\bar{h}}(x_{\text{PE}}) = \sum_{l=1}^{L_{\bar{h}}} a_{l,\bar{h}} \sin(\omega_{l,\bar{h}} t) + \sum_{k=1}^{K_{\bar{h}}} c_{k,\bar{h}} \cos(\bar{\omega}_{k,\bar{h}} t) + e_{\bar{h}}, \quad (5.31)$$

$\forall \bar{h} \in \{1, \dots, h\}$, mit $a_{l,\bar{h}}, c_{k,\bar{h}} \in \mathbb{R}_{\neq 0}$, $e_{\bar{h}} \in \mathbb{R}$ und

$$\omega_{l,\bar{h}} = \sum_{j=1}^m b_{j,l,\bar{h}} \omega_j, \quad (5.32)$$

$$\bar{\omega}_{k,\bar{h}} = \sum_{j=1}^m d_{j,k,\bar{h}} \omega_j, \quad (5.33)$$

wobei $b_{j,l,\bar{h}}, d_{j,k,\bar{h}} \in \mathbb{R}$. Zudem gilt für die oberen Summengrenzen $L_{\bar{h}}, K_{\bar{h}} \in \mathbb{N}_{\geq 0}$

$$L_{\bar{h}} + K_{\bar{h}} \geq 1, \quad \forall \bar{h} \in \{1, \dots, h\}. \quad (5.34)$$

Beweis:

Der Beweis ist in Anhang C.1 skizziert. \square

Im Folgenden werden durch eine Analyse der Frequenzen $\underline{\omega}_{l,\bar{h}}$ nach (5.32) und $\bar{\omega}_{k,\bar{h}}$ nach (5.33), die aus der nichtlinearen Transformation $\dot{\phi}(\mathbf{x}_{\text{PE}})$ resultieren, ausschließende Bedingungen an die Frequenzen ω hergeleitet, sodass $\dot{\phi}(\mathbf{x}_{\text{PE}})$ PE ist. Hierbei wird das nachfolgende notwendige und hinreichende Kriterium für mehrdimensionale PE-Signale genutzt.

Lemma 5.3

Sei $\sigma(t) : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^h$. Das Signal $\sigma(t)$ ist genau dann PE, wenn $\alpha^\top \sigma(t)$ für jeden konstanten Vektor $\alpha \in \mathbb{R}^h$ mit $\alpha \neq \mathbf{0}$ mindestens eine Spektrallinie im Frequenzspektrum aufweist.

Beweis:

Das Signal $\sigma(t)$ ist genau dann PE, wenn $\alpha^\top \sigma(t)$ für jeden konstanten Vektor $\alpha \in \mathbb{R}^h$, $\alpha \neq \mathbf{0}$, PE ist [NA05, Lemma 6.2]. Des Weiteren ist ein skalares Signal $\sigma_s(t) : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$ genau dann PE, wenn es mindestens eine Spektrallinie aufweist [NA05, Sublemma 6.1]. Mit $\sigma_s(t) := \alpha^\top \sigma(t)$ folgt somit die Aussage von Lemma 5.3. \square

Basierend auf Lemma 5.3 wird nun eine Menge Ω definiert, sodass für jeden Vektor $\omega \in \Omega$ die Erfüllung der PE-Bedingung von $\sigma = \dot{\phi}(\mathbf{x}_{\text{PE}})$ garantiert wird. Die nachfolgende Proposition charakterisiert diese Menge Ω .

Proposition 5.1

Seien $\phi(\mathbf{x})$ und \mathbf{x}_{PE} in Übereinstimmung zu Annahme 5.3 und Annahme 5.4 gewählt. Dann ist \mathbf{x}_{PE} SR bezüglich $\dot{\phi}(\mathbf{x}_{\text{PE}})$ und somit $\sigma = \dot{\phi}(\mathbf{x}_{\text{PE}})$ PE $\forall \omega \in \Omega$, wenn eine nichtleere Menge $\Omega \subseteq \mathbb{R}^m$ existiert, sodass $\forall \omega \in \Omega$ und jeden konstanten Vektor $\alpha \in \mathbb{R}^h$ mit $\alpha \neq \mathbf{0}$

$$\alpha^\top \sigma = \alpha^\top \dot{\phi}(\mathbf{x}_{\text{PE}}) = \sum_{l=1}^{L^{(\alpha)}} a_l^{(\alpha)} \cos(\underline{\omega}_l^{(\alpha)} t) + \sum_{k=1}^{K^{(\alpha)}} c_k^{(\alpha)} \sin(\bar{\omega}_k^{(\alpha)} t) \quad (5.35)$$

mit $L^{(\alpha)} + K^{(\alpha)} \geq 1$ und

$$\begin{aligned} a_l^{(\alpha)} &\neq 0, & \underline{\omega}_l^{(\alpha)} &\neq 0, & \underline{\omega}_{l_1}^{(\alpha)} &\neq \underline{\omega}_{l_2}^{(\alpha)}, & \underline{\omega}_{l_1}^{(\alpha)} &\neq -\underline{\omega}_{l_2}^{(\alpha)}, \\ c_k^{(\alpha)} &\neq 0, & \bar{\omega}_k^{(\alpha)} &\neq 0, & \bar{\omega}_{k_1}^{(\alpha)} &\neq \bar{\omega}_{k_2}^{(\alpha)}, & \bar{\omega}_{k_1}^{(\alpha)} &\neq -\bar{\omega}_{k_2}^{(\alpha)}, \end{aligned}$$

$\forall l_{1/2} \in \{1, \dots, L^{(\alpha)}\}$, $l_1 \neq l_2$ und $\forall k_{1/2} \in \{1, \dots, K^{(\alpha)}\}$, $k_1 \neq k_2$ gilt. Der Index $\langle \alpha \rangle$ kennzeichnet hierbei eine Abhängigkeit von α .

Beweis:

Das Signal auf der rechten Seite von (5.35) weist unter den in Proposition 5.1 gegebenen Bedingungen mindestens eine Spektrallinie im Frequenzspektrum auf. Dies wird, neben dem Auftreten mindestens einer Frequenz ($L^{(\alpha)} + K^{(\alpha)} \geq 1$), durch die lineare Unabhängigkeit von Sinusfunktionen (bzw. Kosinusfunktionen) unterschiedlicher Frequenzen sowie der linearen Unabhängigkeit einer Sinus- zur Kosinusfunktion derselben Frequenz gewährleistet. Da (5.35) für jeden konstanten Vektor $\alpha \in \mathbb{R}^h$, $\alpha \neq \mathbf{0}$, gilt, folgt mit Lemma 5.3, dass $\sigma = \dot{\phi}(\mathbf{x}_{\text{PE}})$ PE ist. \square

Um nun eine Anregungstrajektorie $\mathbf{x}_{\text{PE}}(t)$ zu finden, die SR bezüglich $\dot{\phi}(\mathbf{x}_{\text{PE}})$ ist, genügt es gemäß Proposition 5.1, eine nichtleere Menge Ω zu finden, die Bedingungen an die Frequenzen ω in \mathbf{x}_{PE} so formuliert, dass (5.35) erfüllt ist. Die analytischen Zusammenhänge (5.32) und (5.33) zwischen den Frequenzen ω_j in \mathbf{x}_{PE} und den durch die nichtlineare Transformation eingeführten Frequenzen $\underline{\omega}_{l,\bar{h}}$ und $\bar{\omega}_{k,\bar{h}}$ dienen hierzu im Folgenden als Ausgangspunkt. Dies lässt sich motivieren, da die zeitliche Ableitung von (5.31) nur den Typ der trigonometrischen Funktionen sowie deren Koeffizienten verändert. Die in $\sigma_{\bar{h}} = \dot{\phi}_{\bar{h}}(\mathbf{x}_{\text{PE}})$, $\forall \bar{h} \in \{1, \dots, h\}$, auftretenden Frequenzen folgen direkt aus (5.31). Durch die zeitliche Ableitung von $\phi_{\bar{h}}(\mathbf{x}_{\text{PE}})$ entfällt außerdem die Konstante $e_{\bar{h}}$ und muss daher nicht weiter berücksichtigt werden. Zudem sind die genauen Werte der Koeffizienten $a_{l,\bar{h}}$ und $c_{k,\bar{h}}$ in (5.31) irrelevant für die Erfüllung von (5.35). Dies gilt, da gemäß Proposition 5.1 aufgrund von $L^{(\alpha)} + K^{(\alpha)} \geq 1$ für mindestens eine Frequenz eine Amplitude ungleich null resultiert.

Im Folgenden wird nun eine mögliche Menge Ω an Frequenzbedingungen gemäß Proposition 5.1 berechnet. Hierzu werden in den Abschnitten 5.4.1 und 5.4.2 zunächst die Hilfsmengen $\Omega^{(1)}$ und $\Omega^{(2)}$ definiert und vereinfacht. In Abschnitt 5.4.3 folgt schließlich die Menge

$$\Omega = \Omega^{(1)} \cap \Omega^{(2)} \quad (5.36)$$

an Frequenzbedingungen sowie die theoretische Aussage, dass die durch diese Menge Ω gegebenen Bedingungen die Erfüllung von (5.35) sicherstellen.

5.4.1 Hilfsmenge $\Omega^{(1)}$

Die nachfolgende Definition formuliert durch die Hilfsmenge $\Omega^{(1)}$ Bedingungen an ω . Diese garantieren, dass sich, unabhängig von den exakten Werten der Koeffizienten und Argumente, für keinen Vektor $\omega \in \Omega^{(1)}$ die trigonometrischen Funktionen innerhalb eines Elements $\phi_{\bar{h}}(\mathbf{x}_{\text{PE}})$ und somit in $\sigma_{\bar{h}}$ kompensieren.

Definition 5.6 ($\Omega^{(1)}$: Einzigartigkeit der Frequenzen innerhalb jedes Elements $\phi_{\bar{h}}(\mathbf{x}_{\text{PE}})$, $\forall \bar{h} \in \{1, \dots, h\}$)

Die Menge $\Omega^{(1)}$ an Bedingungen an die Frequenzvariable ω sei durch

$$\Omega^{(1)} := \left\{ \omega : \bigwedge_{\bar{h}=1}^h \left(\bigwedge_{l_1=1}^{L_{\bar{h}}-1} \bigwedge_{l_2 > l_1}^{L_{\bar{h}}} (\omega_{l_1, \bar{h}} \neq \omega_{l_2, \bar{h}} \wedge \omega_{l_1, \bar{h}} \neq -\omega_{l_2, \bar{h}}) \right. \right. \\ \left. \left. \wedge \bigwedge_{k_1=1}^{K_{\bar{h}}-1} \bigwedge_{k_2 > k_1}^{K_{\bar{h}}} (\bar{\omega}_{k_1, \bar{h}} \neq \bar{\omega}_{k_2, \bar{h}} \wedge \bar{\omega}_{k_1, \bar{h}} \neq -\bar{\omega}_{k_2, \bar{h}}) \right) \right\} \quad (5.37)$$

definiert.

Somit sind für jedes $\omega \in \Omega^{(1)}$ nach Definition 5.6 in jedem Element $\phi_{\bar{h}}(\mathbf{x}_{\text{PE}})$, d. h. für jedes $\bar{h} \in \{1, \dots, h\}$, die Beträge der Frequenzen der Sinusterme einzigartig, Gleiches gilt für die Kosinusterme. Die Menge $\Omega^{(1)}$ kann wie folgt vereinfacht werden.

Lemma 5.4 (Vereinfachung der Menge $\Omega^{(1)}$)

Die durch Definition 5.6 beschriebene Menge $\Omega^{(1)}$ lässt sich zu

$$\Omega^{(1)} = \left\{ \omega : \bigwedge_{z=1}^{Z_1} \mathbf{c}_z^T \omega \neq 0 \right\}, \quad Z_1 \in \mathbb{N}_{\geq 0}, \quad (5.38)$$

umformulieren.

Beweis:

Ausgehend von (5.37) wird zunächst die Komplementärmenge $\bar{\Omega}^{(1)}$ gebildet. Anschließendes Anwenden der De Morganschen Gesetze [BSMM13, S. 332] führt auf eine Menge, die durch die Disjunktion linearer algebraischer Gleichungen der Form $\omega_{l_1, \bar{h}} = \omega_{l_2, \bar{h}}$ gegeben ist. Diese Gleichungen können in Matrixschreibweise $\bar{\mathbf{c}}_z^T \omega = 0$ mit $\bar{\mathbf{c}}_z \in \mathbb{R}^m$ überführt werden, wobei z einen Laufindex darstellt.

Durch Skalierung der Zeilenvektoren $\bar{\mathbf{c}}_z$ können diese Ausdrücke jeweils auf reduzierte Zeilenstufenform gebracht werden. Die Transformation der Koeffizientenmatrix eines homogenen linearen Gleichungssystems auf reduzierte Zeilenstufenform verändert dessen Lösungsmenge nicht [KB18, S. 22]. Dies gilt auch für $\bar{\mathbf{c}}_z^T \omega = 0$ und dessen reduzierte Zeilenstufenform $\mathbf{c}_z^T \omega = 0$. Werden anschließend alle mehrfach auftretenden Bedingungen eliminiert und abermals die De Morganschen Gesetze angewandt, so resultiert (5.38). Da diese mehrfach auftretenden Bedingungen die Menge $\Omega^{(1)}$ nicht verändern, ist $\Omega^{(1)}$ (5.37) identisch zu (5.38). \square

5.4.2 Hilfsmenge $\Omega^{(2)}$

Durch (5.35) in Proposition 5.1 werden Bedingungen an die gewichtete Summe der Elemente von $\sigma = \dot{\phi}(\mathbf{x}_{\text{PE}})$ formuliert. Daher genügt eine separierte Betrachtung der Elemente $\phi_{\bar{h}}(\mathbf{x}_{\text{PE}})$, $\bar{h} \in \{1, \dots, h\}$, nicht und die Hilfsmenge $\Omega^{(2)}$ analysiert im Folgenden Frequenzen verschiedener Einträge von $\phi(\mathbf{x}_{\text{PE}})$. Bevor die Menge $\Omega^{(2)}$ an Frequenzbedingungen formal definiert wird, werden zunächst die Hilfsmatrizen \mathbf{P}^f und \mathbf{T}^f gegeben. Die Spalten der Matrix \mathbf{P}^f beinhalten dabei Kombinationen der Frequenzen zwischen den verschiedenen Elementen $\phi_{\bar{h}}(\mathbf{x}_{\text{PE}})$ und \mathbf{T}^f indiziert den jeweils zugehörigen Typ der trigonometrischen Funktion.

Definition 5.7 (Repräsentation von Frequenz tupeln durch \mathbf{P}^f und \mathbf{T}^f)

Seien

$$W_{\bar{h}} := \left\{ \underline{\omega}_{1,\bar{h}}, \dots, \underline{\omega}_{L_{\bar{h}},\bar{h}}, \bar{\omega}_{1,\bar{h}}, \dots, \bar{\omega}_{K_{\bar{h}},\bar{h}} \right\}, \quad (5.39)$$

$\bar{h} \in \{1, \dots, h\}$, Mengen, welche die Frequenzen nach (5.32) und (5.33) beinhalten. Dann lassen sich mit $\zeta_{\bar{h}}^{(i)} \in W_{\bar{h}}$, $\bar{h} \in \{1, \dots, h\}$,

$$N_{\text{P}} = \prod_{\bar{h}=1}^h (L_{\bar{h}} + K_{\bar{h}}) \quad (5.40)$$

unterschiedliche Tupel $\mathcal{P}_i := \left(\zeta_1^{(i)}, \dots, \zeta_h^{(i)} \right)$, $i = 1, \dots, N_{\text{P}}$, definieren. Weiterhin sei \mathbf{P}^f eine $h \times N_{\text{P}}$ -Matrix, deren N_{P} Spalten die Einträge der Tupel $\mathcal{P}_1, \dots, \mathcal{P}_{N_{\text{P}}}$ enthalten. Die Information über den jeweils zugehörigen Typ der trigonometrischen Funktion wird in der Matrix $\mathbf{T}^f \in \mathbb{R}^{h \times N_{\text{P}}}$ codiert. Hierbei wird das (\bar{h}, n_{P}) -te Element von \mathbf{T}^f zu null gesetzt ($T_{\bar{h}, n_{\text{P}}}^f = 0$), wenn das (\bar{h}, n_{P}) -te Element $P_{\bar{h}, n_{\text{P}}}^f$ von \mathbf{P}^f das Argument einer Sinusfunktion darstellt. Ist es hingegen das Argument einer Kosinusfunktion, so gilt $T_{\bar{h}, n_{\text{P}}}^f = 1$.

Gemäß Definition 5.7 lassen sich folglich o. B. d. A. die Matrizen

$$\mathbf{P}^f = \begin{bmatrix} \underline{\omega}_{1,1} & \cdots & \underline{\omega}_{1,1} & \cdots & \underline{\omega}_{1,1} & \cdots & \underline{\omega}_{1,1} & \cdots & \bar{\omega}_{K_1,1} & \cdots & \bar{\omega}_{K_1,1} \\ \underline{\omega}_{1,2} & \cdots & \underline{\omega}_{1,2} & \cdots & \bar{\omega}_{K_2,2} & \cdots & \bar{\omega}_{K_2,2} & \cdots & \bar{\omega}_{K_2,2} & \cdots & \bar{\omega}_{K_2,2} \\ \vdots & \ddots & \vdots \\ \underline{\omega}_{1,h} & \cdots & \bar{\omega}_{K_h,h} & \cdots & \underline{\omega}_{1,h} & \cdots & \bar{\omega}_{K_h,h} & \cdots & \underline{\omega}_{1,h} & \cdots & \bar{\omega}_{K_h,h} \end{bmatrix} \quad (5.41)$$

und

$$\mathbf{T}^f = \begin{bmatrix} 0 & \cdots & 0 & \cdots & 0 & \cdots & 0 & \cdots & 1 & \cdots & 1 \\ 0 & \cdots & 0 & \cdots & 1 & \cdots & 1 & \cdots & 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & \cdots & 0 & \cdots & 1 & \cdots & 0 & \cdots & 1 \end{bmatrix} \quad (5.42)$$

definieren. Damit lässt sich nun $\Omega^{(2)}$ formulieren.

Definition 5.8 ($\Omega^{(2)}$): **Bedingungen an die Frequenzen in P^f**

Mit P^f (5.41) und T^f (5.42) nach Definition 5.7 sei

$$\begin{aligned}
 \Omega^{(2)} := & \left\{ \omega : \bigvee_{n_p=1}^{N_p} \underbrace{\left(\bigwedge_{\bar{h}=1}^h P_{\bar{h},n_p}^f \neq 0 \right)}_I \right. \\
 & \wedge \underbrace{\bigwedge_{\bar{h}_1=1}^{h-1} \bigwedge_{\bar{h}_2 \in \bar{H}} \left(P_{\bar{h}_1,n_p}^f \neq P_{\bar{h}_2,n_p}^f \wedge P_{\bar{h}_1,n_p}^f \neq -P_{\bar{h}_2,n_p}^f \right)}_{II} \\
 & \wedge \underbrace{\bigwedge_{\substack{\bar{h}_{1/2}=1 \\ \bar{h}_1 \neq \bar{h}_2}}^h \bigwedge_{\bar{l} \in \bar{L}} \left(P_{\bar{h}_1,n_p}^f \neq \omega_{\bar{l},\bar{h}_2} \wedge P_{\bar{h}_1,n_p}^f \neq -\omega_{\bar{l},\bar{h}_2} \right)}_{III_a} \\
 & \left. \wedge \underbrace{\bigwedge_{\substack{\bar{h}_{1/2}=1 \\ \bar{h}_1 \neq \bar{h}_2}}^h \bigwedge_{\bar{k} \in \bar{K}} \left(P_{\bar{h}_1,n_p}^f \neq \bar{\omega}_{\bar{k},\bar{h}_2} \wedge P_{\bar{h}_1,n_p}^f \neq -\bar{\omega}_{\bar{k},\bar{h}_2} \right)}_{III_b} \right\} \quad (5.43)
 \end{aligned}$$

mit den Indexmengen

$$\begin{aligned}
 \bar{H} &:= \bar{H}(\bar{h}_1, n_p) := \left\{ \bar{h}_2 \in \{1, \dots, h\} : \bar{h}_2 > \bar{h}_1 \wedge T_{\bar{h}_1, n_p}^f = T_{\bar{h}_2, n_p}^f \right\}, \\
 \bar{L} &:= \bar{L}(\bar{h}_1, \bar{h}_2, n_p) := \left\{ \bar{l} \in \{1, \dots, L_{\bar{h}_2}\} : \omega_{\bar{l}, \bar{h}_2} \neq P_{\bar{h}_2, n_p}^f \wedge T_{\bar{h}_1, n_p}^f = 0 \right\}, \\
 \bar{K} &:= \bar{K}(\bar{h}_1, \bar{h}_2, n_p) := \left\{ \bar{k} \in \{1, \dots, K_{\bar{h}_2}\} : \bar{\omega}_{\bar{k}, \bar{h}_2} \neq P_{\bar{h}_2, n_p}^f \wedge T_{\bar{h}_1, n_p}^f = 1 \right\}.
 \end{aligned} \quad (5.44)$$

Somit muss nach Definition 5.8 mindestens eine der N_p Frequenzkombinationen, die durch die N_p Spalten von P^f (5.41) gegeben sind, die Bedingungen innerhalb der äußersten Disjunktion in (5.43) erfüllen. Die zu I zusammengefassten Bedingungen fordern, dass die durch die n_p -te Spalte von P^f beschriebenen Frequenzterme nicht verschwinden. Bedingung II garantiert, dass keine zwei trigonometrischen Funktionen desselben Typs in der betrachteten Spalte von P^f dieselbe Frequenz aufweisen. Bedingung III = III_a \wedge III_b gewährleistet schließlich, dass der Frequenzterm $P_{\bar{h}_1, n_p}^f$ der betrachteten Frequenzkombination, d. h. der n_p -ten Spalte von P^f , ungleich zu den Frequenztermen desselben trigonometrischen Typs in den anderen Elementen $\phi_{\bar{h}_2}(\mathbf{x}_{PE})$, $\bar{h}_2 \neq \bar{h}_1$, ist. Ungleichheit impliziert hier, genau wie in Definition 5.6, auch die Ungleichheit zu Frequenztermen mit gegensätzlichem Vorzeichen.

Analog zu Lemma 5.4 kann die Menge $\Omega^{(2)}$ (5.43) vereinfacht werden. Hierzu wird zunächst die Komplementärmenge $\bar{\Omega}^{(2)}$, die aus der Anwendung der De Morganschen Gesetze resultiert, mithilfe von Algorithmus 5.2 vereinfacht. Daraus ergibt sich¹⁰⁰

$$\bar{\Omega}^{(2)} = \left\{ \omega : \bigwedge_{n_p=1}^{N_p} \left(\bigvee_{\bar{h}=1}^h P_{\bar{h},n_p} = 0 \right. \right. \\ \bigvee \bigvee_{\bar{h}_1=1}^{h-1} \bigvee_{\bar{h}_2 \in \bar{h}} (P_{\bar{h}_1,n_p} = P_{\bar{h}_2,n_p} \vee P_{\bar{h}_1,n_p} = -P_{\bar{h}_2,n_p}) \\ \bigvee \bigvee_{\substack{\bar{h}_{1/2}=1 \\ \bar{h}_1=\bar{h}_2}}^h \bigvee_{\bar{l} \in \bar{L}} (P_{\bar{h}_1,n_p} = \omega_{\bar{l},\bar{h}_2} \vee P_{\bar{h}_1,n_p} = -\omega_{\bar{l},\bar{h}_2}) \\ \left. \left. \bigvee \bigvee_{\substack{\bar{h}_{1/2}=1 \\ \bar{h}_1=\bar{h}_2}}^h \bigvee_{\bar{k} \in \bar{K}} (P_{\bar{h}_1,n_p} = \bar{\omega}_{\bar{k},\bar{h}_2} \vee P_{\bar{h}_1,n_p} = -\bar{\omega}_{\bar{k},\bar{h}_2}) \right) \right\} \quad (5.45a)$$

$$= \left\{ \omega : \underbrace{\left(\bigvee_{n_p=1} (\cdot) = (\cdot) \right) \wedge \left(\bigvee_{n_p=2} (\cdot) = (\cdot) \right)}_{\text{IV}} \wedge \bigwedge_{n_p=3}^{N_p} \left(\bigvee_{n_p=3} (\cdot) = (\cdot) \right) \right\} \quad (5.45b)$$

$$= \left\{ \omega : \left(\bigvee_z C_z \omega = \mathbf{0} \right) \wedge \bigwedge_{n_p=3}^{N_p} \left(\bigvee_{n_p=3} (\cdot) = (\cdot) \right) \right\} \quad (5.45c)$$

$$= \left\{ \omega : \bigvee_{z=1}^{Z_2} C_z \omega = \mathbf{0} \right\}. \quad (5.45d)$$

Für den Spezialfall, dass im letzten Vereinfachungsschritt von $\bar{\Omega}^{(2)}$, d. h. für $k_{it} \geq N_p - 1$, $\omega_j = \omega_j, \forall j \in \{1, \dots, m\}$, folgt, ist die Aussage $\bar{\Omega}^{(2)}$ für beliebige $\omega \in \mathbb{R}^m$ stets wahr und es resultiert $\bar{\Omega}^{(2)} = \mathbb{R}^m$ und somit $\Omega^{(2)} = \emptyset$. Andernfalls ist $\bar{\Omega}^{(2)}$ durch die Disjunktion von $Z_2 \in \mathbb{N}_{\geq 1}$ Bedingungen der Form $C_z \omega = \mathbf{0}, z = 1, 2, \dots, Z_2$, gegeben (vgl. (5.45d)). Dies wird im nachfolgenden Lemma formal gezeigt.

¹⁰⁰ Für den Spezialfall $N_p = 1$ entfällt die Konjunktion mit den durch $n_p = k_{it} + 1$ indizierten Bedingungen in Schritt 1 von Algorithmus 5.2 und (5.45d) folgt aus (5.45a).

Algorithmus 5.2 Vereinfachung von $\bar{\Omega}^{(2)}$

1: **Initialisiere** die Iterationsvariable $k_{\text{it}} := 1$.

Schritt 1

2: Aufstellen homogener linearer Gleichungssysteme, indem die Konjunktion der mit $n_p = k_{\text{it}}$ und $n_p = k_{\text{it}} + 1$ indizierten Bedingungen unter Anwendung des Distributivgesetzes expandiert wird. Falls $k_{\text{it}} + 1 > N_p$ gilt, entfällt diese Konjunktion.

Schritt 2

3: Lösen der linearen Gleichungssysteme nach ω . Hierbei werden frei wählbare Variablen zu $\omega_j = \omega_j$ gesetzt und die Bedingungen durch $\bar{C}_z \omega = \mathbf{0}$ ($\bar{C}_z \in \mathbb{R}^{m \times m}$) mit dem Laufindex z ausgedrückt.

4: **if** $\omega_j = \omega_j$ für jedes $j \in \{1, \dots, m\}$ **then**

5: **if** $k_{\text{it}} \geq N_p - 1$ **then**

6: **return** $\bar{\Omega}^{(2)} = \mathbb{R}^m$

7: **else**

8: Setze $k_{\text{it}} := k_{\text{it}} + 2$ und gehe zu Schritt 1.

9: **end if**

10: **else**

11: Setze $k_{\text{it}} := k_{\text{it}} + 1$ und gehe zu Schritt 3.

12: **end if**

Schritt 3

13: Transformation der Koeffizientenmatrizen \bar{C}_z auf reduzierte Zeilenstufenform, Entfernen aller Nullzeilen. Resultierende Matrizen werden mit C_z ($C_z \in \mathbb{R}^{\bar{m} \times m}$, $\bar{m} < m$) bezeichnet.

Schritt 4

14: Entfernen mehrfach auftretender Bedingungen (identischer C_z). Alle bisher expandierten Konjunktionen werden anschließend zu $\bigvee_z C_z \omega = \mathbf{0}$ zusammengefasst.

Schritt 5

15: **while** $k_{\text{it}} \leq N_p - 1$ **do**

16: Stelle lineare homogene Gleichungssysteme durch die Expansion der Konjunktion zwischen dem Ergebnis aus Schritt 4 und den durch $n_p = k_{\text{it}} + 1$ indizierten Bedingungen auf. Wende die Schritte 2–4 auf diese Gleichungssysteme an.

17: **end while**

18: **return** $\bar{\Omega}^{(2)}$

Lemma 5.5 (Vereinfachung der Menge $\Omega^{(2)}$)

Falls Algorithmus 5.2 die Menge $\bar{\Omega}^{(2)} = \mathbb{R}^m$ ausgibt, so folgt $\Omega^{(2)} = \emptyset$. Andernfalls ist (5.43) identisch zu

$$\Omega^{(2)} = \left\{ \omega : \bigwedge_{z=1}^{Z_2} C_z \omega \neq \mathbf{0} \right\} \quad (5.46)$$

mit $Z_2 \in \mathbb{N}_{\geq 1}$.

Beweis:

Um zu beweisen, dass (5.46) dieselbe Menge wie (5.43) definiert, muss sowohl die Gültigkeit der Vereinfachungsschritte 3 und 4 nach Algorithmus 5.2 als auch die Korrektheit der Schlussfolgerungen im Fall $\omega_j = \omega_j, \forall j \in \{1, \dots, m\}$, in Schritt 2 betrachtet werden.

Zunächst wird Schritt 3 analysiert. Die Transformation der Koeffizientenmatrix eines homogenen linearen Gleichungssystems auf reduzierte Zeilenstufenform verändert dessen Lösungsmenge nicht. Des Weiteren entsprechen Nullzeilen in \bar{C}_z dem Fall $\omega_j = \omega_j$ und stellen somit Bedingungen dar, die durch $\Omega^{(2)}$ nicht erfüllt werden können. Folglich können Nullzeilen in \bar{C}_z entfernt werden. Schritt 4 ist gerechtfertigt, da mehrfach auftretende Bedingungen keine weiteren Informationen hinsichtlich der Definition der gesuchten Menge beinhalten und somit redundant sind.

Nun werde o. B. d. A. der Fall betrachtet, in dem sich in Schritt 2 für die durch IV gekennzeichneten Bedingungen in (5.45) $\omega_j = \omega_j$ für alle $j \in \{1, \dots, m\}$ ergibt. Dann ist Teil IV für jedes beliebige $\omega \in \mathbb{R}^m$ unabhängig von anderen Bedingungen der Form $\bar{C}_z \omega = \mathbf{0}$ in IV wahr. Somit können diese anderen Bedingungen vernachlässigt werden und Schritt 1 muss mit $n_p = 3$ und $n_p = 4$ ausgeführt werden. Da sich diese Überlegungen ebenfalls auf die Wiederholungen in Schritt 5 übertragen lassen und der Spezialfall, dass $\omega_j = \omega_j, \forall j \in \{1, \dots, m\}$, in der letzten Iteration $k_{it} \geq N_p - 1$ resultiert, explizit berücksichtigt wird, folgt die Gültigkeit von Lemma 5.5. \square

5.4.3 Frequenzbedingungen Ω

Die aus dem Schnitt der Hilfsmengen $\Omega^{(1)}$ und $\Omega^{(2)}$ resultierenden Frequenzbedingungen Ω lassen sich, wie im folgenden Lemma gezeigt, weiter vereinfachen.

Lemma 5.6 (Frequenzbedingungen Ω)

Die Schnittmenge zwischen $\Omega^{(1)}$ nach (5.38) und $\Omega^{(2)}$ nach (5.46) mit $Z_2 \in \mathbb{N}_{\geq 1}$ ergibt

$$\begin{aligned} \Omega &= \Omega^{(1)} \cap \Omega^{(2)} = \left\{ \omega : \bigwedge_{z=1}^{Z_1+Z_2} C_z \omega \neq \mathbf{0} \right\} \\ &= \left\{ \omega : \bigwedge_{z=1}^Z C_z \omega \neq \mathbf{0} \right\}, \quad Z \in \mathbb{N}_{\geq 1}, \end{aligned} \quad (5.47)$$

mit $Z \leq Z_1 + Z_2$. Die letzte Gleichung resultiert, indem jede Matrix C_{z_2} entfernt wird, für die bereits eine andere Matrix C_{z_1} , $z_1 \neq z_2$, in der Konjunktion (5.47) enthalten ist und außerdem $M \in \mathbb{R}^{m_1 \times m_2}$ mit $m_1 \leq m_2$ existiert, sodass $C_{z_1} = M C_{z_2}$ gilt.

Beweis:

Die Schnittmenge $\Omega = \Omega^{(1)} \cap \Omega^{(2)}$ ergibt $Z_1 + Z_2$ Bedingungen, die durch die Koeffizientenmatrizen C_z formuliert sind¹⁰¹. Um die Anzahl der Bedingungen zu reduzieren, kann jede Matrix C_{z_2} vernachlässigt werden, deren Bedingungen bereits durch die Matrix $C_{z_1} = MC_{z_2}$, $z_1 \neq z_2$, sichergestellt sind. Diese Redundanz der Matrix C_{z_2} gilt, da aus $C_{z_1}\omega \neq \mathbf{0} \Leftrightarrow MC_{z_2}\omega \neq \mathbf{0}$ direkt $C_{z_2}\omega \neq \mathbf{0}$ folgt. \square

Nach diesen Vorbereitungen kann schließlich die Hauptaussage dieses Kapitels formuliert werden, die belegt, dass die Menge Ω , welche Bedingungen an die Frequenzen ω in x_{PE} formuliert, eine Menge im Sinne von Proposition 5.1 darstellt.

Satz 5.2 (PE von $\sigma = \dot{\phi}(x_{PE}) \forall \omega \in \Omega$)

Seien $\phi(x)$ und x_{PE} entsprechend Annahme 5.3 und Annahme 5.4 gewählt und Ω wie in Lemma 5.6 definiert. Falls $\Omega \neq \emptyset$ gilt, so garantiert jeder beliebige Vektor $\omega \in \Omega$, dass das Signal $x_{PE}(t)$ SR bezüglich $\dot{\phi}(\cdot)$ ist und somit $\sigma = \dot{\phi}(x_{PE})$ PE ist.

Beweis:

Nach Lemma 5.2 lässt sich $\phi(x_{PE})$ durch (5.31) ausdrücken. Darauf basierend lassen sich die Hilfsmengen $\Omega^{(1)}$ nach (5.37) und $\Omega^{(2)}$ nach (5.43) definieren. Da Lemma 5.4, Lemma 5.5 und Lemma 5.6 lediglich Äquivalenztransformationen anwenden und die Mengen $\Omega^{(1)}$, $\Omega^{(2)}$ und Ω nicht verändern, muss im Folgenden nur gezeigt werden, dass die Bedingungen in (5.37) und (5.43) hinreichend sind, damit (5.35) gilt. Dann folgt nach Proposition 5.1, dass $x_{PE}(t)$ für jeden Vektor $\omega \in \Omega$ ein SR-Signal bezüglich $\dot{\phi}(\cdot)$ darstellt und somit $\sigma = \dot{\phi}(x_{PE})$ PE ist.

Die zeitliche Ableitung von (5.31) liefert

$$\sigma_{\bar{h}} = \sum_{l=1}^{L_{\bar{h}}} a_{l,\bar{h}} \omega_{l,\bar{h}} \cos(\omega_{l,\bar{h}} t) - \sum_{k=1}^{K_{\bar{h}}} c_{k,\bar{h}} \bar{\omega}_{k,\bar{h}} \sin(\bar{\omega}_{k,\bar{h}} t), \quad (5.48)$$

$\forall \bar{h} \in \{1, \dots, h\}$. Mit $\alpha \in \mathbb{R}^h$ folgt aus (5.48)

$$\begin{aligned} \alpha^\top \sigma &= \alpha_1 \left(\sum_{l=1}^{L_1} a_{l,1} \omega_{l,1} \cos(\omega_{l,1} t) - \sum_{k=1}^{K_1} c_{k,1} \bar{\omega}_{k,1} \sin(\bar{\omega}_{k,1} t) \right) \\ &+ \dots \\ &+ \alpha_h \left(\sum_{l=1}^{L_h} a_{l,h} \omega_{l,h} \cos(\omega_{l,h} t) - \sum_{k=1}^{K_h} c_{k,h} \bar{\omega}_{k,h} \sin(\bar{\omega}_{k,h} t) \right) \\ &= \sum_{\bar{h}=1}^h \left(\alpha_{\bar{h}} \left(\sum_{l=1}^{L_{\bar{h}}} a_{l,\bar{h}} \omega_{l,\bar{h}} \cos(\omega_{l,\bar{h}} t) - \sum_{k=1}^{K_{\bar{h}}} c_{k,\bar{h}} \bar{\omega}_{k,\bar{h}} \sin(\bar{\omega}_{k,\bar{h}} t) \right) \right). \end{aligned} \quad (5.49)$$

¹⁰¹ Hierbei werden die Zeilenvektoren c_z^\top in (5.38) ebenfalls als Koeffizientenmatrizen C_z aufgefasst.

O. B. d. A. werde nun für die beispielhaft betrachtete $(L_h + 1)$ -te Spalte

$$\mathcal{P}_{L_h+1} = \left(\zeta_1^{\langle L_h+1 \rangle}, \dots, \zeta_h^{\langle L_h+1 \rangle} \right) \quad (5.50)$$

von \mathbf{P}^f in (5.41) mit $\zeta_{\bar{h}}^{\langle L_h+1 \rangle} = \omega_{1,\bar{h}}, \forall \bar{h} \in \{1, \dots, h-1\}$, und $\zeta_h^{\langle L_h+1 \rangle} = \bar{\omega}_{1,h}$ angenommen, dass die Bedingungen I, II und III in (5.43) erfüllt sind. Ausgehend von (5.49) werden die trigonometrischen Funktionen, welche die entsprechenden Frequenzen dieser Spalte als Argumente beinhalten, separiert, und die verbleibenden trigonometrischen Ausdrücke innerhalb der Klammern $\alpha_{\bar{h}}(\cdot), \forall \bar{h} \in \{1, \dots, h\}$, zu $\varepsilon_{\bar{h}}$ zusammengefasst. Damit folgt

$$\begin{aligned} \alpha^\top \sigma &= \alpha_1 (a_{1,1} \omega_{1,1} \cos(\omega_{1,1} t) + \varepsilon_1) \\ &+ \dots \\ &+ \alpha_{h-1} (a_{1,h-1} \omega_{1,h-1} \cos(\omega_{1,h-1} t) + \varepsilon_{h-1}) + \alpha_h (-c_{1,h} \bar{\omega}_{1,h} \sin(\bar{\omega}_{1,h} t) + \varepsilon_h) \\ &= \sum_{\bar{h}=1}^{h-1} \alpha_{\bar{h}} a_{1,\bar{h}} \omega_{1,\bar{h}} \cos(\omega_{1,\bar{h}} t) - \alpha_h c_{1,h} \bar{\omega}_{1,h} \sin(\bar{\omega}_{1,h} t) + \sum_{\bar{h}=1}^h \alpha_{\bar{h}} \varepsilon_{\bar{h}}. \end{aligned} \quad (5.51)$$

Die Bedingungen in $\Omega^{(1)}$ nach (5.37) stellen sicher, dass die separierten trigonometrischen Funktionen in (5.51) innerhalb der Klammern $\alpha_{\bar{h}}(\cdot)$ nicht durch trigonometrische Ausdrücke derselben Art eliminiert werden. Dies folgt aus der geforderten Einzigartigkeit der Frequenzen in jedem Element $\phi_{\bar{h}}(\mathbf{x}_{PE})$ gemäß (5.31) (siehe Definition 5.6). Einzigartigkeit meint hierbei, dass Frequenzen, die Argumente gleichartiger trigonometrischer Funktionen sind, ungleich zueinander und zu ihrer negierten Frequenz sind. Darüber hinaus garantieren die Bedingungen in $\Omega^{(1)}$, dass in dem Element $\phi_{\bar{h}}(\mathbf{x}_{PE})$ $L_{\bar{h}}$ Sinus- und $K_{\bar{h}}$ Kosinusfunktionen existent sind und zwei Sinus- bzw. Kosinusfunktionen sich nicht auf eine einzelne Sinus- bzw. Kosinusfunktion reduzieren lassen. Hierdurch wird gewährleistet, dass die Kardinalität jeder Menge $W_{\bar{h}}, \bar{h} \in \{1, \dots, h\}$, (vgl. (5.39)) $L_{\bar{h}} + K_{\bar{h}}$ entspricht. Dies ist relevant, damit die Matrix \mathbf{P}^f der Frequenztuplel nach Definition 5.7 aufgestellt werden kann und Bedingungen an die Spalten von \mathbf{P}^f abgeleitet werden können.

Die durch III gekennzeichneten Bedingungen in $\Omega^{(2)}$ (5.43) gewährleisten, dass die trigonometrische Funktion, die innerhalb der Klammern $\alpha_{\bar{h}_1}(\cdot)$ in (5.51) separiert wurde, nicht durch trigonometrische Funktionen in $\varepsilon_{\bar{h}_2}$ innerhalb der Klammer $\alpha_{\bar{h}_2}(\cdot)$ ($\forall \bar{h}_1, \bar{h}_2 \in \{1, \dots, h\}, \bar{h}_1 \neq \bar{h}_2$) eliminiert wird. Auch dies gilt wegen der Einzigartigkeit der auftretenden Frequenzen.

Zusammenfassend können aufgrund der Bedingungen in $\Omega^{(1)}$ und der Bedingungen III in $\Omega^{(2)}$ die separierten trigonometrischen Ausdrücke in (5.51), unabhängig von den konkreten Werten von $\alpha_{\bar{h}}, a_{l,\bar{h}}$ und $c_{k,\bar{h}}$, nicht durch die Funktionen in $\varepsilon_{\bar{h}}, \forall \bar{h} \in \{1, \dots, h\}$, eliminiert werden. Daher genügt es, zu zeigen, dass die Bedingungen, die durch I und II in $\Omega^{(2)}$ (5.43) gekennzeichnet sind, unter der Annahme $\varepsilon_{\bar{h}} = 0, \forall \bar{h} \in \{1, \dots, h\}$, zu einem Signal $\alpha^\top \sigma$ nach (5.35) führen. Falls $\varepsilon_{\bar{h}} \neq 0$ gilt, können die in (5.51) separierten trigonometrischen Funktionen nicht ausgelöscht werden, es können höchstens zusätzliche Terme hinzukommen.

Im Folgenden gelte $\varepsilon_{\bar{h}} = 0, \forall \bar{h} \in \{1, \dots, h\}$. Aufgrund der Bedingungen I in $\Omega^{(2)}$ verschwinden die Frequenzen der trigonometrischen Ausdrücke, die in (5.51) separiert wurden, nicht und

es folgt $\underline{\omega}_{1,\bar{h}} \neq 0$, $a_{1,\bar{h}}\underline{\omega}_{1,\bar{h}} \neq 0$, $\forall \bar{h} \in \{1, \dots, h-1\}$, sowie $\bar{\omega}_{1,h} \neq 0$ und $c_{1,h}\bar{\omega}_{1,h} \neq 0$. Falls $\varepsilon_{\bar{h}} = 0 \forall \bar{h} \in \{1, \dots, h\}$ gilt, sind diese Bedingungen notwendig für (5.35). Dies lässt sich wie folgt einsehen. Wäre eine Frequenz null, so existierte ein $\alpha \neq \mathbf{0}$, sodass $\alpha^\top \sigma = 0$ folgt. Würde beispielsweise das \bar{h} -te Element von σ die verschwindende Frequenz beinhalten, so wäre ein solches α gegeben, indem alle Einträge von α außer dem \bar{h} -ten zu null gesetzt würden. Mit den durch Π gekennzeichneten Bedingungen in $\Omega^{(2)}$ bilden die separierten trigonometrischen Ausdrücke $a_{1,1}\underline{\omega}_{1,1} \cos(\underline{\omega}_{1,1}t), \dots, a_{1,h-1}\underline{\omega}_{1,h-1} \cos(\underline{\omega}_{1,h-1}t), c_{1,h}\bar{\omega}_{1,h} \sin(\bar{\omega}_{1,h}t)$ in (5.51) ein linear unabhängiges Funktionensystem. Somit erfüllt ihre Linearkombination (5.35), falls mindestens einer der Koeffizienten ungleich null ist, also $\alpha \neq \mathbf{0}$ gilt. Im Spezialfall $\varepsilon_{\bar{h}} = 0 \forall \bar{h} \in \{1, \dots, h\}$ sind die Bedingungen Π in $\Omega^{(2)}$ wiederum ebenfalls notwendig zur Erfüllung von (5.35). Um dies einzusehen, sei o. B. d. A. $\underline{\omega}_{1,1} = \pm \underline{\omega}_{1,h-1}$. Dann folgt $\alpha^\top \sigma = 0$ für $\alpha^\top = \left[\mp \frac{1}{a_{1,1}} \quad 0 \quad \dots \quad 0 \quad \frac{1}{a_{1,h-1}} \quad 0 \right]$.

Somit sind für den Fall $\varepsilon_{\bar{h}} = 0$ die Bedingungen I und II in $\Omega^{(2)}$ notwendig und hinreichend dafür, dass (5.35) für jedes beliebige $\alpha \neq \mathbf{0}$ gilt. Wie zuvor gezeigt, sind die zusätzlichen Bedingungen im Fall $\varepsilon_{\bar{h}} \neq 0$ hinreichend dafür, dass sich trigonometrische Funktionen nicht aufheben. Da die beispielhaft gewählte Frequenzkombination, d. h. die Spalte aus \mathbf{P}^f im hier erfolgten Beweis, austauschbar ist und mindestens eine solche Spalte existiert, führt jeder Frequenzvektor $\omega \in \Omega = \Omega^{(1)} \cap \Omega^{(2)}$ zu einem Signal $\alpha^\top \sigma$, für das (5.35) gilt. Schlussendlich folgt aus Proposition 5.1, dass x_{PE} SR bezüglich der Transformation $\dot{\phi}(\cdot)$ ist und somit σ die PE-Eigenschaft erfüllt. \square

Satz 5.2 liefert ein hinreichendes Kriterium für die Wahl der Frequenzen ω in x_{PE} , sodass das Signal x_{PE} SR bezüglich $\dot{\phi}(\cdot)$ ist. Das folgende Lemma adressiert die Annahme in Satz 5.2, dass Ω eine nichtleere Menge ist.

Lemma 5.7

Falls der sukzessive Vereinfachungsprozess von $\bar{\Omega}^{(2)}$ nach Algorithmus 5.2 in der letzten Iteration ($k_{\text{it}} \geq N_{\text{p}} - 1$) nicht zur allgemeinen Lösung $\omega_j = \omega_j, \forall j \in \{1, \dots, m\}$, führt, ergibt sich Ω wie in (5.47) und es gilt $\Omega \neq \emptyset$.

Beweis:

Tritt $\omega_j = \omega_j, \forall j \in \{1, \dots, m\}$, in der letzten Iteration $k_{\text{it}} \geq N_{\text{p}} - 1$ von $\bar{\Omega}^{(2)}$ in Algorithmus 5.2 nicht auf, so folgt nach Lemma 5.5 die nichtleere Menge $\Omega^{(2)}$ durch (5.46). Zudem ist nach Lemma 5.6 Ω durch (5.47) gegeben. Unter der Annahme, jede Zeile c_z^\top der Koeffizientenmatrizen C_z erfülle die Ungleichung $c_z^\top \omega \neq 0$, lässt sich eine Teilmenge Ω_{I} von Ω definieren, wobei die durch Ω_{I} beschriebenen Bedingungen einschränkender als die Bedingungen in Ω sind, d. h.

$$\Omega_{\text{I}} = \left\{ \omega : \bigwedge_z c_z^\top \omega \neq 0 \right\} \subset \Omega = \left\{ \omega : \bigwedge_{z=1}^Z C_z \omega \neq \mathbf{0} \right\}. \quad (5.52)$$

Da $c_z^\top \omega = 0$ eine $(m - 1)$ -dimensionale Hyperebene im \mathbb{R}^m beschreibt, ist deren Lebesgue-Maß null und die Hyperebene definiert eine Lebesguesche Nullmenge [Coh13, S. 146]. Zudem führt jede abzählbare Vereinigung Lebesguescher Nullmengen wieder auf eine Lebesguesche Nullmenge [Coh13, Proposition 1.2.4]. Daraus folgt, dass Ω_1 den gesamten Raum \mathbb{R}^m mit Ausnahme einer Lebesgueschen Nullmenge enthält und somit $\Omega_1 \neq \emptyset$ sowie $\Omega \neq \emptyset$ gilt. \square

Bemerkung 5.3

Sollte in der letzten Iteration von $\bar{\Omega}^{(2)}$ nach Algorithmus 5.2 $\omega_j = \omega_j, \forall j \in \{1, \dots, m\}$, resultieren, d. h. sollten für die Komplementärmenge $\bar{\Omega}^{(2)}$ keinerlei Einschränkungen an die Wahl der Frequenzen gelten, so ist es ratsam, zusätzliche Frequenzvariablen in den Elementen der Anregungstrajektorie x_{PE} einzuführen, bis eine nichtleere Menge Ω resultiert.

Für den Spezialfall $N_P = 1$, d. h., falls P^f nach Definition 5.7 nur eine Spalte besitzt, sind die Frequenzbedingungen aus Satz 5.2 notwendig und hinreichend. Zudem bleibt die PE-Eigenschaft von $\sigma = \dot{\phi}(x_{PE})$ bei einer Skalierung von x_{PE} erhalten. Diese Aussagen sind in Lemma 5.8 bzw. Lemma 5.9 formalisiert.

Lemma 5.8

Seien x_{PE} und $\phi(x)$ wie in Satz 5.2 gewählt und Ω wie in Lemma 5.6 gegeben. Falls $N_P = 1$ gilt, ist $\sigma = \dot{\phi}(x_{PE})$ genau dann PE, wenn $\omega \in \Omega$.

Beweis:

Für $N_P = 1$ ergibt sich $\Omega^{(1)} = \mathbb{R}^m$. Da hierbei zudem die durch III in $\Omega^{(2)}$ (5.43) gekennzeichneten Bedingungen entfallen, spielen lediglich die Teilbedingungen I und II von $\Omega^{(2)}$ bei der Definition von Ω eine Rolle. Außerdem impliziert $N_P = 1$ in (5.51) $\varepsilon_{\bar{h}} = 0, \forall \bar{h} \in \{1, \dots, h\}$. Für diesen Fall sind gemäß dem Beweis von Satz 5.2 die Bedingungen I und II in $\Omega^{(2)}$ notwendig und hinreichend dafür, dass (5.35), aber insbesondere auch $\alpha^\top \sigma \neq 0$, für jedes beliebige $\alpha \neq 0$ gilt. Mit Lemma 5.3 folgt schließlich die gesuchte Aussage. \square

Lemma 5.9 (Skalierung von x_{PE})

Seien x_{PE} und Ω entsprechend Satz 5.2 gegeben. Dann ist

$$\bar{x}_{PE} = \begin{bmatrix} \nu_1 & 0 & \dots & 0 \\ 0 & \nu_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \nu_n \end{bmatrix} x_{PE} \quad (5.53)$$

für jeden Vektor $\omega \in \Omega$ und $\nu_j \in \mathbb{R}_{\neq 0}, \forall j \in \{1, \dots, n\}$, ein SR-Signal bezüglich $\dot{\phi}(\cdot)$.

Beweis:Siehe Anhang C.2. □

Während zu Beginn von Abschnitt 5.4 der Spielerindex i zunächst vernachlässigt wurde, sollen nun wieder alle N Spieler berücksichtigt werden. Dabei können sich deren Basisfunktionsvektoren $\phi_i(\mathbf{x})$, $i \in \mathcal{N}$, im Allgemeinen unterscheiden. Für ein gegebenes \mathbf{x}_{PE} nach Annahme 5.4 wird für jeden der N Basisfunktionsvektoren $\phi_i(\mathbf{x})$ nach Annahme 5.3 eine Menge an Frequenzbedingungen nach Satz 5.2 berechnet. Daraus ergeben sich N Mengen Ω_i . Wird deren Schnittmenge analog zu Lemma 5.6 gebildet, folgt die Menge

$$\Omega = \bigcap_{i=1}^N \Omega_i. \quad (5.54)$$

Jeder Vektor $\boldsymbol{\omega} \in \Omega$ stellt dann sicher, dass $\mathbf{x}_{\text{PE}}(t)$ SR bezüglich $\dot{\phi}_i(\cdot)$, $\forall i \in \mathcal{N}$, ist.

Satz 5.2 stellt das Hauptergebnis dieses Kapitels dar und liefert eine neuartige und generell anwendbare Aussage bezüglich PE-Signalen, die aus nichtlinearen, polynomiellen Transformationen resultieren. Die durch die Menge Ω definierten Frequenzbedingungen können für ADP-Systeme angewandt werden, welche die Erfüllung der PE-Bedingung nach Definition 5.4 erfordern. Insbesondere liefert Satz 5.2 eine Lösung für Problem 5.1, da $\mathbf{x}_{\text{PE}}(t)$ mit beliebigem $\boldsymbol{\omega} \in \Omega$ ein $\mathbf{x}(t)$ im Sinne von Problem 5.1 darstellt. Da aus praktischer Sicht die Überprüfung, ob die PE-Eigenschaft bei vorliegenden Signalen tatsächlich erfüllt ist, bislang nur unzureichend geklärt ist, beschäftigt sich der nächste Abschnitt mit dieser Fragestellung.

5.5 Signal zur Überprüfung der Erfüllung der PE-Eigenschaft

Im Folgenden werden zwei Eigenwertsignale vorgestellt, mit deren Hilfe eine Antwort auf die Frage, ob ein vorliegendes Signal die PE-Eigenschaft erfüllt, gegeben werden kann.

Lemma 5.10 (Eigenwertsignal $\lambda_{\min,1}(t)$ zur Überprüfung der PE-Eigenschaft)

Das Signal $\boldsymbol{\sigma}(t) : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^h$ ist genau dann PE $\forall t \geq t_0$, wenn Konstanten $\alpha > 0$ und $T > 0$ existieren, sodass $\forall t \geq t_0$

$$\lambda_{\min,1}(t) := \lambda_{\min}(\Xi_1(t)) := \lambda_{\min} \left(\int_t^{t+T} \boldsymbol{\sigma}(\tau) \boldsymbol{\sigma}^T(\tau) d\tau \right) \geq \alpha > 0 \quad (5.55)$$

gilt. Hierbei bezeichnet $\lambda_{\min}(\cdot)$ den kleinsten Eigenwert einer Matrix¹⁰².

¹⁰² Da es sich bei $\Xi_1(t)$ um eine symmetrische Matrix handelt, sind deren Eigenwerte alle reell und der kleinste Eigenwert ist eindeutig definiert.

Beweis:

Nach Definition 5.4 ist $\sigma(t)$ genau dann PE, wenn

$$\Xi_1(t) := \int_t^{t+T} \sigma(\tau)\sigma^\top(\tau) d\tau \succeq \alpha I \quad (5.56)$$

gilt. Da $\Xi_1(t)$ nach (5.56) $\forall t$ symmetrisch ist, existiert eine reguläre Matrix $M(t)$, sodass

$$\Xi_1(t) = M(t)G(t)M^{-1}(t) \quad (5.57)$$

gilt, wobei $G(t)$ eine Diagonalmatrix mit den Eigenwerten von $\Xi_1(t)$ ist [BSMM13, S. 325]. Die PE-Eigenschaft von $\sigma(t)$ ist folglich äquivalent zu

$$\Xi_1(t) = M(t)G(t)M^{-1}(t) \succeq \alpha I \quad (5.58)$$

und aufgrund der Regularität von $M(t)$ zu

$$G(t) \succeq \alpha I. \quad (5.59)$$

Sei o. B. d. A. das (j, j) -te Element von $G(t)$ gerade $\lambda_{\min}(\Xi_1(t))$ und $e_j(t)$ der j -te Einheitsvektor. Dann folgt

$$\lambda_{\min,1}(t) = \lambda_{\min}(\Xi_1(t)) = e_j^\top(t)G(t)e_j(t) \geq \alpha. \quad (5.60)$$

□

Lemma 5.11 (Eigenwertsignal $\lambda_{\min,2}(t)$ zur Überprüfung der PE-Eigenschaft)

Wenn $\forall t \geq t_0$ das Signal $\sigma(t) : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^h$ PE ist, so steigt $\forall t \geq t_0$ das Eigenwertsignal

$$\lambda_{\min,2}(t) := \lambda_{\min}(\Xi_2(t)) := \lambda_{\min} \left(\int_{t_0}^t \sigma(\tau)\sigma^\top(\tau) d\tau \right) \quad (5.61)$$

monoton und es gilt $\lambda_{\min,2}(kT + t_0) \geq k\alpha$, $k \in \mathbb{N}_{\geq 0}$, $\alpha > 0$, $T > 0$.

Beweis:

Sei $\sigma(t)$ PE. Unter Verwendung von Definition 5.4 folgt daher

$$\Xi_2(kT + t_0) = \sum_{j=0}^{k-1} \underbrace{\int_{jT+t_0}^{(j+1)T+t_0} \sigma(\tau)\sigma^\top(\tau) d\tau}_{\succeq \alpha I} \succeq k\alpha I. \quad (5.62)$$

Analog zu (5.57)–(5.59) resultiert aus (5.62)

$$\lambda_{\min,2}(kT + t_0) = \lambda_{\min}(\Xi_2(kT + t_0)) \geq k\alpha. \quad (5.63)$$

Zudem steigt (5.61) monoton, da für $kT + t_0 < t < (k + 1)T + t_0$

$$\Xi_2(t) = \Xi_2(kT + t_0) + \underbrace{\int_{kT+t_0}^t \boldsymbol{\sigma}(\tau) \boldsymbol{\sigma}^\top(\tau) \, d\tau}_{\geq 0} \succeq \Xi_2(kT) \quad (5.64)$$

gilt. □

Lemma 5.10 liefert mit dem Eigenwertsignal $\lambda_{\min,1}(t)$ und (5.55) ein notwendiges und hinreichendes Kriterium für die Erfüllung der PE-Eigenschaft. Für dessen Berechnung müsste jedoch ein festes T im Voraus bestimmt werden, was in der Simulation nicht praktikabel ist. Daher wird zunächst $\lambda_{\min,2}(t)$ betrachtet, um einen möglichen Kandidaten für T zu ermitteln. Vorteilhaft ist, dass das Signal $\lambda_{\min,2}(t)$ direkt zur Laufzeit berechnet werden kann. Es stellt nach Lemma 5.11 jedoch lediglich ein notwendiges Kriterium zur Erfüllung der PE-Eigenschaft dar. Nach Wahl eines numerischen Schwellwerts α_t wird der Zeitpunkt T_t bestimmt, zu dem erstmalig $\lambda_{\min,2} \geq \alpha_t$ gilt. T_t wird dann als Kandidat zur Überprüfung der Ungleichung (5.55) mit $T = T_t$ und $\alpha = \alpha_t$ verwendet. Ist (5.55) erfüllt, so ist $\boldsymbol{\sigma}(t)$ nach Lemma 5.10 PE. Sollte hingegen $\lambda_{\min,1}(t)$ den Schwellwert α_t unterschreiten, so wird T_t iterativ erhöht. Bleibt auch hierbei $\lambda_{\min,1}(t)$ unterhalb von α_t , so wurde die Einhaltung der PE-Eigenschaft für $\boldsymbol{\sigma}(t)$ nicht nachgewiesen. Sollte das Eigenwertsignal $\lambda_{\min,2}(t)$ sättigendes Verhalten aufweisen, so beweist Lemma 5.11, dass $\boldsymbol{\sigma}(t)$ nicht PE ist.

5.6 Anregungssignale für ADP-basierte Differenzialspiele

Die bisher präsentierten Hauptbeiträge des vorliegenden Kapitels formulieren Bedingungen an den Systemzustand \boldsymbol{x} , um Konvergenz des Policy-Evaluation-Schrittes der Policy Iteration nach Algorithmus 5.1 zu gewährleisten. Während Satz 5.2 zeigt, dass hierbei eine genaue Analyse möglich ist, stellt die Frage nach geeigneten Stellgrößen \boldsymbol{u}_i , $i \in \mathcal{N}$, um diese Bedingungen an den Systemzustand \boldsymbol{x} zu garantieren, insbesondere bei unbekannter Systemdynamik ein grundlegendes Problem dar. Auch für die Identifikation einer a priori unbekannt Systemdynamik wäre wiederum eine geeignete Anregung erforderlich (vgl. Abschnitt 2.3), sodass ein komplett analytischer Anregungssignalentwurf der fehlenden Systemkenntnis grundsätzlich widerspricht. Um jedoch zu Analysezwecken die Umsetzung von Simulationsbeispielen in Abschnitt 5.7 zu ermöglichen, wird die folgende Annahme verwendet¹⁰³.

Annahme 5.5

Für den Anregungssignalentwurf seien $\boldsymbol{f}(\boldsymbol{x})$, $\boldsymbol{g}(\boldsymbol{x})$, $\boldsymbol{\phi}_i(\boldsymbol{x})$, \boldsymbol{R}_{ii} und $\hat{\boldsymbol{w}}_i^{[l]}$, $\forall i \in \mathcal{N}$, bekannt.

¹⁰³ Wenngleich die Spieler $i \in \mathcal{N}$ entsprechend Definition 5.1 weiterhin keine Kenntnis der Regelgesetze $\hat{\boldsymbol{\mu}}_j^{[l]}$, $j \in \mathcal{N}$, $j \neq i$, besitzen, wird allein für den Entwurf von $\boldsymbol{u}_{\text{ex}}$ angenommen, dass $\hat{\boldsymbol{w}}_i^{[l]}$, $\forall i \in \mathcal{N}$, bekannt sind.

Im Folgenden wird ein additives Anregungssignal \mathbf{u}_{ex} gesucht, sodass die Critic-Gewichtsfehler $\tilde{\mathbf{w}}_i^{[l+1]} = \mathbf{w}_i^{[l+1]} - \hat{\mathbf{w}}_i^{[l+1]}$, $\forall i \in \mathcal{N}$, im Policy-Evaluation-Schritt (vgl. Algorithmus 5.1) exponentiell konvergieren. Hierzu ist nach Lemma 5.1 die Erfüllung der PE-Eigenschaft von σ_i , $\forall i \in \mathcal{N}$, notwendig und hinreichend. Um eine Eingriffsmöglichkeit zur Systemanregung zu bieten, wird ein Anregungssignal \mathbf{u}_{ex} verwendet, das die im Folgenden formulierte Problemstellung lösen soll.

Problem 5.2

Sei Annahme 5.5 erfüllt. Gesucht ist ein Anregungssignal

$$\mathbf{u}_{\text{ex}} = [\mathbf{u}_{\text{ex},1}^\top \quad \mathbf{u}_{\text{ex},2}^\top \quad \dots \quad \mathbf{u}_{\text{ex},N}^\top]^\top, \quad (5.65)$$

sodass σ_i , $\forall i \in \mathcal{N}$, PE ist. Dabei sind $\mathbf{u}_{\text{ex},i}$ Anregungssignale, die nach Abbildung 5.1 auf die den N Spielern zur Verfügung stehenden Eingänge addiert werden, d. h. es gilt $\mathbf{u}_i = \hat{\boldsymbol{\mu}}_i^{[l]}(\mathbf{x}) + \mathbf{u}_{\text{ex},i}$.

Wie in Abbildung 5.1 ersichtlich ist, wirkt jeder Spieler $i \in \mathcal{N}$ durch das jeweilige Regelgesetz $\hat{\boldsymbol{\mu}}_i^{[l]}(\mathbf{x})$ in der l -ten Iteration der Policy Iteration nach Algorithmus 5.1 auf das System ein. Basierend auf dem durch (5.15) gegebenen Gradientenabstieg werden die Critic-Gewichte $\hat{\mathbf{w}}_i^{[l+1]}$, $\forall i \in \mathcal{N}$, bis zur Konvergenz des Policy-Evaluation-Schrittes adaptiert. Anschließend

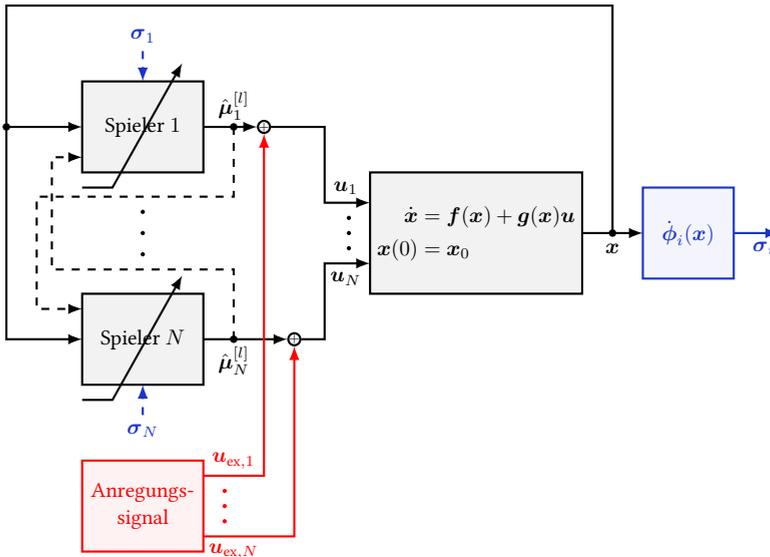


Abbildung 5.1: Regelkreisstruktur des eingangsaффinen Nicht-Nullsummen-Differenzialspiels mit Anregungssignal \mathbf{u}_{ex} (rot). Das Signal σ_i (blau) wird von jedem Spieler individuell berechnet und ist nicht als Systemgröße messbar, hier der Vollständigkeit halber jedoch dargestellt (Abbildung nach [Kar19]).

werden die Regelgesetze $\hat{\mu}_i^{[l+1]}(\mathbf{x})$, $\forall i \in \mathcal{N}$, gemäß (5.24) angepasst. Die Anregung des Systems erfolgt hierbei nicht über zusätzliche Anregungssignale der einzelnen Spieler, sondern soll durch das externe Anregungssignal \mathbf{u}_{ex} sichergestellt werden.

Die nachfolgende Bemerkung skizziert, weshalb häufig verwendetes weißes Rauschen als Anregungssignal \mathbf{u}_{ex} (vgl. Abschnitt 2.3) aus theoretischer Sicht zwar plausibel, jedoch praktisch ungeeignet erscheint.

Bemerkung 5.4 (Weißes Rauschen als Anregungssignal \mathbf{u}_{ex})

Nach Proposition 5.1 ist $\sigma = \dot{\phi}(\mathbf{x}_{\text{PE}})$ PE $\forall \omega \in \Omega$, wenn es gelingt, eine Menge Ω zu bestimmen, für die $\forall \omega \in \Omega$ und $\forall \alpha \neq \mathbf{0}$ folgt, dass $\alpha^\top \sigma_i \neq 0$ gilt. Wird das Anregungssignal \mathbf{u}_{ex} in (5.65) durch weißes Rauschen realisiert, so erscheint es für vollständig steuerbare Systeme und sinnvoll gestellte Probleme, d. h. unter der Annahme, dass eine ausreichende Anregung grundsätzlich möglich ist, unwahrscheinlich, dass sich alle Frequenzen durch die Transformation $\dot{\phi}_i(\mathbf{x}_{\text{PE}})$ gegenseitig aufheben und $\alpha^\top \sigma_i = 0$ resultiert.

Jedoch sei zu erwähnen, dass eine Anregung technischer Systeme mit weißem Rauschen häufig nicht erwünscht ist. So sind beispielsweise hohe mechanische Belastungen zu erwarten, zudem dämpfen Systeme mit Tiefpassverhalten einen wesentlichen Teil der Anregungsenergie. Die Verwendung weißen Rauschens als Anregungssignal bietet also keine Freiheitsgrade, zusätzliche problem- oder systemspezifische Anforderungen zu berücksichtigen.

Im Folgenden werden die theoretischen Ergebnisse aus Abschnitt 5.4 und insbesondere Satz 5.2 verwendet, um ein Anregungssignal \mathbf{u}_{ex} zu entwerfen und somit Simulationen zu ermöglichen. Der prinzipielle Ablauf zum Entwurf von \mathbf{u}_{ex} ist in Abbildung 5.2 gezeigt. In einem ersten Schritt wird hierbei $\mathbf{x}_{\text{PE}}(t)$ in Übereinstimmung mit Annahme 5.4 gewählt. Die konkreten Frequenzen ω sind zu diesem Zeitpunkt noch frei. Mithilfe von Satz 5.2 werden dann die Signale $\dot{\phi}_i(\mathbf{x}_{\text{PE}})$, $\forall i \in \mathcal{N}$, analysiert und die Mengen Ω_i berechnet, sodass für jedes $\omega \in \Omega_i$ folgt, dass \mathbf{x}_{PE} SR bezüglich $\dot{\phi}_i(\cdot)$ ist. Anschließend können verbliebene Freiheitsgrade gewählt werden. Einerseits beinhaltet dies die konkrete Wahl der Frequenzen in \mathbf{x}_{PE} zu einem beliebigen $\omega_c \in \Omega = \bigcap_i \Omega_i$, andererseits kann durch Wahl von $\nu \in \mathbb{R}^n$, $\nu_j \neq 0$, eine Skalierung nach Lemma 5.9 erfolgen. Hierdurch können problemspezifische Anforderungen, wie beispielsweise die Tiefpasscharakteristik eines Systems, berücksichtigt werden. Schließlich wird im letzten Schritt ein Anregungssignal \mathbf{u}_{ex} basierend auf

$$\bar{\mathbf{x}}_{\text{PE}} = \text{diag}(\nu) \mathbf{x}_{\text{PE}}|_{\omega=\omega_c} \quad (5.66)$$

berechnet. Dieser letzte Schritt wird im Folgenden für das Beispiel eines exakt zustandslinearisierbaren Systems anhand eines flachheitsbasierten Vorsteuerungsansatzes betrachtet.

In einem ersten Schritt soll untersucht werden, ob alle $p = \sum_{i=1}^N p_i$ Systemeingänge zur Anregung benötigt werden. Zu diesem Zweck wird $\mathbf{g}(\mathbf{x})$ in (5.1) analysiert, um Spalten von

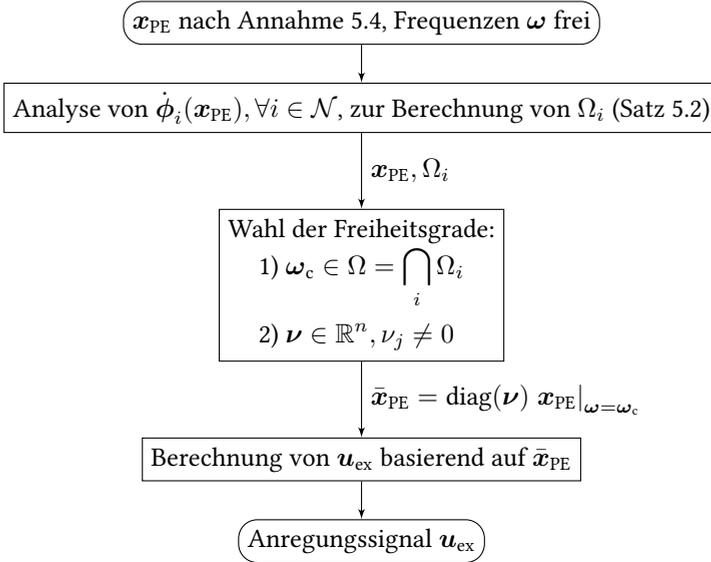


Abbildung 5.2: Ablaufdiagramm zum Design eines Anregungssignals u_{ex} .

$g(x)$ zu identifizieren, die durch einen funktionalen Zusammenhang mittels anderer Spalten ausgedrückt werden können. Eine Umsortierung der Spalten von $g(x)$ führt dann zu

$$\bar{g}(x) = [g_{\text{I}}(x) \quad g_{\text{II}}(x)], \quad (5.67)$$

$g_{\text{I}}: \mathbb{R}^n \rightarrow \mathbb{R}^{n \times p_{\text{I}}}$, $g_{\text{II}}: \mathbb{R}^n \rightarrow \mathbb{R}^{n \times p_{\text{II}}}$ und der funktionalen Beziehung $\gamma(x)$, sodass $\forall x \in \mathcal{X}$ der Zusammenhang

$$g_{\text{II}}(x) = g_{\text{I}}(x)\gamma(x) \quad (5.68)$$

besteht und $g_{\text{I}}(x)$ maximalen Spaltenrang aufweist, d. h. $\text{Rang}(g_{\text{I}}(x)) = p_{\text{I}}$ gilt. Werden die Systemeingänge in analoger Weise wie die Spalten von $g(x)$ umsortiert, ergibt sich

$$\bar{u} = \begin{bmatrix} u_{\text{I}} \\ u_{\text{II}} \end{bmatrix} \quad (5.69)$$

mit $\bar{u} \in \mathbb{R}^p$, $u_{\text{I}} \in \mathbb{R}^{p_{\text{I}}}$ und $u_{\text{II}} \in \mathbb{R}^{p_{\text{II}}}$. Da aufgrund von (5.68) eine Anregung des Systems über die zu u_{II} zusammengefassten Eingänge keinen Mehrwert liefert, wird das Anregungssignal $u_{\text{ex,II}}$ der entsprechenden Eingänge zu null gesetzt. Für die Systemanregung werden lediglich die Eingänge u_{I} verwendet, sodass $u_{\text{ex,I}}$ entworfen werden muss und sich das Anregungssignal u_{ex} schließlich aus der Rücksortierung von $\bar{u}_{\text{ex}} = [u_{\text{ex,I}}^{\text{T}} \quad u_{\text{ex,II}}^{\text{T}}]^{\text{T}} = [u_{\text{ex,I}}^{\text{T}} \quad \mathbf{0}_{1 \times p_{\text{II}}}]^{\text{T}}$ ergibt¹⁰⁴.

¹⁰⁴ Dieses Vorgehen, um die Anzahl der Eingänge zur Systemanregung zu reduzieren, ist generell anwendbar und lässt sich somit beispielsweise auch bei klassischer Anregung durch weißes Rauschen nutzen.

Das globale Anregungssignal \mathbf{u}_{ex} (bzw. $\mathbf{u}_{\text{ex},1}$) kann zu Simulationszwecken mithilfe eines inversionsbasierten Vorsteuerungsansatzes (siehe beispielsweise [DCP96], [Hir79], [Dev02]) entworfen werden. Im Folgenden wird für das Beispiel eines (global) exakt zustandslinearisierbaren Systems (vgl. [Ada18, Definition 29]), das folglich auch differenziell flach¹⁰⁵ ist [Ada18, Satz 67], eine flachheitsbasierte Vorsteuerung vorgestellt und analysiert. Dazu werde zunächst allgemein

$$\dot{\mathbf{x}} = \mathbf{f}_g(\mathbf{x}) + \mathbf{g}_I(\mathbf{x})\mathbf{u}_I \quad (5.70)$$

mit

$$\mathbf{f}_g(\mathbf{x}) = \mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})\hat{\boldsymbol{\mu}}^{[l]}(\mathbf{x}) \quad (5.71)$$

betrachtet. Somit entspricht \mathbf{f}_g der Dynamik des durch $\hat{\boldsymbol{\mu}}^{[l]}(\mathbf{x})$ geschlossenen Regelkreises des Systems (5.1). Dieses System lässt sich auf Byrnes-Isidori-Normalform (BINF) [Ada18, S. 354] ohne interne Dynamik transformieren, falls eine (ggf. fiktive) Ausgangsfunktion $\mathbf{h}(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^{p_I}$ mit vollem totalen relativen Grad¹⁰⁶ [Ada18, Definition 28] ($\delta = \sum_j \delta_j = n$) existiert [Zim84]. Nach [Isi89, S. 245] ist das System in BINF dann durch p_I Subsysteme der Form

$$\begin{aligned} \dot{z}_{\text{B},1,j} &= z_{\text{B},2,j}, \\ &\vdots \\ \dot{z}_{\text{B},\delta_j-1,j} &= z_{\text{B},\delta_j,j}, \\ \dot{z}_{\text{B},\delta_j,j} &= \underbrace{L_{\mathbf{f}_g}^{\delta_j} h_j(\mathbf{x})}_{\beta_j(\mathbf{x})} + \underbrace{\left[L_{\mathbf{g}_{1,1}} L_{\mathbf{f}_g}^{\delta_j-1} h_j(\mathbf{x}) \quad \dots \quad L_{\mathbf{g}_{1,p_I}} L_{\mathbf{f}_g}^{\delta_j-1} h_j(\mathbf{x}) \right]}_{\boldsymbol{\alpha}_j^T(\mathbf{x})} \mathbf{u}_I, \end{aligned} \quad (5.72)$$

$j \in \{1, \dots, p_I\}$, gegeben. Dabei bezeichnet $h_j(\mathbf{x})$ das j -te Element von $\mathbf{h}(\mathbf{x})$, L die Lie-Ableitung¹⁰⁷, δ_j den relativen Grad von (5.70) bezüglich $h_j(\mathbf{x})$ und $\mathbf{g}_{I,k}(\mathbf{x})$, $k \in \{1, \dots, p_I\}$, die k -te Spalte von $\mathbf{g}_I(\mathbf{x})$. Der Diffeomorphismus

$$\begin{aligned} \mathbf{z}_{\text{B}} &= \mathbf{t}(\mathbf{x}) \\ &= [t_{1,1}(\mathbf{x}) \quad \dots \quad t_{\delta_1,1}(\mathbf{x}) \quad t_{1,2}(\mathbf{x}) \quad \dots \quad t_{\delta_2,2}(\mathbf{x}) \quad t_{1,p_I}(\mathbf{x}) \quad \dots \quad t_{\delta_{p_I},p_I}(\mathbf{x})]^T, \end{aligned} \quad (5.73)$$

¹⁰⁵ Ein System $\dot{\mathbf{x}} = \mathbf{f}_{\text{dyn}}(\mathbf{x}, \mathbf{u})$ mit $\text{Rang} \left(\frac{\partial \mathbf{f}_{\text{dyn}}(\mathbf{x}, \mathbf{u})}{\partial \mathbf{u}} \right) = \dim(\mathbf{u})$ heißt flach (vgl. [FLMR95], [Ada18, Kapitel 3.2]), wenn ein (realer oder fiktiver) flacher Ausgang $\mathbf{y} = \mathbf{h}_{\Gamma}(\mathbf{x}, \mathbf{u}, \dot{\mathbf{u}}, \dots, \mathbf{u}^{(\chi_1)})$, $\chi_1 \in \mathbb{N}$, χ_1 endlich, $\dim(\mathbf{y}) = \dim(\mathbf{u})$, existiert, sodass $\exists \chi_2 \in \mathbb{N}$, χ_2 endlich, mit $\mathbf{x} = \Gamma_1(\mathbf{y}, \dot{\mathbf{y}}, \dots, \mathbf{y}^{(\chi_2)})$ und $\mathbf{u} = \Gamma_2(\mathbf{y}, \dot{\mathbf{y}}, \dots, \mathbf{y}^{(\chi_2+1)})$ [Ada18, Definition 27], [Lév09, Definition 6.1].

¹⁰⁶ Der totale relative Grad δ eines Systems ist auch als Differenzordnung eines Systems bekannt. Eine solche Ausgangsfunktion lässt sich bei einfachen Systemen häufig durch Betrachtung der Systemmatrizen ermitteln, ansonsten kann im Fall der betrachteten eingangsaффinen Systeme die in [Zim84] beschriebene Methode verwendet werden.

¹⁰⁷ $L_{\mathbf{f}}^j h(\mathbf{x}) := L_{\mathbf{f}} L_{\mathbf{f}}^{j-1} h(\mathbf{x})$ bezeichnet die j -te Lie-Ableitung von $h(\mathbf{x})$ bezüglich $\mathbf{f}(\mathbf{x})$. Dabei gilt per Definition $L_{\mathbf{f}} h(\mathbf{x}) := \frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} \mathbf{f}(\mathbf{x})$.

der den Systemzustand \mathbf{x} auf BINF transformiert, ergibt sich zu

$$\begin{aligned} z_{B,1,j} &= t_{1,j}(\mathbf{x}) = h_j(\mathbf{x}), \\ z_{B,2,j} &= t_{2,j}(\mathbf{x}) = L_{\mathbf{f}_g} h_j(\mathbf{x}), \\ &\vdots \\ z_{B,\delta_j,j} &= t_{\delta_j,j}(\mathbf{x}) = L_{\mathbf{f}_g}^{\delta_j-1} h_j(\mathbf{x}). \end{aligned} \quad (5.74)$$

Werden nun die Elemente $z_{B,1,j}, \forall j \in \{1, \dots, p_l\}$, verwendet, um

$$\mathbf{y}_f = \mathbf{h}_f(\mathbf{z}_B) = [z_{B,1,1} \quad z_{B,1,2} \quad \dots \quad z_{B,1,p_l}]^T \quad (5.75)$$

als Kandidat für einen flachen Ausgang zu wählen, so folgt

$$\begin{aligned} \mathbf{z}_B &= [z_{B,1,1} \quad \dot{z}_{B,1,1} \quad \dots \quad z_{B,1,1}^{(\delta_1-1)} \quad \dots \quad z_{B,1,p_l} \quad \dot{z}_{B,1,p_l} \quad \dots \quad z_{B,1,p_l}^{(\delta_{p_l}-1)}]^T \\ &=: \Psi_1(\mathbf{y}_f, \dot{\mathbf{y}}_f, \dots, \mathbf{y}_f^{(\delta_{\max}-1)}) \end{aligned} \quad (5.76)$$

mit $\delta_{\max} = \max\{\delta_1, \dots, \delta_{p_l}\}$. Aus (5.72) folgt

$$\begin{aligned} \mathbf{u}_I &= D^{-1}(\mathbf{x}) \left(-\boldsymbol{\beta}(\mathbf{x}) + \begin{bmatrix} \dot{z}_{B,\delta_1,1} \\ \vdots \\ \dot{z}_{B,\delta_{p_l},p_l} \end{bmatrix} \right) = D^{-1}(\mathbf{x}) \left(-\boldsymbol{\beta}(\mathbf{x}) + \begin{bmatrix} \mathbf{y}_{f,1}^{(\delta_1)} \\ \vdots \\ \mathbf{y}_{f,p_l}^{(\delta_{p_l})} \end{bmatrix} \right) \\ &= \Psi_2(\mathbf{y}_f, \dot{\mathbf{y}}_f, \dots, \mathbf{y}_f^{(\delta_{\max})}), \end{aligned} \quad (5.77)$$

wobei die Vektoren $\mathbf{d}_j^T(\mathbf{x})$ zur sogenannten Entkopplungsmatrix $\mathbf{D}(\mathbf{x})$ [Ada18, S. 373] zusammengefasst werden, sich der Vektor $\boldsymbol{\beta}(\mathbf{x})$ aus $\beta_j(\mathbf{x})$ zusammensetzt und im letzten Schritt $\mathbf{x} = \mathbf{t}^{-1}(\Psi_1(\cdot))$ genutzt wird. Falls nun $\det(\mathbf{D}(\mathbf{x})) \neq 0$ gilt, so existiert Ψ_2 in (5.77) und der gewählte Kandidat $\mathbf{h}_f(\mathbf{z}_B)$ stellt tatsächlich einen flachen Ausgang dar, da Ψ_1 und Ψ_2 existieren. Da Ψ_2 nach [FLMR95] die Systemdynamik invertiert, kann direkt eine Vorsteuerung entworfen werden. Hierzu wird für die hier betrachtete Problemstellung die gewählte Anregungstrajektorie $\bar{\mathbf{x}}_{PE}$ mit dem Diffeomorphismus $\mathbf{t}(\cdot)$ in BINF überführt und mithilfe von $\mathbf{y}_f|_{\bar{\mathbf{x}}_{PE}} = \mathbf{h}_f(\mathbf{z}_B|_{\bar{\mathbf{x}}_{PE}}) = \mathbf{h}_f(\mathbf{t}(\bar{\mathbf{x}}_{PE}))$ das Anregungssignal $\mathbf{u}_{ex,I}$ aus (5.77) berechnet. Schließlich resultiert das Anregungssignal \mathbf{u}_{ex} aus der Rücksortierung von $\bar{\mathbf{u}}_{ex} = [\mathbf{u}_{ex,I}^T \quad \mathbf{u}_{ex,II}^T]^T$, wobei $\mathbf{u}_{ex,II} = \mathbf{0}$ gilt.

Für die nachfolgenden Analysen bezeichne

$$\tilde{\mathbf{x}} = \mathbf{t}^{-1} \left(\Psi_1 \left(\mathbf{y}_f|_{\bar{\mathbf{x}}_{PE}}, \dot{\mathbf{y}}_f|_{\bar{\mathbf{x}}_{PE}}, \dots, \mathbf{y}_f^{(\delta_{\max}-1)}|_{\bar{\mathbf{x}}_{PE}} \right) \right) \quad (5.78)$$

den idealen Systemzustand, der mit exakter Systeminversion und ohne Anfangsfehler des Zustands, d. h. ohne Einschwingvorgänge, resultieren würde. Die Definition der Flachheit [Ada18, Definition 27] stellt insbesondere durch $\dim(\mathbf{y}_f) = \dim(\mathbf{u}_I)$ sicher, dass die Ausgänge

$y_{f,i}, i \in \{1, \dots, p_f\}$, differenziell unabhängig sind. Dies ermöglicht, dass die Wahl der (fiktiven) Ausgangsgröße \mathbf{y}_f beliebig ist und durch die Vorsteuerung auch erzielt werden kann (vgl. [Ada18, S. 210 f.]). Falls jedoch nicht der gesamte Zustandsvektor \mathbf{x} den flachen Ausgang \mathbf{y}_f darstellt, so können im Allgemeinen die Elemente des Zustandsvektors nicht unabhängig voneinander vorgegeben werden, sondern können differenzielle Abhängigkeiten aufweisen. In diesem Fall kann $\tilde{\mathbf{x}} = \mathbf{x}_{\text{PE}}$ im Allgemeinen nicht erreicht werden.

Im Folgenden wird daher formal untersucht, unter welchen Annahmen trotz einer Abweichung des idealen Zustands $\tilde{\mathbf{x}}$ von $\bar{\mathbf{x}}_{\text{PE}}$, beispielsweise durch differenzielle Abhängigkeiten, sowie möglichen Abweichungen, die der tatsächliche Systemzustand \mathbf{x} von $\tilde{\mathbf{x}}$ haben kann¹⁰⁸, die Erfüllung der PE-Eigenschaft von $\sigma_i(t) \forall i \in \mathcal{N}$ gewährleistet werden kann. Hierzu gelte die folgende Annahme.

Annahme 5.6

Für $\tilde{\mathbf{x}}$ sei

$$\dot{\sigma}_i := \sigma_i(\tilde{\mathbf{x}}(t)) = \dot{\phi}_i(\tilde{\mathbf{x}}(t)) = \mathbf{M}_i \frac{d\mathbf{v}_{\text{freq}}(t)}{dt} \quad (5.79)$$

mit $\mathbf{M}_i \in \mathbb{R}^{h_i \times h_{\text{freq}}}$, $h_i \leq h_{\text{freq}}$, und $\mathbf{v}_{\text{freq}}(t) \in \mathbb{R}^{h_{\text{freq}}}$, wobei $\mathbf{v}_{\text{freq}}(t)$ ein Vektor ist, dessen Elemente Sinus- und Kosinusfunktionen unterschiedlicher Frequenzen (d. h. unterschiedlich für identische trigonometrische Funktionen) darstellen.

Der tatsächliche Systemzustand \mathbf{x} , der unter Verwendung der Regelgesetze $\hat{\boldsymbol{\mu}}_1^{[l]}(\mathbf{x}), \dots, \hat{\boldsymbol{\mu}}_N^{[l]}(\mathbf{x})$ und des Anregungssignals \mathbf{u}_{ex} resultiert, führe zu

$$\sigma_i(t) = \check{\sigma}_i(t) + \varepsilon_{1,i}(t) + \varepsilon_{2,i}(t), \quad \forall i \in \mathcal{N}, \quad (5.80)$$

mit stetigen, beschränkten Signalen $\varepsilon_{1,i}(t), \varepsilon_{2,i}(t) : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{h_i}$. Zudem gelte $\|\varepsilon_{1,i}(t)\|_2 \leq \bar{\varepsilon}_{\sigma,i}, \forall t$, und $\varepsilon_{2,i}(t) \rightarrow \mathbf{0}$ für $t \rightarrow \infty$.

Bemerkung 5.5

Aufgrund von (5.79) muss für die nachfolgenden Aussagen somit nicht $\tilde{\mathbf{x}} = \bar{\mathbf{x}}_{\text{PE}}$ gelten, sondern differenzielle Abhängigkeiten können berücksichtigt werden. Ob (5.79) gilt, lässt sich anhand von Ψ_1 , dem Diffeomorphismus $\mathbf{t}(\cdot)$ zur Transformation auf BINF und dem flachen Ausgang \mathbf{y}_f a priori überprüfen. Zudem entsprechen die in \mathbf{v}_{freq} auftretenden Frequenzen Linearkombinationen der Einträge von ω_c . Weiterhin berücksichtigt (5.80) durch $\varepsilon_{1,i}(t)$ Modellungenauigkeiten, beispielsweise der Regelgesetze oder der Systemdynamik, und durch $\varepsilon_{2,i}(t)$ transiente Vorgänge. Beides kann zu Abweichungen zwischen \mathbf{x} und $\tilde{\mathbf{x}}$ führen.

Das nachfolgende Lemma liefert eine theoretische Aussage hinsichtlich der Erfüllung der PE-Bedingung für das aus dem tatsächlichen Zustand $\mathbf{x}(t)$ resultierende Signal $\sigma_i(\mathbf{x}(t))$.

¹⁰⁸ Insbesondere transiente Vorgänge während des Einschwingvorgangs können zu $\mathbf{x} \neq \tilde{\mathbf{x}}$ führen.

Lemma 5.12

Sei $\bar{x}_{\text{PE}}(t) = \text{diag}(\boldsymbol{\nu}) \mathbf{x}_{\text{PE}}(t)|_{\omega=\omega_c}$ wie in Lemma 5.9 mit entsprechender Wahl von $\omega_c \in \Omega$ und $\boldsymbol{\nu}$ gegeben. Unter Annahme 5.6 ist $\boldsymbol{\sigma}_i(t) = \boldsymbol{\sigma}_i(\mathbf{x}(t)) \forall i \in \mathcal{N}$ PE, falls $\bar{\varepsilon}_{\sigma,i}$ hinreichend klein ist und $\text{Rang}(\mathbf{M}_i) = h_i$ gilt.

Beweis:

Der Beweis erfolgt in zwei Schritten. Im ersten Schritt wird gezeigt, dass $\boldsymbol{\sigma}_i(t)$ PE ist, falls $\check{\boldsymbol{\sigma}}_i(t)$ die PE-Eigenschaft erfüllt. Im zweiten Schritt wird bewiesen, dass $\check{\boldsymbol{\sigma}}_i(t)$ PE ist. Mit (5.80) folgt aus (5.19)

$$\begin{aligned} & \frac{1}{T} \int_t^{t+T} |(\check{\boldsymbol{\sigma}}_i(\tau) + \boldsymbol{\varepsilon}_{1,i}(\tau) + \boldsymbol{\varepsilon}_{2,i}(\tau))^\top \mathbf{e}| \, d\tau \\ & \geq \frac{1}{T} \int_t^{t+T} |(\check{\boldsymbol{\sigma}}_i(\tau))^\top \mathbf{e}| \, d\tau - \frac{1}{T} \int_t^{t+T} |(\boldsymbol{\varepsilon}_{1,i}(\tau))^\top \mathbf{e}| \, d\tau - \frac{1}{T} \int_t^{t+T} |(\boldsymbol{\varepsilon}_{2,i}(\tau))^\top \mathbf{e}| \, d\tau \\ & \geq \alpha_{1,i} - \bar{\varepsilon}_{\sigma,i} - \frac{1}{T} M_\varepsilon(T). \end{aligned} \quad (5.81)$$

Die letzte Ungleichheit in (5.81) ergibt sich, wenn $\check{\boldsymbol{\sigma}}_i$ PE ist und aufgrund der Beschränktheit von $\boldsymbol{\varepsilon}_{1,i}$ sowie dem Konvergenzverhalten von $\boldsymbol{\varepsilon}_{2,i}$. Letzteres stellt

$$\int_t^{t+T} |(\boldsymbol{\varepsilon}_{2,i}(\tau))^\top \mathbf{e}| \, d\tau \leq M_\varepsilon(T) \quad (5.82)$$

sicher, wobei $M_\varepsilon(T)$ eine von T abhängige obere Schranke darstellt. Da $M_\varepsilon(T)$ aufgrund von $\boldsymbol{\varepsilon}_{2,i}(t) \rightarrow \mathbf{0}$ ($t \rightarrow \infty$) mit steigendem T in die Sättigung übergeht, existiert \bar{T} , sodass $\alpha_{1,i} > \bar{\varepsilon}_{\sigma,i} + \frac{1}{T} M_\varepsilon(T)$, $\forall T > \bar{T}$, gilt, wenn $\bar{\varepsilon}_{\sigma,i}$ hinreichend klein ist. Nach (5.19) folgt, dass $\boldsymbol{\sigma}_i$ ein PE-Signal ist, wenn $\check{\boldsymbol{\sigma}}_i$ PE ist. Dass $\check{\boldsymbol{\sigma}}_i$ ein PE-Signal ist, folgt schlussendlich aus (5.79) unter Anwendung von [NA05, Lemma 6.1], da $\mathbf{v}_{\text{freq}}(t)$ PE ist und \mathbf{M}_i maximalen Zeilenrang besitzt. Die PE-Eigenschaft von $\mathbf{v}_{\text{freq}}(t)$ ist dabei ein direktes Resultat aus Lemma 5.3. \square

In Abschnitt 5.4 wurden zentrale theoretische Zusammenhänge zwischen dem Systemzustand \mathbf{x} und der Erfüllung der PE-Eigenschaft von $\boldsymbol{\sigma}_i$ analysiert und somit Problem 5.1 gelöst. Darauf aufbauend wurde in Abschnitt 5.6, ausgehend von Problem 5.2, ein mögliches Entwurfsverfahren für Anregungssignale ADP-basierter Differenzialspiele vorgestellt, welches für die Umsetzung der nachfolgenden Simulationsbeispiele genutzt wird.

5.7 Simulationsbeispiel zur Anregung von ADP-basierten Differenzialspielen

In diesem Abschnitt wird zunächst ein Beispiel eines eingangsaffinen Nicht-Nullsummen-Differenzialspiels vorgestellt. Anschließend werden gemäß dem in Abbildung 5.2 gezeigten

Ablauf Anregungssignale konstruiert. Diese Anregungssignale sowie ein Vergleichssignal, das weißes Rauschen approximiert, werden anschließend zur Anregung eines Beispielsystems verwendet. Die Adaption der Critic-Gewichte wird hierfür in der Simulation betrachtet.

5.7.1 Beispielproblem

Um die Konvergenz der Critic-Gewichte anhand von Simulationen zu untersuchen, wird zunächst ein beispielhaftes Differenzialspiel definiert. Damit die durch \mathbf{w}_i^* , $\forall i \in \mathcal{N}$, beschriebene Lösung zu Vergleichszwecken bekannt ist, wird ein Beispielsystem mit $N = 2$ Spielern, das mithilfe des Converse-HJB-Ansatzes [NP96] konstruiert wird, verwendet¹⁰⁹. Nach Wahl der optimalen Value Functions zu

$$V_1^*(\mathbf{x}) := \frac{1}{2}x_1^2 + x_2^2, \quad V_2^*(\mathbf{x}) := \frac{1}{4}x_1^2 + \frac{1}{2}x_2^2 \quad (5.83)$$

und der Gütefunktionale der beiden Spieler zu

$$\begin{aligned} J_1(\mathbf{x}_0, \mu_1, \mu_2) &:= \int_0^\infty \underbrace{2(x_1^2 + x_2^2)}_{=:q_1(\mathbf{x})} + \underbrace{2(\mu_1(\mathbf{x}))^2}_{=: \mu_1^\top \mathbf{R}_{11} \mu_1} + \underbrace{2(\mu_2(\mathbf{x}))^2}_{=: \mu_2^\top \mathbf{R}_{12} \mu_2} d\tau, \\ J_2(\mathbf{x}_0, \mu_1, \mu_2) &:= \int_0^\infty \underbrace{x_1^2 + x_2^2}_{=:q_2(\mathbf{x})} + \underbrace{(\mu_1(\mathbf{x}))^2}_{=: \mu_1^\top \mathbf{R}_{21} \mu_1} + \underbrace{(\mu_2(\mathbf{x}))^2}_{=: \mu_2^\top \mathbf{R}_{22} \mu_2} d\tau \end{aligned} \quad (5.84)$$

(vgl. (5.2)) liefert der Converse-HJB-Ansatz das nichtlineare, eingangsaффine System

$$\begin{aligned} \dot{\mathbf{x}} &= \underbrace{\begin{bmatrix} -2x_1 + x_2 \\ -x_2 - \frac{1}{2}x_1 + \frac{1}{4}x_2 ((\cos(2x_1) + 2)^2 + (\sin(4x_1^2) + 2)^2) \end{bmatrix}}_{=: \mathbf{f}(\mathbf{x})} \\ &+ \underbrace{\begin{bmatrix} 0 \\ \cos(2x_1) + 2 \end{bmatrix}}_{=: \mathbf{g}_1(\mathbf{x})} u_1 + \underbrace{\begin{bmatrix} 0 \\ \sin(4x_1^2) + 2 \end{bmatrix}}_{=: \mathbf{g}_2(\mathbf{x})} u_2 \end{aligned} \quad (5.85)$$

(vgl. (5.1)) als mögliche Lösung¹¹⁰. Der initiale Systemzustand wird zu $\mathbf{x}_0 = [0 \ 0]^\top$ gesetzt¹¹¹. Die Basisfunktionen

$$\phi_1(\mathbf{x}) = \phi_2(\mathbf{x}) := [x_1^2 \ x_1 x_2 \ x_2^2]^\top, \quad (5.86)$$

¹⁰⁹ Ähnliche Testsysteme werden auch in der Literatur konstruiert, siehe beispielsweise [VL11], [LLW14].

¹¹⁰ Das sich im Zuge des Converse-HJB-Ansatzes ergebende Gleichungssystem ist unterbestimmt, da den beiden skalaren HJB-Gleichungen der Spieler mit den Funktionen $\mathbf{f}(\mathbf{x})$, $\mathbf{g}_1(\mathbf{x})$ und $\mathbf{g}_2(\mathbf{x})$ (jeweils $\in \mathbb{R}^2$) insgesamt sechs unbekannte Größen gegenüberstehen.

¹¹¹ Eine initiale Auslenkung des Systemzustands würde eine zusätzliche Anregung des Systems bedeuten und wird daher zugunsten der Vergleichbarkeit der nachfolgend untersuchten Anregungssignale vermieden. Generell muss jedoch nicht $\mathbf{x}_0 = \mathbf{0}$ gelten.

die Annahme 5.3 erfüllen, können die optimalen Value Functions (5.83) parametrieren und die als unbekannt angenommenen gesuchten Nash-Gewichte ergeben sich zu

$$\boldsymbol{w}_1^* = \begin{bmatrix} \frac{1}{2} \\ 0 \\ 1 \end{bmatrix}, \quad \boldsymbol{w}_2^* = \begin{bmatrix} \frac{1}{4} \\ 0 \\ \frac{1}{2} \end{bmatrix}. \quad (5.87)$$

5.7.2 Konstruktion geeigneter Anregungssignale

Um, wie in Abschnitt 5.6 beschrieben und in Abbildung 5.2 veranschaulicht, ein Anregungssignal $\boldsymbol{u}_{\text{ex}}$ zur Lösung von Problem 5.2 zu konstruieren, wird zunächst $\boldsymbol{x}_{\text{PE}}$ nach Annahme 5.4 gewählt. Anschließend wird die Menge Ω an Frequenzbedingungen nach Satz 5.2 berechnet. Nach Wahl der verbliebenen Freiheitsgrade $\boldsymbol{\omega}_c \in \Omega$ und $\boldsymbol{\nu} \in \mathbb{R}^n, \nu_j \neq 0$, wird schließlich $\boldsymbol{u}_{\text{ex}}$ mithilfe eines flachheitsbasierten Vorsteuerungsentwurfs berechnet.

5.7.2.1 Wahl geeigneter Anregungstrajektorien $\boldsymbol{x}_{\text{PE}}$

In Übereinstimmung mit Annahme 5.4 wird

$$\boldsymbol{x}_{\text{PE}} = \begin{bmatrix} \sin(\omega_1 t) + \sin(\omega_2 t) \\ \sin(\omega_3 t) \end{bmatrix} \quad (5.88)$$

gewählt, weshalb im Folgenden Bedingungen an $\boldsymbol{\omega} = [\omega_1 \ \omega_2 \ \omega_3]^\top$ untersucht werden. Die sich ergebenden Frequenzbedingungen Ω nach Satz 5.2 sind, wie in Lemma 5.6 gezeigt, durch Z Matrizen \boldsymbol{C}_z definiert und durch

$$\Omega = \left\{ \boldsymbol{\omega} : \bigwedge_{z=1}^Z \boldsymbol{C}_z \boldsymbol{\omega} \neq \mathbf{0} \right\}, \quad Z \in \mathbb{N}_{\geq 1}, \quad (5.89)$$

gegeben. Die Matrizen \boldsymbol{C}_z sind in Tabelle 5.1 aufgelistet¹¹². Anschaulich betrachtet ist jeder Frequenzvektor $\boldsymbol{\omega}$ geeignet, der keine der $Z = 21$ Gleichungen $\boldsymbol{C}_z \boldsymbol{\omega} = \mathbf{0}$, $z = 1, \dots, Z$, erfüllt. Durch die in Tabelle 5.1 gegebenen Matrizen \boldsymbol{C}_z werden 21 Ebenen $\boldsymbol{C}_z \boldsymbol{\omega} = \mathbf{0}$ im \mathbb{R}^3 beschrieben, die alle durch den Koordinatenursprung verlaufen und auf denen der Frequenzvektor $\boldsymbol{\omega}$ nach (5.89) nicht liegen darf. In Abbildung 5.3 sind exemplarisch sechs dieser Ebenen veranschaulicht. Die Menge Ω , für die nach Satz 5.2 folgt, dass $\boldsymbol{\sigma} = \dot{\boldsymbol{\phi}}(\boldsymbol{x}_{\text{PE}}) \forall \boldsymbol{\omega} \in \Omega$ PE ist, wird durch den Raum \mathbb{R}^3 ohne diese 21 Ebenen beschrieben.

In Übereinstimmung mit Lemma 5.7 existieren somit geeignete Frequenzen $\boldsymbol{\omega} \in \Omega$, da $\Omega \neq \emptyset$ gilt. Durch die Wahl der verbliebenen Freiheitsgrade $\boldsymbol{\omega}_c \in \Omega$ und $\boldsymbol{\nu} \in \mathbb{R}^n, \nu_j \neq 0$, nach Lemma 5.9 folgt somit

$$\bar{\boldsymbol{x}}_{\text{PE}}(t) = \begin{bmatrix} \nu_1 & 0 \\ 0 & \nu_2 \end{bmatrix} \begin{bmatrix} \sin(\omega_1 t) + \sin(\omega_2 t) \\ \sin(\omega_3 t) \end{bmatrix}. \quad (5.90)$$

¹¹² Die Berechnung der Matrizen \boldsymbol{C}_z , welche die Frequenzbedingungen Ω nach Abschnitt 5.4 definieren, erfolgt mithilfe des Computeralgebraprogramms MAPLE 2018.

$C_1 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$	$C_8 = \begin{bmatrix} 1 & 0 & -3 \end{bmatrix}$	$C_{15} = \begin{bmatrix} 1 & -1 & -2 \end{bmatrix}$
$C_2 = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$	$C_9 = \begin{bmatrix} 0 & 1 & -1 \end{bmatrix}$	$C_{16} = \begin{bmatrix} 0 & 1 & 1 \end{bmatrix}$
$C_3 = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}$	$C_{10} = \begin{bmatrix} 1 & -1 & 2 \end{bmatrix}$	$C_{17} = \begin{bmatrix} 1 & 1 & -2 \end{bmatrix}$
$C_4 = \begin{bmatrix} 1 & 1 & 2 \end{bmatrix}$	$C_{11} = \begin{bmatrix} 0 & 1 & -3 \end{bmatrix}$	$C_{18} = \begin{bmatrix} 1 & 0 & 3 \end{bmatrix}$
$C_5 = \begin{bmatrix} 0 & 1 & 3 \end{bmatrix}$	$C_{12} = \begin{bmatrix} 1 & -1 & 0 \end{bmatrix}$	$C_{19} = \begin{bmatrix} 1 & 1 & 0 \end{bmatrix}$
$C_6 = \begin{bmatrix} 1 & -\frac{1}{3} & 0 \end{bmatrix}$	$C_{13} = \begin{bmatrix} 1 & \frac{1}{3} & 0 \end{bmatrix}$	$C_{20} = \begin{bmatrix} 1 & -3 & 0 \end{bmatrix}$
$C_7 = \begin{bmatrix} 1 & 3 & 0 \end{bmatrix}$	$C_{14} = \begin{bmatrix} 1 & 0 & 1 \end{bmatrix}$	$C_{21} = \begin{bmatrix} 1 & 0 & -1 \end{bmatrix}$

Tabelle 5.1: Matrizen C_z der Frequenzbedingungen Ω für das Simulationsbeispiel nach Abschnitt 5.7.

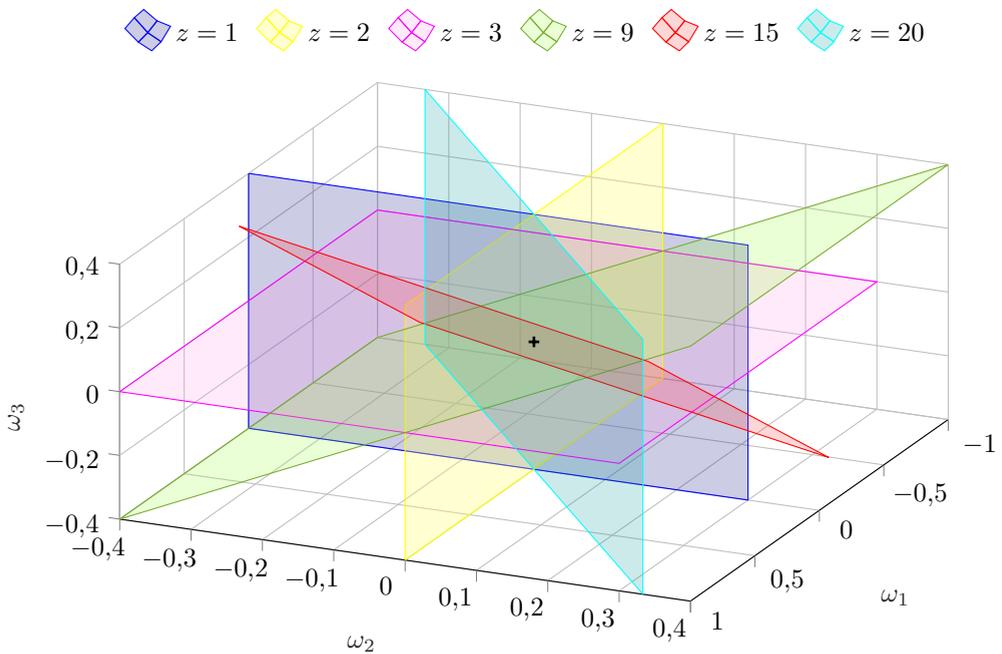


Abbildung 5.3: Grafische Veranschaulichung der Frequenzmenge Ω . Der Übersichtlichkeit halber wurden beispielhaft sechs der 21 Ebenen, welche durch C_z in Tabelle 5.1 definiert sind, und nach (5.89) verbotene Frequenzen beschreiben, skizziert. Der Koordinatenursprung ist durch das schwarze Kreuz gekennzeichnet.

5.7.2.2 Berechnung von u_{ex} basierend auf x_{PE}

Um, wie in Abschnitt 5.6 beschrieben, mithilfe eines flachheitsbasierten Vorsteuerungsentwurfs das Anregungssignal u_{ex} zu erhalten, wird zunächst das durch (5.85) gegebene Beispielsystem untersucht. Aufgrund von

$$\underbrace{g_2(x)}_{=:g_{\text{II}}(x)} = \underbrace{g_1(x)}_{=:g_{\text{I}}(x)} \underbrace{\frac{\sin(4x_1^2) + 2}{\cos(2x_1) + 2}}_{=: \gamma(x)} \quad (5.91)$$

(vgl. (5.68)) wird $g_{\text{I}}(x) = g_1(x)$ und $g_{\text{II}}(x) = g_2(x)$ gewählt und der zweite Eingang des Anregungssignals $u_{\text{ex}} = \begin{bmatrix} u_{\text{ex},1} & u_{\text{ex},2} \end{bmatrix}^{\top}$ wird zur Anregung nicht benötigt. Daher wird $u_{\text{ex},2} = u_{\text{ex,II}} = 0$ gesetzt und im Folgenden $u_{\text{ex},1} = u_{\text{ex,I}}$ betrachtet. Durch $h(x) = x_1$ ist aufgrund von

$$\begin{aligned} L_{g_1} h(x) &= 0 \\ L_{g_1} L_f h(x) &= \cos(2x_1) + 2 =: D(x) \end{aligned} \quad (5.92)$$

ein (fiktiver) Ausgang mit maximalem relativen Grad ($\delta = n = 2$) gegeben. Da außerdem $D(x) \neq 0$ gilt, ist das betrachtete System flach, eine flachheitsbasierte Vorsteuerung existiert [Lév09, S. 143 f.], und aus dem Sollverlauf der Anregungstrajektorie (5.90) ergibt sich das Anregungssignal $u_{\text{ex},1}$ ¹¹³. Die Gültigkeit von (5.79) in Annahme 5.6 kann, wie in Bemerkung 5.5 beschrieben, a priori bestätigt werden. Zudem gilt $\text{Rang}(M_i) = h_i$. Dies ist in Anhang C.3 ausgeführt. Somit resultiert eine Vorsteuerung der Form $u_{\text{ex},1}(t, \hat{w}_1^{[t]}, \hat{w}_2^{[t]})$, die eine exakte Systeminversion in Abhängigkeit der aktuellen Reglergewichte erlaubt.

5.7.3 Simulationsergebnisse

Neben den in Abschnitt 5.7.2 konstruierten Anregungssignalen $u_{\text{ex},1,j}$, $j \in \mathbb{N}$, die aus unterschiedlichen Beispielsignalen $\bar{x}_{\text{PE},j}$ resultieren, soll der Einfluss weißen Rauschens als Anregungssignal untersucht werden. Dazu wird das Signal $u_{\text{ex},1,w}(t)$ verwendet, das in der Simulation durch eine Pseudozufallsfolge mit $\Delta t = 0,01$ s und anschließendem Halteglied erzeugt wird, sodass eine mittelwertfreie weiße Gaußverteilung mit einer Standardabweichung von 2,8 approximiert wird¹¹⁴. Im Folgenden wird die Konvergenz der Critic-Gewichte $\hat{w}_i^{[t]}$, $i \in \mathcal{N}$, für diese unterschiedlichen Anregungssignale untersucht. Eine erste Untersuchung in Abschnitt 5.7.3.1 verwendet hierbei bezüglich der Maximalamplitude von σ_i normierte Lernraten. In einem zweiten Schritt werden in Abschnitt 5.7.3.2 Anregungssignale mit gleicher mittlerer Leistung bei Nutzung identischer Lernraten betrachtet.

¹¹³ Die analytischen Berechnungen erfolgen mithilfe des Computeralgebraprogramms MAPLE 2018. Die anschließenden numerischen Simulationen erfolgen in MATLAB.

¹¹⁴ Der Wert 2,8 der Standardabweichung ist so gewählt, dass die mittlere Signalleistung von $u_{\text{ex},1,w}(t)$ über der Simulationsdauer T_{sim} der mittleren Signalleistung von $u_{\text{ex},1,1}(t)$ entspricht. Dies spielt im späteren Vergleich der Effizienz von Anregungssignalen derselben mittleren Signalleistung eine Rolle.

5.7.3.1 Normierte Lernraten

Um den Einfluss unterschiedlicher Signalamplituden von σ_i zu Vergleichbarkeitszwecken der verschiedenen Anregungssignale auszugleichen, werden in einer ersten Simulation normierte Lernraten

$$\eta_{\text{norm},i} = \frac{\eta_v}{\sigma_{i,\max}^2}, \quad \forall i \in \{1, 2\}, \quad (5.93)$$

verwendet. Hierbei stellt $\eta_v > 0$ eine Konstante dar, die für alle Anregungssignale identisch gewählt ist, und $\sigma_{i,\max} \geq \|\sigma_i\|_2$ wird simulativ für jedes der verwendeten Anregungssignale ermittelt. Mit der normierten Lernrate $\eta_{\text{norm},i}$ aus (5.93) folgt für die Abschätzung der Konvergenzrate

$$\rho_i = 1 - \frac{2T_i\eta_v}{(1 + T_i\eta_v)^2} \frac{\alpha_{1,i}^2}{\sigma_{i,\max}^2}, \quad \forall i \in \{1, 2\}, \quad (5.94)$$

(vgl. (5.22)). Da nach (5.19) größere Amplituden von σ_i und somit größere Werte von $\sigma_{i,\max}$ in gleichem Maße zu einer Erhöhung von $\alpha_{1,i}$ führen, gleicht die Normierung der Lernrate unterschiedliche Signalamplituden von σ_i aus. Die Abschätzung der Konvergenzgeschwindigkeit ρ_i des Gradientenabstiegs nach (5.15) hängt (neben der für alle Anregungen identisch gewählten Größe η_v) somit nur noch von der Qualität der Erfüllung der PE-Bedingung ab, da ein höherer Grad $\alpha_{1,i}$ der Erfüllung der PE-Bedingung (bei gleichbleibendem $\sigma_{i,\max}$ und gleicher Zeitkonstante T_i) zu schnellerer Konvergenz führt.

Im Folgenden werden die drei Beispielsignale

$$\begin{aligned} \bar{x}_{\text{PE},1} &= \frac{1}{4} \begin{bmatrix} \sin(1s^{-1}t) + \sin(2s^{-1}t) \\ \sin(3s^{-1}t) \end{bmatrix}, \\ \bar{x}_{\text{PE},2} &= \frac{1}{4} \begin{bmatrix} \sin(\frac{1}{2}s^{-1}t) + \sin(1s^{-1}t) \\ \sin(2s^{-1}t) \end{bmatrix}, \\ \bar{x}_{\text{PE},3} &= 2\bar{x}_{\text{PE},2} \end{aligned} \quad (5.95)$$

gewählt (vgl. (5.90)), wobei gemäß Tabelle 5.1 $\omega \in \Omega$ gilt. Die zugehörigen Anregungssignale $u_{\text{ex},1,j}(t, \hat{w}_1^{[j]}, \hat{w}_2^{[j]})$, $\forall j \in \{1, 2, 3\}$, werden damit, wie in Abschnitt 5.6 und Abschnitt 5.7.2.2 beschrieben, berechnet. Die initialen Critic-Gewichte werden zu

$$\hat{w}_1^{[0]} = \begin{bmatrix} 1,783 \\ -2,33 \\ 2,215 \end{bmatrix}, \quad \hat{w}_2^{[0]} = \begin{bmatrix} 0,8916 \\ -1,165 \\ 1,107 \end{bmatrix} \quad (5.96)$$

gewählt. Die Regelgesetze $\hat{\mu}_i^{[0]}(\mathbf{x})$, $\forall i \in \mathcal{N}$, werden aus den in (5.96) gegebenen Gewichten über den Zusammenhang (5.24) berechnet. Diese Initialisierung entspricht stabilisierenden initialen Regelgesetzen, die jedoch nicht der Nash-Lösung (vgl. (5.87)) entsprechen. Mit dem

exemplarisch gewählten Wert $\eta_\nu = 10^4$ ergeben sich die in Tabelle 5.2 gezeigten normierten Lernraten $\eta_{\text{norm},i}$, $\forall i \in \mathcal{N}$, nach (5.93). Der Policy-Evaluation-Schritt erfolgt mithilfe des durch (5.15) beschriebenen Gradientenabstiegs. Als Abbruchbedingung der Policy Evaluation wird im Folgenden

$$\left\| \hat{\boldsymbol{w}}_i^{[l]}(t - \tau) - \hat{\boldsymbol{w}}_i^{[l]}(t) \right\|_2 < 10^{-3}, \quad (5.97)$$

$\forall i \in \mathcal{N}$, mit dem Designparameter $\tau = 40$ s, verwendet¹¹⁵. Der Policy-Improvement-Schritt geschieht anschließend mithilfe des Critic-Gewichts $\hat{\boldsymbol{w}}_i^{[l+1]}$ des jeweiligen Spielers $i \in \mathcal{N}$ basierend auf (5.24).

Abbildung 5.4 zeigt den Verlauf der Critic-Gewichte $\hat{\boldsymbol{w}}_i^{[l]}$, $i \in \{1, 2\}$, exemplarisch für die Verwendung der Anregungssignale $u_{\text{ex},1,1}$ und $u_{\text{ex},1,w}$. Die resultierenden Fehlernormen $\|\boldsymbol{w}_i^* - \hat{\boldsymbol{w}}_i\|_2$, $i \in \{1, 2\}$, für die Anregungssignale $u_{\text{ex},1,j}$ mit $j \in \{1, 2, 3\}$ und $u_{\text{ex},1,w}$ sind in Abbildung 5.5 gezeigt. Für alle vier betrachteten Anregungssignale konvergieren die adaptiven Optimalregler gegen die Nash-Lösung. Die Konvergenzzeiten T_{konv} , ab denen $\|\hat{\boldsymbol{w}}_i(t)\|_2 < 10^{-3}$, $\forall i \in \mathcal{N}$, gilt, sind Tabelle 5.2 zu entnehmen.

Das Eigenwertsignal $\lambda_{\min,2}(t)$ zur Überprüfung der PE-Eigenschaft nach Lemma 5.11, das für die vier verwendeten Anregungssignale resultiert, ist in Abbildung 5.6 gezeigt. Gemäß dem in Abschnitt 5.5 beschriebenen Verfahren können für den gewählten numerischen Schwellwert $\alpha_t = 10^{-4}$ für alle vier Anregungssignale Zeitpunkte T_t gefunden werden, für welche die Ungleichung (5.55) mit $T = T_t$ und $\alpha = \alpha_t$ erfüllt ist. Diese Zeitkonstanten T_t sind in Tabelle 5.2 gegeben. Die Eigenwertsignale $\lambda_{\min,1}(t)$ in Abbildung 5.7 zeigen, dass Lemma 5.10 mit $\alpha = 10^{-4}$ und $T = T_t$ nach Tabelle 5.2 erfüllt ist, und somit die Signale σ_i , $\forall i \in \{1, 2\}$, für die betrachteten Anregungssignale tatsächlich die PE-Bedingung erfüllen.

5.7.3.2 Gleiche mittlere Anregungssignalleistung

Während eine Normierung der Lernraten des Gradientenabstiegs nach (5.93) insbesondere Unterschiede der Amplitude von σ_i zu Vergleichszwecken ausgleicht, soll zudem ein weiteres

	$u_{\text{ex},1,1}(t)$	$u_{\text{ex},1,2}(t)$	$u_{\text{ex},1,3}(t)$	$u_{\text{ex},1,w}(t)$
$\eta_{\text{norm},i}, \forall i \in \{1, 2\}$	36,50	80	5,99	41,49
T_{konv}	241 s	497 s	490 s	3311 s
T_t	1,14 s	3,61 s	1,54 s	11,56 s

Tabelle 5.2: Normierte Lernraten $\eta_{\text{norm},i}$ nach (5.93), Konvergenzzeiten T_{konv} und Zeitpunkte T_t als Kandidaten für (5.55) mit $T = T_t$ und $\alpha = \alpha_t = 10^{-4}$ für die Anregungssignale $u_{\text{ex},1,j}(t)$ ($\forall j \in \{1, 2, 3\}$) und $u_{\text{ex},1,w}(t)$.

¹¹⁵ Grundsätzlich ist τ hierbei ein frei wählbarer Designparameter. Da nach Lemma 5.1 jedoch zunächst nur eine Verbesserung der geschätzten Critic-Gewichte $\hat{\boldsymbol{w}}_i^{[l+1]}$ für ein Intervall der Länge T_i gewährleistet wird, sollte τ nicht kleiner als T_i gewählt werden.

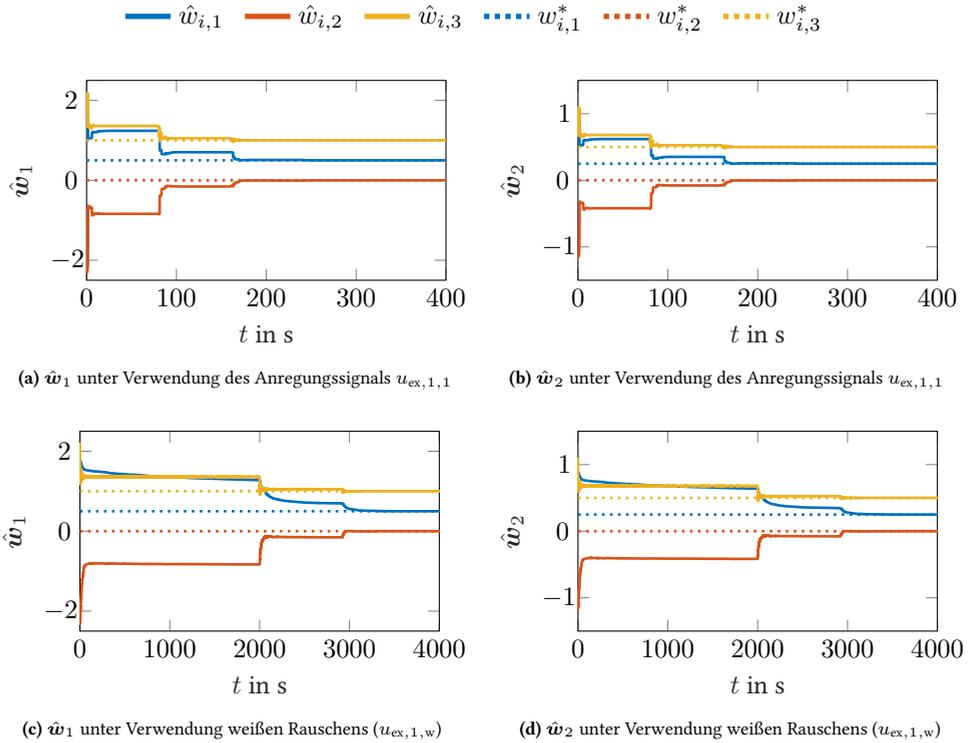


Abbildung 5.4: Verlauf der Critic-Gewichte \hat{w}_1 und \hat{w}_2 unter Verwendung des Anregungssignals $u_{ex,1,1}$ (oben) bzw. weißen Rauschens $u_{ex,1,w}$ (unten) mit nach Tabelle 5.2 normierten Lernraten. Zu beachten ist insbesondere die um Faktor 10 unterschiedlich skalierte Zeitachse.

Kriterium zur Untersuchung der Effizienz der Anregungssignale betrachtet werden. Hierbei wird ein Anregungssignal als effizienter im Vergleich zu einem anderen Signal bezeichnet, wenn es bei identischer mittlerer Signalleistung

$$\frac{1}{T_{sim}} \int_0^{T_{sim}} \mathbf{u}_{ex}^T(t) \mathbf{u}_{ex}(t) dt \quad (5.98)$$

und gleicher Lernrate η_i zu schnellerer Konvergenz der Critic-Gewichte \hat{w}_i führt. Somit wird mit diesem Bewertungskriterium insbesondere bestraft, wenn ein Anregungssignal $\mathbf{u}_{ex}(t)$ Signalanteile besitzt, die, beispielsweise durch eine mögliche Tiefpasscharakteristik des Systems, nicht zu einer Auslenkung des Systemzustands und somit zur Anregung von σ_i beitragen.

Daher sollen im Folgenden die Anregungssignale $u_{ex,1,1}$ und $u_{ex,1,w}$, welche die gleiche mittlere Signalleistung nach (5.98) aufweisen und in Abbildung 5.8 gezeigt sind, bei identischer Lernrate $\eta_i = 1$, $i \in \mathcal{N}$, hinsichtlich der Konvergenz der Critic-Gewichte \hat{w}_i untersucht werden. Wie Abbildung 5.9 veranschaulicht, lenkt, bei gleicher mittlerer Signalleistung, das Anregungssignal $u_{ex,1,1}$ den Systemzustand \mathbf{x} stärker aus als weißes Rauschen. Für diese

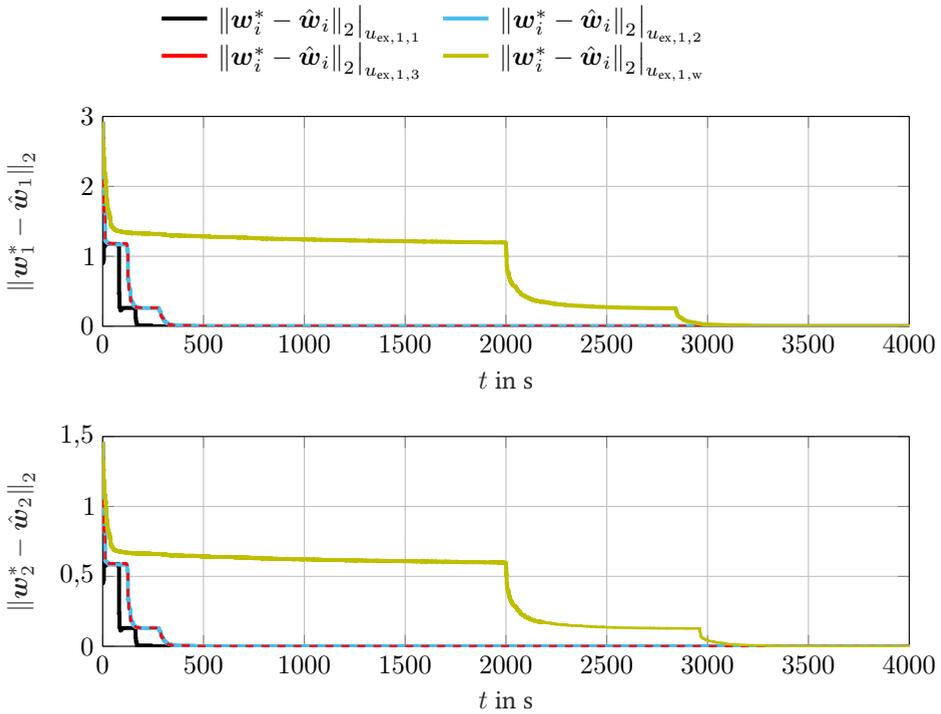


Abbildung 5.5: Norm des Fehlers der Critic-Gewichte $\|w_i^* - \hat{w}_i\|_2, i \in \{1, 2\}$, für die Anregungssignale $u_{ex,1,j}, j \in \{1, 2, 3\}$, (Sinussignale) und $u_{ex,1,w}$ (weißes Rauschen) mit nach Tabelle 5.2 normierten Lernraten.

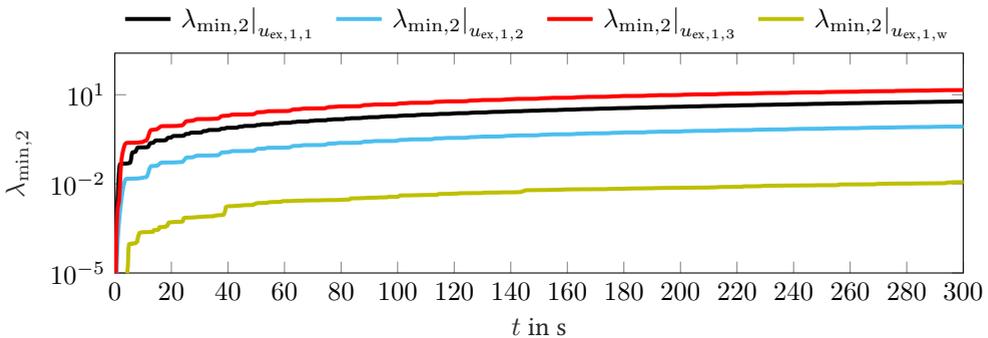


Abbildung 5.6: Eigenwertsignal $\lambda_{\min,2}(t)$ zur Überprüfung der PE-Eigenschaft nach Lemma 5.11. Zu beachten ist die logarithmische Darstellung der Ordinate.

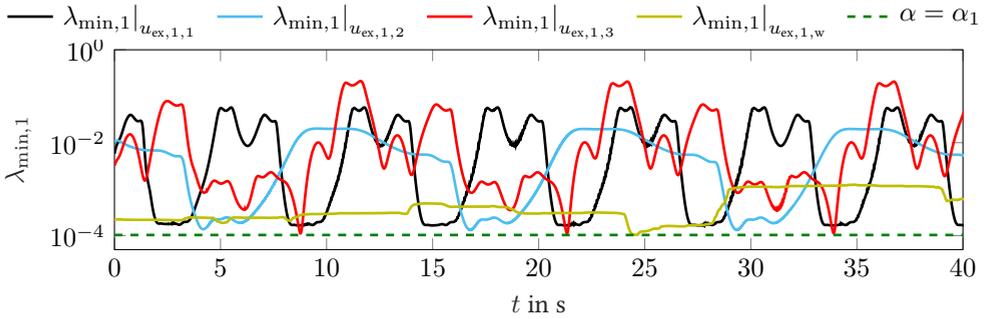


Abbildung 5.7: Eigenwertensignal $\lambda_{\min,1}(t)$ zur Überprüfung der PE-Eigenschaft nach Lemma 5.10. Der Übersichtlichkeit halber sind nur die ersten 30 s gezeigt. Die gestrichelte horizontale Linie entspricht dem gewählten Schwellwert $\alpha_t = 10^{-4}$. Zu beachten ist die logarithmische Darstellung der Ordinate.

beiden Anregungssignale wird die Konvergenz eines einzelnen Policy-Evaluation-Schrittes im Folgenden betrachtet. Die initialen Critic-Gewichte $\hat{w}_i^{[0]}$, $i \in \mathcal{N}$, und Regelgesetze $\hat{\mu}_i^{[0]}(\mathbf{x})$ werden wie in Abschnitt 5.7.3.1 gewählt.

Der Verlauf der Critic-Gewichtsfehler $\left\| \mathbf{w}_i^{[1]} - \hat{\mathbf{w}}_i^{[1]} \right\|_2$, $i \in \mathcal{N}$, ist in Abbildung 5.10 gegeben. Bei der hier untersuchten Lernrate von $\eta_i = 1$ genügt die Simulationsdauer von 4000 s für die Anregung mit weißem Rauschen, d. h. unter Verwendung von $u_{\text{ex},1,\text{w}}$, nicht, um Konvergenz der Critic-Gewichte $\hat{w}_i^{[1]}$ zu erreichen. So gilt in diesem Fall nach $t = 4000$ s

$$\left\| \mathbf{w}_1^{[1]} - \hat{\mathbf{w}}_1^{[1]} \right\|_2 = 0,31 \quad \text{und} \quad \left\| \mathbf{w}_2^{[1]} - \hat{\mathbf{w}}_2^{[1]} \right\|_2 = 0,16. \quad (5.99)$$

Unter Verwendung des Anregungssignals $u_{\text{ex},1,1}$ folgt hingegen bereits nach $t = 338$ s, dass $\left\| \mathbf{w}_i^{[1]} - \hat{\mathbf{w}}_i^{[1]} \right\|_2 < 0,01$, $\forall i \in \{1, 2\}$, gilt.

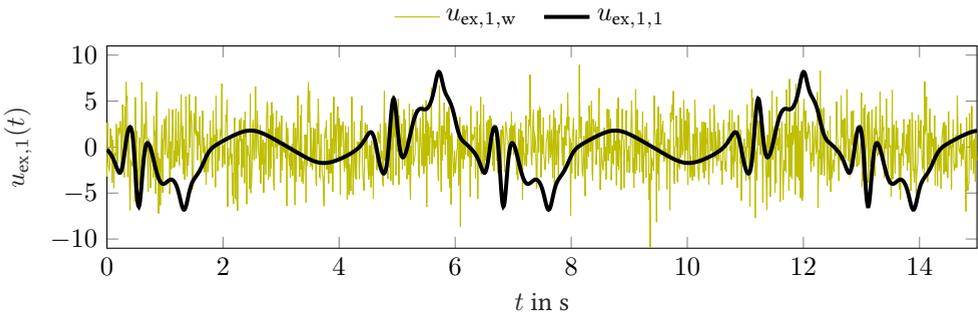


Abbildung 5.8: Anregungssignal $u_{\text{ex},1}$ für $\bar{x}_{\text{PE},1}$ nach (5.95) (d. h. $u_{\text{ex},1,1}$) und $u_{\text{ex},1,\text{w}}$ (d. h. weißes Rauschen) für die ersten 15 s der Simulation. Durch die Standardabweichung von 2,8 bei der Erzeugung von $u_{\text{ex},1,\text{w}}$ haben beide Anregungssignale die gleiche mittlere Signalleistung.

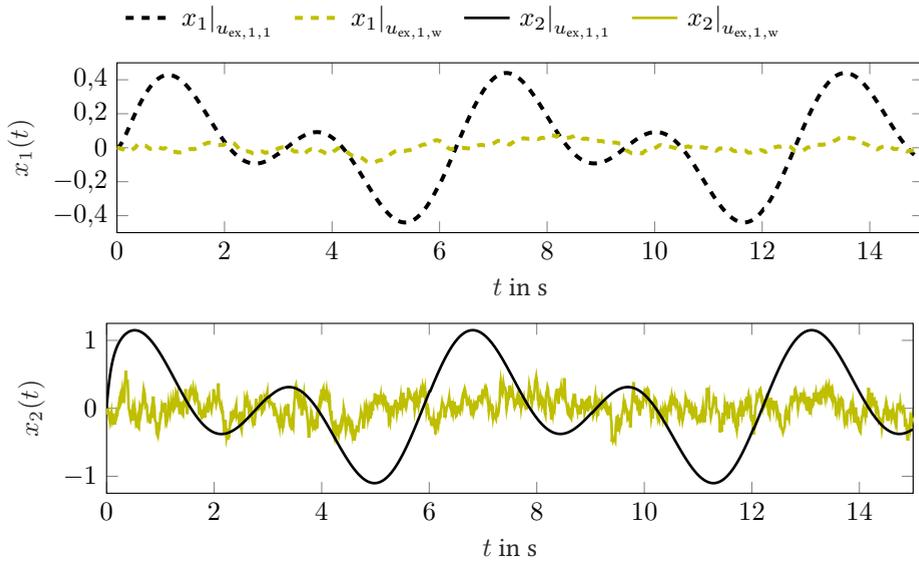


Abbildung 5.9: Verlauf der Zustandsgröße $x(t)$ für die Anregungssignale $u_{ex,1,1}$ (Sinussignale) und $u_{ex,1,w}$ (weißes Rauschen), welche die gleiche mittlere Anregungssignalleistung aufweisen.

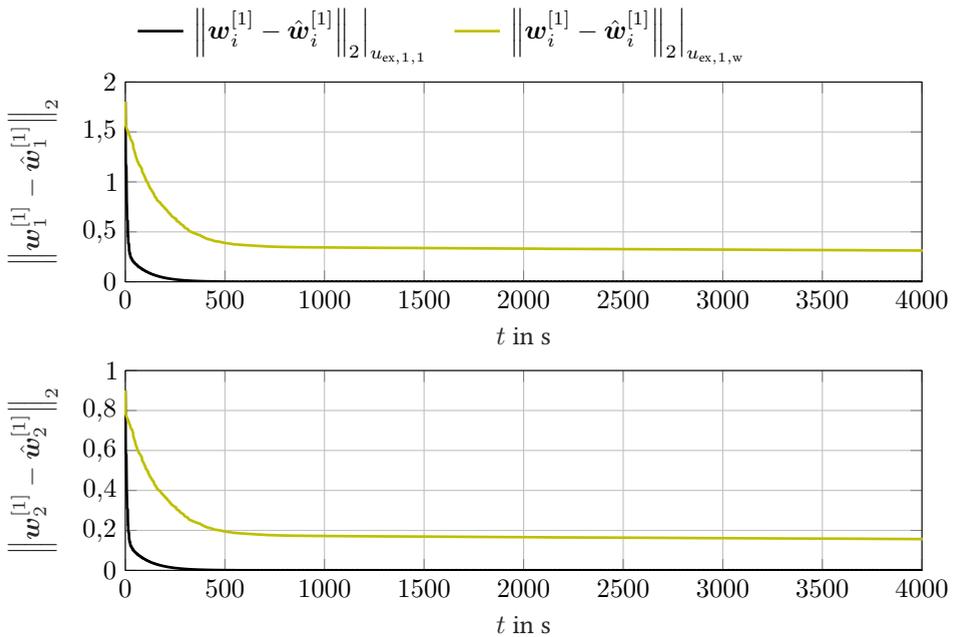


Abbildung 5.10: Verlauf der Critic-Gewichtsfehlnorm $\|w_i^{[1]} - \hat{w}_i^{[1]}\|_2$, $i \in \mathcal{N}$, unter Verwendung des Anregungssignals $u_{ex,1,1}$ bzw. weißen Rauschens $u_{ex,1,w}$, welche die gleiche mittlere Anregungssignalleistung aufweisen, bei identischen Lernraten $\eta_i = 1$.

5.8 Diskussion

Die Erfüllung der PE-Eigenschaft spielt für die Konvergenz ADP-basierter Regelungsmethoden eine zentrale Rolle. Durch den nichtlinearen Zusammenhang zwischen dem Systemzustand \boldsymbol{x} und den für den Adaptionvorgang relevanten Signalen $\boldsymbol{\sigma}_i$ sind die in der Literatur vorhandenen theoretischen Aussagen zu linearen Transformationen von PE-Signalen nicht ausreichend, um Rückschlüsse darüber zu ziehen, für welche Zustandstrajektorien das Signal $\boldsymbol{\sigma}_i$ die PE-Eigenschaft erfüllt.

Die vorliegende Arbeit schließt diese Lücke, indem für polynomielle Basisfunktionen $\phi(\boldsymbol{x})$ erstmalig ein allgemeingültiges hinreichendes Kriterium an den Systemzustand \boldsymbol{x} präsentiert wird, sodass $\boldsymbol{\sigma}_i$ die PE-Bedingung erfüllt. Proposition 5.1 formuliert hierzu zunächst allgemeine Forderungen an eine Menge Ω an Frequenzbedingungen, um die Erfüllung der PE-Bedingung von $\boldsymbol{\sigma}_i$ zu gewährleisten. Durch Satz 5.2 ist schließlich eine konkrete Entwurfsmethode für eine solche Menge Ω gegeben. Damit gilt für jeden beliebigen Vektor $\boldsymbol{\omega} \in \Omega$, dass das Anregungstrajektorienignal $\boldsymbol{x}_{\text{PE}}(t)$ SR bezüglich $\dot{\phi}(\cdot)$ und somit $\boldsymbol{\sigma} = \dot{\phi}(\boldsymbol{x}_{\text{PE}})$ PE ist. Da die genauen Amplituden der Sinus- und Kosinusterme in der Anregungstrajektorie $\boldsymbol{x}_{\text{PE}}$ bei der Analyse nicht berücksichtigt wurden, sind die durch Ω beschriebenen Bedingungen (vgl. Proposition 5.1 und Satz 5.2) hinreichend. Somit können bei geeigneten Amplitudenverhältnissen gegebenenfalls $\boldsymbol{\omega} \notin \Omega$ existieren, für die $\boldsymbol{\sigma}_i$ dennoch die PE-Bedingung erfüllt. Für den praktischen Entwurf von Anregungssignalen ist dies jedoch unerheblich. Die Simulationsergebnisse bestätigen die Konvergenz einer gradientenbasierten Policy Evaluation, wenn $\boldsymbol{\sigma}_i \forall i \in \mathcal{N}$ die PE-Bedingung erfüllt. Dies gilt sowohl für eine Anregung mit weißem Rauschen als auch für Anregungssignale, die, wie in Abschnitt 5.6 beschrieben, auf Basis der theoretischen Erkenntnisse aus Satz 5.2 konstruiert sind. Für die praktische Überprüfung der Erfüllung der PE-Bedingung wurden in Abschnitt 5.5 zudem die Eigenwertsignale $\lambda_{\min,2}$ und $\lambda_{\min,1}$ präsentiert.

Die Untersuchungen des vorliegenden Kapitels beschränken sich aufgrund von Annahme 5.3 bislang auf polynomielle Basisfunktionen $\phi(\boldsymbol{x})$. Für den Fall, dass diese Annahme verletzt ist, könnte eine Taylor-Approximation von $\phi(\boldsymbol{x})$ zukünftig weitere Analysen ermöglichen. Bei hinreichend kleinem Restterm dieser Approximation könnten dann, ähnlich wie in Lemma 5.12, Rückschlüsse über die Erfüllung der PE-Bedingung gezogen werden. Auch sind die theoretischen Analysen bislang auf den Fall der zeitkontinuierlichen PE-Bedingung (2.27) begrenzt. Sie könnten jedoch als Ausgangspunkt für die Untersuchung der zeitdiskreten PE-Bedingung (2.29) dienen.

Aus den theoretischen Analysen dieses Kapitels lassen sich schließlich wesentliche Erkenntnisse schlussfolgern, die für den Entwurf von Anregungssignalen für ADP-basierte Regler hilfreich sind. So kann bereits eine geringe Anzahl an Frequenzen $\boldsymbol{\omega} \in \Omega$ im Systemzustand \boldsymbol{x} ausreichend sein, damit $\boldsymbol{\sigma}_i$ PE ist. Hierbei lässt Satz 5.2 Freiheiten bei der Wahl der Frequenzen und Amplituden in $\boldsymbol{x}_{\text{PE}}(t)$. Anstatt, wie bei der Anregung mit weißem Rauschen, alle Frequenzen gleich stark im Anregungssignal zu verwenden, können diese Freiheitsgrade dazu dienen,

die Anregungssignalleistung effizient zu nutzen. Auch könnten die Freiheitsgrade gegebenenfalls für eine zukünftige Optimierung der Konvergenzeigenschaften verwendet werden. Liegt darüber hinaus sogar (ggf. teilweise) Vorwissen über das Übertragungsverhalten des betrachteten Systems vor (beispielsweise über dessen Tiefpass- oder Bandpasscharakteristik), so können die Anregungsfrequenzen entsprechend gewählt werden.

Zusammenfassend verdeutlicht dieses Kapitel, dass die präsentierten Anregungssignale, die ihre Signalleistung auf niedrigere Bereiche des Frequenzspektrums konzentrieren, deutlich effizienter zur Anregung beitragen können, als dies bei der Verwendung von weißem Rauschen der Fall ist. Dies zeigt sich selbst bei der Nutzung normierter Lernraten durch eine wesentlich schnellere Konvergenz der Critic-Gewichte und somit des Policy-Iteration-Algorithmus (85,0 %–92,7 % Verbesserung, vgl. Tabelle 5.2 und Abbildung 5.5). Bei gleicher mittlerer Anregungssignalleistung nach (5.98) zeigt sich bei identischen Lernraten $\eta_i = 1$ im betrachteten Simulationsbeispiel die Überlegenheit des analytisch berechneten Anregungssignals im Vergleich zu weißem Rauschen sogar noch deutlicher (vgl. Abbildung 5.10). Somit liefert dieses Kapitel Antworten auf die in Abschnitt 2.4.2 formulierte Forschungsfrage 2 hinsichtlich der Bedingungen an Systemzustände, um die PE-Eigenschaft zu erfüllen und Konvergenz der betrachteten ADP-Methode zu gewährleisten.

6 Reale Anwendung ADP-basierter Solltrajektorienfolgeregler

Nach den theoretischen Beiträgen zu zeitdiskreten (Kapitel 3) und zeitkontinuierlichen (Kapitel 4) ADP-kompatiblen Solltrajektorien- und Konvergenzbedingungen (Kapitel 5) werden in diesem Kapitel zwei reale Anwendungsbeispiele ADP-basierter Solltrajektorienfolgeregler betrachtet. Da bei diesen realen Systemen Messungen der Ausgangs- bzw. Zustandsgrößen zu diskreten Zeitpunkten stattfinden, wird eine zeitdiskrete Darstellung verwendet. Unter den in Kapitel 3 präsentierten zeitdiskreten Methoden erscheint insbesondere die in Abschnitt 3.2 vorgestellte ADP-kompatible parametrisierte Referenztrajektorien- und Konvergenzbedingungen (Kapitel 5) Darstellung, da durch die parametrische lokale Approximation des Solltrajektorienverlaufs eine kompakte Darstellung vorliegt (vgl. auch Abschnitt 3.4).

In Abschnitt 6.1 wird ein modellfreier Actor-Critic-Ansatz angewandt, um eine Längsregelung eines realen Fahrzeugs zu realisieren, wobei online, d. h. während der Fahrt, die Reglergewichte adaptiert werden. Motiviert ist dieses Anwendungsbeispiel durch eine wachsende Modellvielfalt und den Wunsch, ausgehend von einer initialen Reglerparametrierung, die beispielsweise aus einem vorhandenen modellbasierten Reglerentwurf resultiert, eine automatisierte Feinabstimmung vornehmen zu können. Durch die Verwendung der parametrischen Darstellung des Sollgeschwindigkeitsverlaufs wird dabei zusätzlich zu einem Ausgangsrückführungsterm ein Vorsteuerterm gelernt. Abschnitt 6.2 betrachtet schließlich ein reales Ball-auf-Platte-System. Für dieses beliebte regelungstechnische Benchmarksystem wird ebenfalls ein modellfreier Trajektorienfolgeregler trainiert und mit einem modellbasierten Ansatz verglichen. Betrachtet werden dabei modellfreie ADP-Solltrajektorienfolgeregler, die den Sollpositionsverlauf einerseits durch ein Polynom zweiten Grades und andererseits durch eine stationäre Sollzustandsvorgabe beschreiben, sowie ein modellbasierter Vergleichsregler. Hierbei offenbaren sich die Vorteile des präsentierten Ansatzes, der ohne aufwendige Modellbildung zu den geringsten Kosten im Sinne des zugrunde liegenden Gütefunktionalis führt.

6.1 Modellfreie, adaptive Längsregelung eines realen Fahrzeugs

In diesem Abschnitt wird die in Kapitel 3 vorgestellte ADP-kompatible parametrisierte Referenztrajektorien- und Konvergenzbedingungen (Kapitel 5) Darstellung verwendet, um einen Geschwindigkeitsregler, der einem vorgegebenen

Geschwindigkeitsprofil folgen soll, mithilfe einer Online-ADP-Methode in einem realen Fahrzeug zu trainieren¹¹⁶. Das verwendete Versuchsfahrzeug ist in Abbildung 6.1 zu sehen.

Bei der ADP-basierten Geschwindigkeitsregelung eines realen Fahrzeugs ergeben sich zwei wesentliche Herausforderungen. Erstens muss der Sollgeschwindigkeitsverlauf durch eine ADP-kompatible Solltrajektorienarstellung repräsentiert und in den ADP-Formalismus integriert werden (vgl. Abschnitt 2.2 und Kapitel 3). Zweitens basiert die Längsdynamik des Versuchsfahrzeugs auf einer Aktuatorik (Antriebsstrang und Bremssystem) mit vergleichsweise langsamen dynamischen Vorgängen und teilweise nicht gemessenen internen Zuständen (vgl. [PRH19]). Anstelle des Systemzustands liegt nur eine Ausgangsgröße in Form der gemessenen Fahrzeuggeschwindigkeit vor. Würden die nicht messbaren Zustände der Aktuatordynamik jedoch schlicht ignoriert, wäre damit die Markov-Annahme verletzt. Die in der Literatur vorhandenen ADP- oder RL-basierten Geschwindigkeitsregler, die zudem meist auf Simulationen basieren oder nicht online lernen (vgl. Abschnitt 2.2.4), ignorieren dieses Problem bislang. Beispielhaft seien [BK18], [NCH08], [DCd11], [PT12], [WXL⁺14] und [WZLD15] zu nennen. Im Folgenden wird hingegen, wie in den Arbeiten von Puccetti et al. [PRH19], [PKRH20] vorgestellt und in Anhang D.1 motiviert, ein FIR-Filter¹¹⁷ verwendet, um den Lernvorgang durch einen geschätzten Hilfszustand zu stützen.

Die zentralen Ideen sind somit die Verwendung der in Abschnitt 3.2 vorgestellten zeitdiskreten ADP-kompatiblen parametrisierten Referenztrajektorie sowie eines FIR-Filters zur Rekonstruktion fehlender Zustandsinformationen. Ersteres ermöglicht, nicht nur die aktuelle Wunschgeschwindigkeit, sondern eine lokale Approximation des aktuellen Wunschgeschwindigkeitspro-



Abbildung 6.1: Versuchsfahrzeug (*BMW 740Li*), das für die modellfreie, adaptive Längsregelung mittels eines Online-ADP-Ansatzes verwendet wurde. Bildquelle: *BMW Press Club* (Zugriff am 19.10.2021).

¹¹⁶ Ergebnisse dieses Abschnitts wurden im Rahmen zweier Konferenzbeiträge veröffentlicht [KPRH20], [PKRH20].

¹¹⁷ (engl.): *finite impulse response filter*. Filter mit endlicher Impulsantwort.

files explizit zu berücksichtigen. Die Off-Policy-Charakteristik (vgl. Abschnitt 2.1.4.4) der verwendeten Actor-Critic-Methode erlaubt zudem einerseits, das Wunschgeschwindigkeitsprofil während des Trainingsvorgangs künstlich durch zusätzliches Anregungsrauschen des Sollverlaufs zu überlagern, und andererseits die Verwendung von Experience Replay [MKS⁺13], [MKS⁺15].

Nachfolgend wird zunächst die zugrunde liegende Problemstellung formuliert, bevor der verwendete modellfreie ADP-Solltrajektorienregelungsalgorithmus mit Zustandsrekonstruktion vorgestellt wird. Anschließend werden Trainingsergebnisse im Realfahrzeug sowie Validierungsfahrten präsentiert. Eine abschließende Diskussion rundet den Abschnitt ab.

6.1.1 Problemstellung

Ziel ist ein datenbasierter, selbstlernender ADP-Algorithmus, der in einem Realfahrzeug umgesetzt wird und online einen Geschwindigkeitsregler adaptiert, der einer ADP-kompatiblen Approximation gewünschter Sollgeschwindigkeitsverläufe optimal im Sinne eines gegebenen Gütefunktional folgt.

Der ADP-Regler verfügt dabei jedoch über kein Modell der Longitudinaldynamik, die Antriebsstrang, Bremssystem und eine unterlagerte Regelung beinhaltet. Wie in Abbildung 6.2 dargestellt, beschreibt diese Longitudinaldynamik den Zusammenhang zwischen der Eingangsgröße u_k , die der angeforderten Beschleunigung im Zeitschritt k entspricht, und der tatsächlichen Beschleunigung a_k , die jedoch nicht gemessen wird¹¹⁸. Stattdessen steht dem ADP-Regler lediglich eine Messung der Geschwindigkeit y_k zur Verfügung.

Zum Zeitschritt k ist der Sollgeschwindigkeitsverlauf auf einem Vorausschauhorizont der Länge n_h durch $y_{r,\text{soll},k}, y_{r,\text{soll},k+1}, \dots, y_{r,\text{soll},k+n_h-1}$ beschrieben. Aus diesem zunächst beliebigen Geschwindigkeitsprofil wird in jedem Zeitschritt k eine nach Definition 3.1 ADP-kompatible lokale Approximation

$$y_r(\mathbf{z}_k^{(\kappa)}) := y_r(\mathbf{z}_k, \kappa) = \mathbf{z}_k^\top \boldsymbol{\rho}(\kappa) \quad (6.1)$$

(vgl. (3.15)) des Solltrajektorienverlaufs berechnet. Dabei gewichtet \mathbf{z}_k die gegebenen Basisfunktionen $\boldsymbol{\rho}(\kappa)$ und κ bezeichnet einen Zeitindex. Gesucht ist ein Regelgesetz $\mu(y_\kappa, \mathbf{z}_k^{(\kappa)})$, das die Value Function

$$\begin{aligned} V^\mu(y_k, \mathbf{z}_k) &= \sum_{\kappa=0}^{\infty} \gamma^\kappa r(y_{k+\kappa}, y_r(\mathbf{z}_k^{(\kappa)}), \mu(y_{k+\kappa}, \mathbf{z}_k^{(\kappa)})) \\ &= \sum_{\kappa=0}^{\infty} \gamma^\kappa \left(Q_Y \left(y_{k+\kappa} - y_r(\mathbf{z}_k^{(\kappa)}) \right)^2 + R \left(\mu \left(y_{k+\kappa}, \mathbf{z}_k^{(\kappa)} \right) \right)^2 \right) \end{aligned} \quad (6.2)$$

¹¹⁸ Abhängig vom Entwurfsziel des Reglers und anderen Faktoren, wie beispielsweise äußeren Störeinflüssen, kann es sinnvoll sein, dass die Stellgröße u_k von der numerischen Differenzierung des Sollgeschwindigkeitsprofils abweicht. Das Ziel des ADP-Reglers ist daher, die optimale Stellgröße u_k bezüglich eines vorgegebenen Gütefunktional zu erlernen.

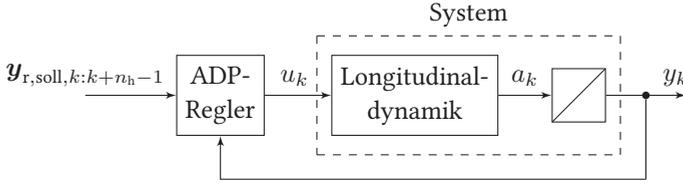


Abbildung 6.2: Struktur der gegebenen Problemstellung zur Längsregelung eines realen Fahrzeugs. Basierend auf der Ausgangsgröße y_k , die der gemessenen Geschwindigkeit des Fahrzeugs entspricht, sowie dem Sollgeschwindigkeitsprofil $\mathbf{y}_{r,soll,k:k+n_h-1} := [y_{r,soll,k} \quad y_{r,soll,k+1} \quad \dots \quad y_{r,soll,k+n_h-1}]$ und einem internen Kostensignal $r(\cdot)$ soll der Regler adaptiert werden.

minimiert (vgl. (3.37)). Hierbei wird mit $Q_y > 0$ und $R > 0$ sowohl die Abweichung der Geschwindigkeit $y_{k+\kappa}$ im Zeitschritt $k + \kappa$ von der approximierten Sollgeschwindigkeit $y_r(z_k^{(\kappa)})$ als auch der Stellaufwand quadratisch bestraft. Eine im Sinne der ADP-kompatiblen Solltrajektorienendarstellung resultierende Verschiebung des Parametervektors z_k um κ Zeitschritte ist dabei durch

$$z_k^{(\kappa)} = z_k D(\kappa) \quad (6.3)$$

(vgl. Definition 3.2) gegeben. $D(\kappa)$ stellt eine Verschiebungsmatrix dar, die so gewählt wird, dass

$$y_r(z_k^{(\kappa+j)}) = y_r(z_k^{(\kappa)}, j) = y_r(z_k, \kappa + j), \quad \forall \kappa, j \in \mathbb{N}_{\geq 0}, \quad (6.4)$$

gilt.

Bemerkung 6.1

Die ADP-kompatible Approximation des Sollgeschwindigkeitsverlaufs in (6.1) durch den endlichdimensionalen Vektor z_k hat zur Folge, dass die Value Function V^μ in (6.2) trotz des unendlichen Optimierungshorizonts durch einen funktionalen Zusammenhang der Form $V^\mu(y_k, z_k)$ beschrieben werden kann (vgl. Proposition 3.1).

6.1.2 Modellfreier ADP-Solltrajektorienfolgeregler mit Zustandsrekonstruktion

In diesem Abschnitt wird der verwendete modellfreie ADP-Solltrajektorienfolgeregler mit Zustandsrekonstruktion und Vorsteuerterm präsentiert. Grundsätzlich finden hierbei zwei wesentliche Modifikationen statt. Erstens wird in Abschnitt 6.1.2.1 eine erweiterte Q-Funktion definiert. Neben der gemessenen Geschwindigkeit y_k hängt diese Q-Funktion vom Vektor z_k ab, der im Zeitschritt k den approximierten Sollgeschwindigkeitsverlauf beschreibt (vgl. (6.1)). Zweitens werden dieser Q-Funktion zusätzlich vergangene Stellgrößen zur Verfügung gestellt, um, wie in Abschnitt 6.1.2.2 beschrieben, einen Hilfszustand \bar{x}_k zu schätzen, der den

Lernprozess stützt¹¹⁹. Die Schätzung dieses Hilfszustands erfolgt mithilfe eines FIR-Filters. Schließlich wird in Abschnitt 6.1.2.3 eine Übersicht über den verwendeten Algorithmus des ADP-basierten Geschwindigkeitsreglers präsentiert.

6.1.2.1 Q-Function und Solltrajektorienapproximation

Die zur Value Function V^μ nach (6.2) gehörende Q-Function ergibt sich zu¹²⁰

$$Q^\mu(y_k, z_k, u_k) = r(y_k, y_r(z_k), u_k) + \gamma Q^\mu(y_{k+1}, z_k^{(1)}, \mu(y_{k+1}, z_k^{(1)})). \quad (6.5)$$

Der zu trainierende Geschwindigkeitsregler basiert im Folgenden auf einer Approximation erster Ordnung (Polynom vom Grad $d = 1$) des Sollgeschwindigkeitsprofils. Somit gilt

$$\rho(\kappa) = \begin{bmatrix} \kappa \Delta t \\ 1 \end{bmatrix}, \quad D(\kappa) = \begin{bmatrix} 1 & \kappa \Delta t \\ 0 & 1 \end{bmatrix}, \quad (6.6)$$

wobei Δt die Abtastzeit des Systems beschreibt¹²¹. Zu Vergleichszwecken¹²² wird ein zweiter Regler betrachtet, der in jedem Zeitschritt k den Sollgeschwindigkeitsverlauf durch eine konstante Approximation (Polynom vom Grad $d = 0$) beschreibt, d. h.

$$\rho(\kappa) = 1, \quad D(\kappa) = 1. \quad (6.7)$$

Analog zu (3.72) mit $\beta = 1$ wird der Referenztrajektorienparameter z_k durch

$$z_k^\top = [z_{1,k} \quad \dots \quad z_{d+1,k}] = [y_{r,\text{soll},k} \quad y_{r,\text{soll},k+1} \quad \dots \quad y_{r,\text{soll},k+n_h-1}] P_{\text{LS}} \quad (6.8)$$

berechnet. Hierbei kann die durch

$$P_{\text{LS}} := \begin{bmatrix} \rho^\top(0) \\ \rho^\top(1) \\ \vdots \\ \rho^\top(n_h - 1) \end{bmatrix} \left(\begin{bmatrix} \rho^\top(0) \\ \rho^\top(1) \\ \vdots \\ \rho^\top(n_h - 1) \end{bmatrix}^\top \begin{bmatrix} \rho^\top(0) \\ \rho^\top(1) \\ \vdots \\ \rho^\top(n_h - 1) \end{bmatrix} \right)^{-1} \quad (6.9)$$

¹¹⁹ Dieser Hilfszustand muss dabei nicht notwendigerweise physikalisch interpretierbar sein. Neben der nicht gemessenen Beschleunigung a_k können weitere interne Größen der Longitudinaldynamik Einfluss auf \bar{x}_k haben.

¹²⁰ Grundsätzlich kann sich der Parameter z_k , der den Solltrajektorienverlauf beschreibt, beliebig mit der Zeit k ändern. Daher würde die Verwendung von z_k und z_{k+1} in (6.5) die Markov-Annahme verletzen. Aufgrund der Verschiebungsmatrix $D(\kappa)$ in (6.3), welche die durch z_k gegebene Approximation des Solltrajektorienverlaufs um einen Zeitschritt propagiert, ist die Markov-Annahme für den Zustandsübergang von $[y_k \quad z_k^\top]^\top$ nach $[y_{k+1} \quad z_k^{(1)\top}]^\top$ gewährleistet. Deshalb wird in (6.5) $z_k^{(1)}$ verwendet (vgl. auch Bemerkung 3.2).

¹²¹ Die verwendete Approximation erster Ordnung hat sich für dieses Anwendungsbeispiel als weniger anfällig gegenüber Oszillationen erwiesen und der zugehörige ADP-Regler ist zudem im Vergleich zu Approximatoren höherer Ordnung einfacher zu trainieren – insbesondere im Hinblick auf die zusätzlich benötigte Zustandsrekonstruktion durch ein FIR-Filter.

¹²² Ziel der vorliegenden Arbeit ist das Aufzeigen der grundsätzlichen Anwendbarkeit der vorgestellten Methode. Für eine ausführliche Diskussion des hier betrachteten Anwendungsbeispiels sei auf die Arbeiten von Puccetti et al. [PRH19], [PKRH20], [KPRH20] verwiesen.

definierte Projektionsmatrix aufgrund der Unabhängigkeit von $y_{r,\text{soll}}$ einmalig im Voraus berechnet und gespeichert werden, um Rechenzeit während des Onlinetrainings und Onlinebetriebs einzusparen. Der Rechenaufwand für die Approximation des Sollverlaufs nach (3.72) reduziert sich somit in jedem Zeitschritt k zu einer einzelnen Matrixmultiplikation.

Bemerkung 6.2

Ein Regler, der auf einer konstanten Sollgeschwindigkeitsapproximation ($d = 0$) basiert, entspricht dem in [PKRH20] präsentierten Ansatz zur ADP-basierten Geschwindigkeitsregelung ohne Vorsteuerterm. Die Verwendung einer linearen Sollgeschwindigkeitsapproximation ($d = 1$) hat hingegen zur Folge, dass dem Regler zusätzlich zur aktuellen Sollgeschwindigkeit die aktuelle Sollbeschleunigung, d. h. die erste Ableitung des Sollgeschwindigkeitsprofils im Zeitschritt k , zur Verfügung gestellt wird. Die Projektion in (6.8) übernimmt dabei unter anderem die Aufgabe einer numerischen Differenziation¹²³.

6.1.2.2 Netzwerkarchitektur der Q-Function-Approximation

Dem ADP-Regler werden, neben y_k und z_k , vergangene Stellgrößen übergeben, aus denen fehlende Zustandsinformation rekonstruiert werden soll. Diese Idee folgt damit dem Vorgehen von Puccetti et al. [PKRH20], [PRH19]. Konkret wird eine Approximation

$$\bar{x}_k := \mathbf{w}_{\text{FIR}}^T \mathbf{u}_{k-h_{\text{FIR}}:k-1} \quad (6.10)$$

aus h_{FIR} vergangenen Stellgrößen

$$\mathbf{u}_{k-h_{\text{FIR}}:k-1} := [u_{k-h_{\text{FIR}}} \quad u_{k-h_{\text{FIR}}+1} \quad \cdots \quad u_{k-1}]^T \quad (6.11)$$

und den zu schätzenden Gewichten $\mathbf{w}_{\text{FIR}} \in \mathbb{R}^{h_{\text{FIR}}}$ verwendet. Dies entspricht einem FIR-Filter. Eine geeignete Länge h_{FIR} dieses Filters ergibt sich grundsätzlich anhand der Einschwingzeit der Impulsantwort der Aktuatordynamik (vgl. [PRH19]). Basierend auf Erfahrungswerten hat sich die Rekonstruktion eines einzelnen Zustands, wie in (6.10) gezeigt, für das hier betrachtete Beispiel der Longitudinaldynamik bewährt [PRH19, Abschnitt VI].

Die verwendete Netzwerkarchitektur, die diesen Rekonstruktionsmechanismus und den durch z_k approximierten Sollgeschwindigkeitsverlauf integriert, ist in Abbildung 6.3 gegeben. Die Eingangsgrößen sind durch die Geschwindigkeit y_k , den Referenzparametervektor z_k , vergangene Stellgrößen $\mathbf{u}_{k-h_{\text{FIR}}:k-1}$ sowie die aktuelle Stellgröße u_k gegeben. Die Schätzung $Q_w^{\mu\theta}(\tilde{x}_k, u_k)$ der Q-Function $Q^\mu(\tilde{x}_k, u_k)$ bildet die Ausgangsgröße. Die Wahl der quadratischen Schicht („qf“) wird in der Arbeit von Puccetti et al. [PRH19] motiviert und folgt aufgrund näherungsweise linearer Dynamik des Systems sowie der quadratischen Form des Gütefunktional nach (6.2). Im nächsten Abschnitt folgt die Vorstellung des verwendeten Reglers sowie dessen Trainingsprozedur.

¹²³ Beispielsweise resultiert für den Fall $d = 1$ und $n_h = 2$ aus (6.9) gerade $\mathbf{P}_{\text{LS}} = \begin{bmatrix} -\Delta t & 1 \\ \Delta t & 0 \end{bmatrix}$ und somit aus (6.8)

$$z_k^T = [(y_{r,\text{soll},k+1} - y_{r,\text{soll},k}) \Delta t \quad y_{r,\text{soll},k}].$$

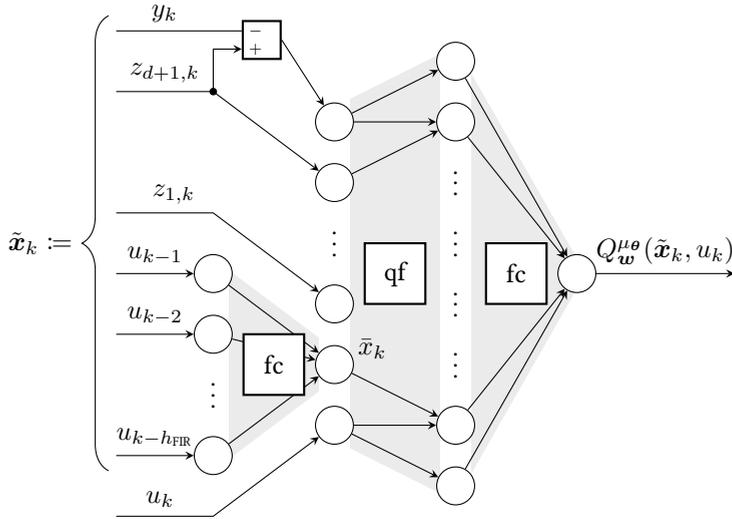


Abbildung 6.3: Verwendete Netzwerkarchitektur für die Approximation der Q-Funktion. Hierbei ist d der Grad des Polynoms der Sollgeschwindigkeitsapproximation mithilfe der Basisfunktionen $\rho(\kappa)$ und des Gewichts z_k . In der ersten Schicht werden einerseits der Geschwindigkeitsfehler $y_r(z_k) - y_k = z_{d+1,k} - y_k$, der Parameter z_k und die aktuelle Stellgröße u_k direkt übergeben. Andererseits wird mithilfe einer vollständig verbundenen Schicht ((engl.): *fully connected layer*, gekennzeichnet durch ‚fc‘) gemäß (6.10) aus h_{FIR} vergangenen Stellgrößen eine Schätzung des Hilfszustands \tilde{x}_k durchgeführt. Die zweite Schicht (gekennzeichnet durch ‚qf‘) berechnet jede multiplikative Kombination aus der vorherigen Schicht und entspricht somit quadratischen Basisfunktionen. Die dritte Schicht ist wiederum eine vollständig verbundene Schicht, welche die geschätzte Q-Funktion ausgibt.

6.1.2.3 Gesamtalgorithmus des ADP-basierten Geschwindigkeitsreglers

Der ADP-basierte Regelungsalgorithmus gliedert sich in die drei Teile *Regelung*, *Datenaufbereitung* und *Training*, die nachfolgend erläutert werden. Dabei erfolgt die *Regelung* mit einer durch Δt gegebenen Abtastzeit. Demgegenüber ist die Updaterate des *Trainings* durch Δl gegeben. Um die Echtzeitanforderung der *Regelung* zu gewährleisten, kann Δl von Δt abweichen. Eine grafische Übersicht über den Gesamtalgorithmus ist in Abbildung 6.4 gegeben. Eine kurze Einführung in während des Trainingsvorgangs verwendete Actor-Critic-Mechanismen ist zudem in Anhang D.2 gegeben.

Regelung

Der Block *Regelung* bildet die Schnittstelle zum realen System. Die Eingangsgrößen sind durch die gemessene Geschwindigkeit y_k und einen von außen vorgegebenen Sollgeschwindigkeitsverlauf $y_{r,\text{soll},k}, \dots, y_{r,\text{soll},k+n_h-1}$ definiert, wohingegen die Stellgröße u_k die Ausgabe darstellt. Innerhalb dieses Blocks wird einerseits der erweiterte Zustand \tilde{x}_k gebildet und andererseits die Stellgröße u_k berechnet.

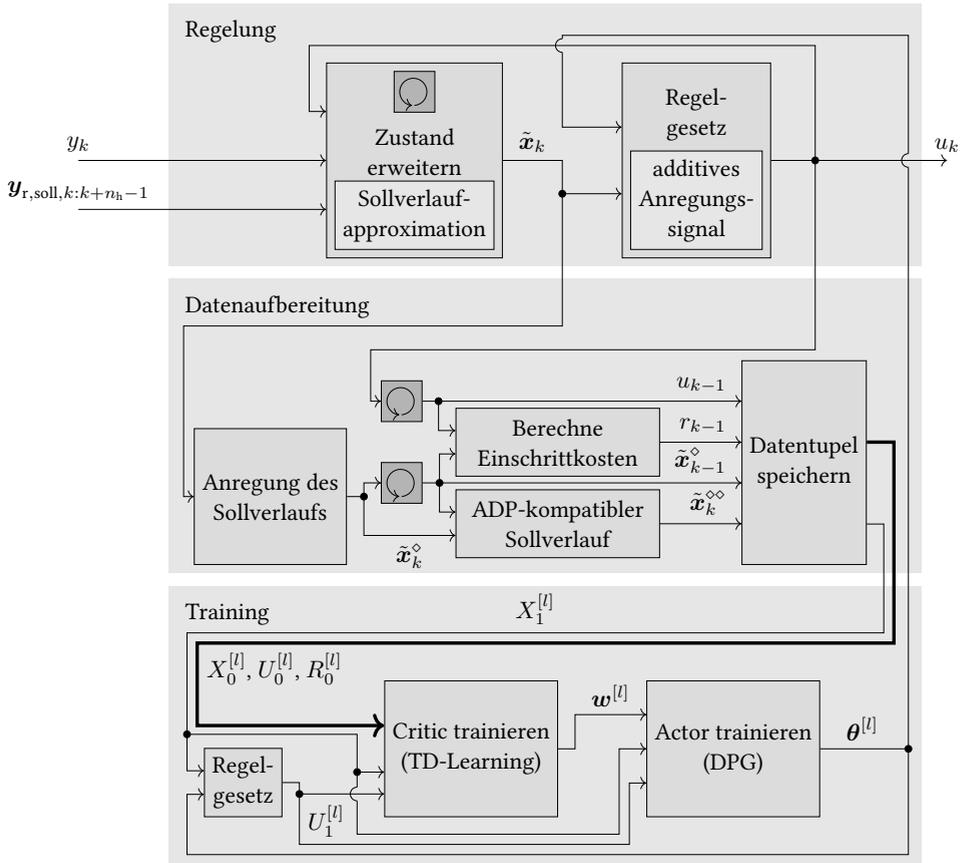


Abbildung 6.4: Übersicht über den Gesamtalgorithmus des ADP-basierten Geschwindigkeitsreglers. Die gezeigte Gesamtstruktur entspricht dem Block *ADP-Regler* in Abbildung 6.2. Der Block *Regelung* übergibt mit der Abtastzeit Δt die Stellgröße u_k an das reale System. Die *Datenaufbereitung* übernimmt Datenvorverarbeitungsschritte, um ADP-kompatible Datentupel zu erzeugen und zu speichern. Im Block *Training* werden diese Datentupel schließlich zum Critic- und Actor-Training genutzt. Dieser Trainingsschritt findet mit der Updaterate Δl statt.

Für die Zustandserweiterung werden zunächst mithilfe eines Ringspeichers Stellgrößen gespeichert, um den Vektor $u_{k-h_{FIR}:k-1}$ vergangener Stellgrößen zu erhalten. Weiterhin wird der Referenzparameter z_k , wie in (6.8) beschrieben, berechnet. Der erweiterte Zustand \tilde{x}_k wird schließlich aus der gemessenen Geschwindigkeit y_k , den Parametern z_k und den vergangenen Stellgrößen $u_{k-h_{FIR}:k-1}$ gebildet (vgl. Abbildung 6.3).

Basierend auf dem aktuellen Reglergewicht $\theta^{[l]}$ und dem erweiterten Zustand \tilde{x}_k wird anschließend mithilfe des Regelgesetzes $\mu_{\theta^{[l]}}(\tilde{x}_k)$ die Stellgröße berechnet. Dabei ist $\mu_{\theta^{[l]}}(\tilde{x}_k)$

durch eine lineare Abhangigkeit des Regelfehlers $y_r(\mathbf{z}_k) - y_k = z_{d+1,k} - y_k$ und, im Fall $d > 0$, die ubrigen Solltrajektorienparameter $z_{d,k}, \dots, z_{1,k}$ gegeben:

$$\mu_{\theta^{[l]} }(\tilde{\mathbf{x}}_k) = [(z_{d+1,k} - y_k) \quad z_{d,k} \quad \dots \quad z_{1,k}]^\top \theta^{[l]}. \quad (6.12)$$

Wahrend des Trainingsvorgangs wird zudem ein Anregungssignal zu $\mu_{\theta^{[l]} }(\tilde{\mathbf{x}}_k)$ addiert, um die Stellgroe u_k zu erhalten. Dieses wird in Abschnitt 6.1.3.1 konkretisiert. Letztlich wird der erweiterte Zustand $\tilde{\mathbf{x}}_k$ sowie die Stellgroe u_k an den Block *Datenaufbereitung* ubergeben.

Datenaufbereitung

Der Block *Datenaufbereitung* ubernimmt Datenvorverarbeitungsschritte, um ADP-kompatible Datentupel zu erzeugen, zu speichern, und dem nachfolgenden *Training* zur Verfugung zu stellen. Insbesondere werden hier auch die Referenzparameter \mathbf{z}_k modifiziert. Diese Modifikation verfolgt zwei Ziele. Einerseits werden die Referenzparameter angeregt, andererseits wird dafur gesorgt, dass die beim Training verwendeten Datentupel $\{X_0^{[l]}, U_0^{[l]}, R_0^{[l]}, X_1^{[l]}\}$ ADP-kompatibel sind.

Ersteres wird durch Addition von Zufallszahlen zu den Polynomkoeffizienten \mathbf{z}_k erreicht. Dieses Vorgehen ist aufgrund der Off-Policy-Charakteristik zulassig (vgl. Anhang D.2). Im Folgenden wird hierzu eine mittelwertfreie Gauverteilung mit Standardabweichung 1 verwendet. In Abbildung 6.4 kennzeichnet $\tilde{\mathbf{x}}_k^\diamond$, dass dieses zusatzliche Anregungssignal zum Referenzparameter \mathbf{z}_k addiert wird, der Teil des erweiterten Zustands $\tilde{\mathbf{x}}_k$ ist.

Ein Pufferspeicher verzogert den erweiterten Zustand $\tilde{\mathbf{x}}_k^\diamond$, sodass $\tilde{\mathbf{x}}_{k-1}^\diamond$ zur Verfugung steht. Die Groe $\tilde{\mathbf{x}}_k^\diamond$ stellt schlielich ADP-Kompatibilitat sicher. $\tilde{\mathbf{x}}_k^\diamond$ entsteht, indem in $\tilde{\mathbf{x}}_k$ der Solltrajektorienparameter \mathbf{z}_k durch $\mathbf{z}_{k-1}^{(1)}$ ersetzt wird, also durch eine um einen Zeitschritt verschobene Version des durch \mathbf{z}_{k-1} beschriebenen Solltrajektorienverlaufs aus $\tilde{\mathbf{x}}_{k-1}^\diamond$ (vgl. (6.3), (6.5) und Bemerkung 3.2). Mithilfe der vergangenen Stellgroe u_{k-1} und $\tilde{\mathbf{x}}_{k-1}^\diamond$ werden zudem die Einschrittkosten r_{k-1} berechnet. Das so entstandene Tupel $\{\tilde{\mathbf{x}}_{k-1}^\diamond, u_{k-1}, r_{k-1}, \tilde{\mathbf{x}}_k^\diamond\}$ wird in einem Ringspeicher gespeichert, der M Datentupel fasst. Sobald dieser Ringspeicher gefullt ist, konnen Batches $\{X_0^{[l]}, U_0^{[l]}, R_0^{[l]}, X_1^{[l]}\}$, die aus M_B Datentupeln bestehen, aus diesem Speicher gezogen und dem nachfolgenden Trainingsalgorithmus ubergeben werden.

Training

In jedem Trainingsschritt wird zunachst das aktuelle Regelgesetz $\mu_{\theta^{[l]}}$ (vgl. (6.12)) fur jedes $\tilde{\mathbf{x}}$ in $X_1^{[l]}$ ausgewertet, um $U_1^{[l]}$ zu berechnen. Damit wird der Trainingsbatch schlielich zu $\{X_0^{[l]}, U_0^{[l]}, R_0^{[l]}, X_1^{[l]}, U_1^{[l]}\}$ vervollstandigt. Dieser Batch¹²⁴ wird dann genutzt, um das Critic-Gewicht w basierend auf dem quadrierten TD-Fehler (D.2) anzupassen. Hierbei wird

¹²⁴ Dieses Vorgehen ist unter dem Begriff *Experience Replay* bekannt (vgl. Anhang D.2).

der Levenberg-Marquardt-Algorithmus [Mor78] mit Norm-Clipping [PMB13] verwendet¹²⁵. Anschließend wird das geschätzte Critic-Gewicht $w^{[l]}$ genutzt, um mithilfe des Deterministic Policy Gradients [SLH⁺14] (vgl. Anhang D.2) die Reglergewichte $\theta^{[l]}$ zu adaptieren¹²⁶. Anzumerken sei an dieser Stelle, dass der Trainingsschritt nicht direkt von externen Signalen abhängt. Die Updaterate Δl , mit der das Training durchgeführt wird, kann somit unabhängig von der Abtastzeit Δt , mit der die Messung der Geschwindigkeit y_k erfolgt und Stellgrößen u_k auf das System appliziert werden, gewählt sein. Eine geeignete Wahl von Δl hängt hierbei insbesondere von der verfügbaren Rechenkapazität ab.

6.1.3 ADP-Solltrajektorienfolgeregelung im Realfahrzeug

In diesem Abschnitt wird gezeigt, dass der vorgestellte ADP-Sollgeschwindigkeitsfolgeregler online in einem Realfahrzeug lernt, vorgegebenen Geschwindigkeitsprofilen zu folgen. Zunächst wird der experimentelle Aufbau beschrieben. Anschließend wird der Online-Trainingsvorgang im Realfahrzeug betrachtet und anhand von Validierungsfahrten ausgewertet. Dabei wird der vorgestellte ADP-Regler, der den Sollgeschwindigkeitsverlauf durch ein Polynom ersten Grades approximiert, mit einem ADP-Regler verglichen, der eine konstante Sollgeschwindigkeit zugrunde legt.

6.1.3.1 Online-Training im Realfahrzeug

Für das Online-Training wird ein auf einem *BMW 740Li* basierendes Versuchsfahrzeug mit Verbrennungsmotor und Automatikgetriebe genutzt (siehe Abbildung 6.1). Dieses Fahrzeug ist mit einer *dSpace Autobox*¹²⁷ mit *DS1007* Prozessorboard¹²⁸ ausgestattet, die über das fahrzeuginterne Bussystem Zugriff auf das Bremssystem, den Antriebsstrang und weitere Regler hat. Diese Plattform kann mehrere Aufgaben in Echtzeit ausführen und stellt unter anderem die Schnittstelle zu dem vorgestellten ADP-basierten Geschwindigkeitsregler dar. So werden Geschwindigkeitsmessungen y_k zur Verfügung gestellt, der lokale Sollgeschwindigkeitsverlauf $y_{r,soll,k}, y_{r,soll,k+1}, \dots, y_{r,soll,k+n_h-1}$ verarbeitet und die resultierende Stellgröße u_k der

¹²⁵ Diese Wahl des Optimierungsalgorithmus hat sich als geeignet erwiesen, auch, da vergleichsweise wenige Hyperparameter angepasst werden müssen [PRH19]. Der Levenberg-Marquardt-Algorithmus stellt eine Trust-Region-Optimierungsmethode dar, die neben derselben Approximation der Hessematrix wie bei der Gauß-Newton-Methode einen zusätzlichen Regularisierungsterm nutzt [NW06, S. 258]. Der Regularisierungsterm wird zu $\lambda_{LM}I$ gesetzt und mittels einfacher Schrittweitensteuerung angepasst. Bei einer Verbesserung des Optimierungsziels wird $\lambda_{LM} = \frac{\lambda_{LM}}{3}$ gesetzt, andernfalls $\lambda_{LM} = 5\lambda_{LM}$. Initial gilt $\lambda_{LM} = 1$. Die Verwendung von Norm Clipping, d. h. die Normierung des Gradienten, sobald dessen Norm einen nutzerdefinierten Schwellwert überschreitet, stellt insbesondere eine Begrenzung der Lernrate des Actor-Gewichts sicher. Da die Adaption des Actor-Gewichts auf einer Schätzung des Critic-Gewichts basiert, sollte die Anpassung des Actor-Gewichts langsamer als die Anpassung des Critic-Gewichts stattfinden.

¹²⁶ Genau wie bei der Anpassung des Critic-Gewichts wird auch hier der Levenberg-Marquardt-Algorithmus [Mor78] mit Norm-Clipping [PMB13] verwendet.

¹²⁷ Datenblatt verfügbar unter: https://www.dspace.com/shared/data/pdf/2019/dSPACE_AutoBox_PHS_Catalog2019.pdf (Zugriff am 19.10.2021).

¹²⁸ Datenblatt verfügbar unter: https://www.dspace.com/shared/data/pdf/2019/dSPACE_DS1007_Catalog2019.pdf (Zugriff am 19.10.2021).

unterlagerten Regelung übergeben. Die Datenaufzeichnung erfolgt mittels eines Ethernet-Anschlusses durch einen PC.

Die verwendeten Hyperparameter sind Tabelle D.1 in Anhang D.3 zu entnehmen. Die Wahl einer kleinen maximalen Norm für den Adaptionsschritt des Actor-Gewichts $\theta^{[l]}$ führt zwar zu einer vergleichsweise langsamen Adaption, soll jedoch Rauschen im Adaptionprozess des Actors reduzieren. Dies ist darin begründet, dass die Adaption von $\theta^{[l]}$ gemäß (D.5) von der erfolgreichen Schätzung des Critic-Gewichts w abhängt. Aufgrund der begrenzten Rechenkapazität der verwendeten Hardware wird zudem eine Updaterate des Trainingsvorgangs von $\Delta l = 0,6 \text{ s}$ gewählt (vgl. Tabelle D.1), während die Abtastzeit des Reglers mit $\Delta t = 0,02 \text{ s}$ der Abtastzeit der Systemausgangsgröße y_k entspricht und somit die gegebenen Echtzeitanforderungen erfüllt.

In beiden Trainingsdurchgängen (Ordnung $d = 1$ bzw. $d = 0$ für die Approximation des Sollgeschwindigkeitsprofils) wird das Fahrzeug im ersten Gang betrieben. Die Longitudinaldynamik, die aus Antriebsstrang, Bremssystem und unterlagerter Regelung besteht, dämpft hohe Frequenzen von u_k und weist daher Tiefpassverhalten auf. Der ADP-Mechanismus, der zum Training des Geschwindigkeitsreglers genutzt wird, erfordert jedoch, wie alle adaptiven Methoden, eine ausreichende Systemanregung, um eine Verbesserung der Critic- und Actor-Gewichte zu erzielen. Da mittelwertfreies, hochfrequentes Rauschen aufgrund des Tiefpassverhaltens, ähnlich wie im Beispiel in Abbildung 5.9 gezeigt, kaum Einfluss auf das System hätte, ist es als Anregungssignal für die Stellgröße ungeeignet. Die Ergebnisse aus Kapitel 5 motivieren daher die Verwendung eines Anregungssignals, das zu einer deutlicheren Beeinflussung des Systems führt. Im konkreten Fall des vorliegenden ADP-basierten Geschwindigkeitsreglers hat sich die Addition eines Zufallssignals zum Reglerausgang, der aus einem initial gegebenen, suboptimalen Regler resultiert, bewährt. Dieses Zufallssignal wird dabei alle 2 s zufällig aus einer Gleichverteilung im Intervall $[-1 \text{ m s}^{-2}, 1 \text{ m s}^{-2}]$ gezogen. Während des Trainingsvorgangs wird die Geschwindigkeit zwischen ca. 6 km h^{-1} und 30 km h^{-1} variiert. Ein beispielhafter Ausschnitt der Geschwindigkeit y_k und der Stellgröße u_k während des Trainingsvorgangs ist Abbildung 6.5 zu entnehmen.

Um in den Trainingstupeln $\{X_0^{[l]}, U_0^{[l]}, R_0^{[l]}, X_1^{[l]}, U_1^{[l]}\}$ auch eine Anregung der Referenztrajektorienparameter z_k zu erreichen, werden, wie in Abschnitt 6.1.2.3 beschrieben, verrauschte Parameter z_{k-1} verwendet. Die Verteilung der während des Trainingsvorgangs genutzten verrauschten Referenzparameter z_{k-1} , aus denen mithilfe von $D(1)$ die ADP-kompatiblen Parameter $z_{k-1}^{(1)}$ erzeugt werden, ist für $d = 1$ in Abbildung 6.6 gezeigt¹²⁹.

Abbildung 6.7 zeigt den Verlauf der Reglergewichte $\theta^{[l]}$ für die beiden Durchgänge mit linearer Referenzapproximation (Ordnung $d = 1$) und konstanter Referenzapproximation (Ordnung $d = 0$) während des Trainingsvorgangs. Die Verstärkung, die den konstanten Anteil der Sollgeschwindigkeitsapproximation beschreibt, verhält sich dabei für beide Fälle ähnlich. Nach etwa 1300 s verbleiben die Reglerparameter in einem kleinen Intervall.

¹²⁹ Für den Fall $d = 0$ entfällt der durch z_1 gegebene lineare Anteil.

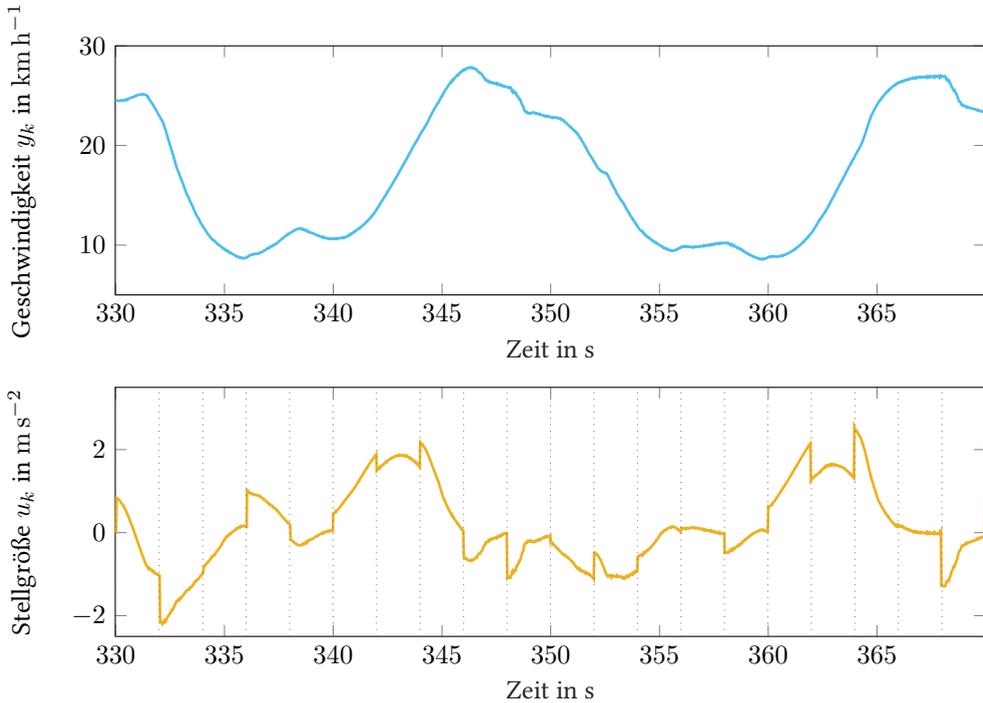


Abbildung 6.5: Beispielhafter Ausschnitt der Geschwindigkeit y_k sowie der Stellgröße u_k während des Trainingsvorgangs. Im Sinne der Systemanregung wird alle 2 s ein Zufallswert aus einer Gleichverteilung auf dem Intervall $[-1 \text{ m s}^{-1}, 1 \text{ m s}^{-1}]$ gezogen, der zum Reglerausgang addiert wird. Die Zeitpunkte, zu denen dieser Zufallswert gezogen wird, sind durch vertikale, gepunktete Linien gekennzeichnet.

6.1.3.2 Auswertung der gelernten Regler im Realfahrzeug

Abschließend sollen die beiden gelernten Regler anhand einer Validierungsfahrt verglichen werden. Dazu werden die Reglergewichte $\theta^{[l]}$ am Ende der beiden Trainingsdurchgänge konstant gehalten. Zudem wird das Anregungssignal entfernt und eine Auswertungsfahrt durchgeführt.

In Abbildung 6.8 ist ersichtlich, dass mit beiden Reglern die gewünschte Sollgeschwindigkeit erreicht wird. Aufgrund des zu minimierenden Gütefunktional (6.2), das unter anderem den Stellaufwand bestraft¹³⁰, sowie der begrenzten Vorausschaufähigkeit der Regler, wird dem vorgegebenen Geschwindigkeitsprofil zwar nicht exakt gefolgt, aufgrund der prädiktiven Eigenschaften des Reglers mit linearer Sollgeschwindigkeitsapproximation ist für $d = 1$ der zeitliche Versatz von y_k bezüglich der Sollgeschwindigkeit im Vergleich zum Regler mit $d = 0$ jedoch merklich reduziert.

¹³⁰ Hierdurch werden insbesondere im Sinne des Fahrkomforts hohe Beschleunigungen vermieden (vgl. [DLL18]).

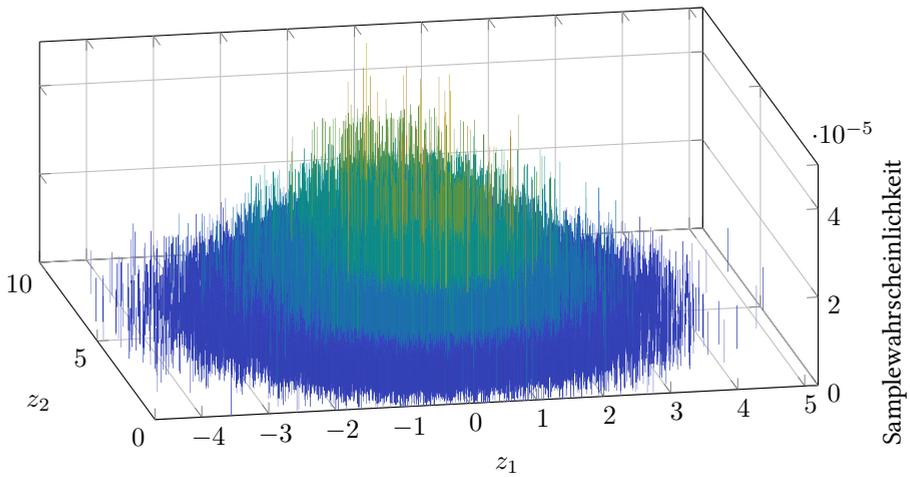


Abbildung 6.6: Verteilung der während des Trainingsvorgangs verwendeten verrauschten Referenzparameter z_{k-1} . Dabei stellt z_2 die Sollgeschwindigkeit (konstanter Anteil in m s^{-1}) und z_1 die Steigung der Sollgeschwindigkeit (linearer Anteil in m s^{-2}) dar.

— konstanter Anteil Ordnung 1 - - - linearer Anteil Ordnung 1
 - - - konstanter Anteil Ordnung 0

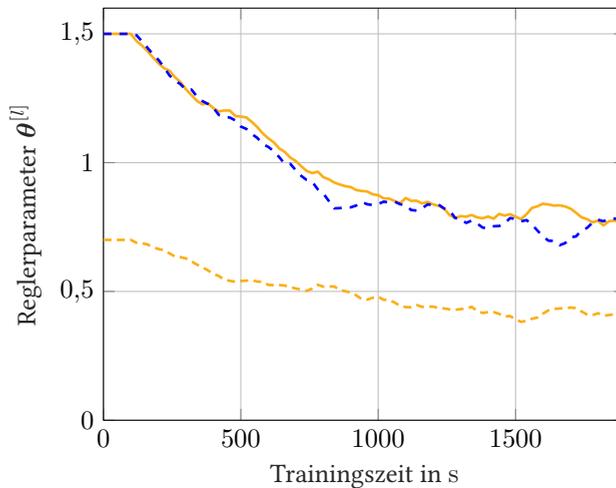


Abbildung 6.7: Reglerparameter $\theta^{[l]}$ während des Trainingsvorgangs. Während der ersten 100 s wird zunächst der Speicher der Größe M mit Datentupeln gefüllt.

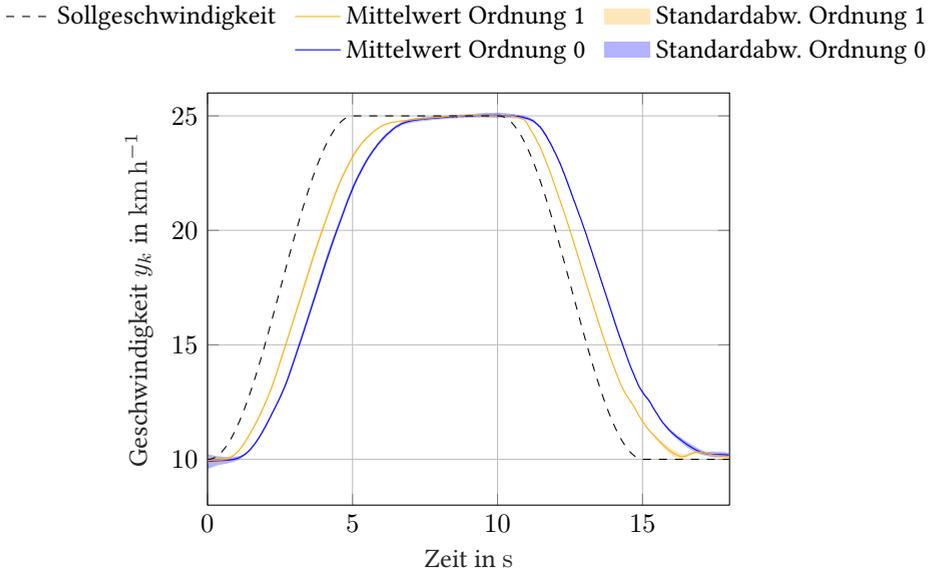


Abbildung 6.8: Sollgeschwindigkeitsverlauf und gemessene Geschwindigkeit y_k während der Auswertungsfahrt nach dem Training. Als beispielhaftes Sollgeschwindigkeitsprofil wurde ein um $17,5 \text{ km h}^{-1}$ verschobener Sinusverlauf mit einer Amplitude von $7,5 \text{ km h}^{-1}$ und einer Periodendauer von 10 s, der an den Extrempunkten für jeweils 5 s konstant gehalten wird, gewählt. Die Ergebnisse sind über 20 Durchgänge gemittelt. Der Regler, der eine lineare Approximation des Sollgeschwindigkeitsverlaufs berücksichtigt ($d = 1$), weist dabei eine geringere Abweichung vom Sollverlauf auf, als der Vergleichsregler mit konstanter Geschwindigkeit ($d = 0$).

Die aktuelle Sollgeschwindigkeitsabweichung, die durch den konstanten Anteil des Polynoms nullten Grades ($d = 0$) bzw. des Polynoms ersten Grades ($d = 1$) berücksichtigt wird, hat bei beiden Reglern einen nahezu identischen Einfluss. Dies ist den konstanten Anteilen der Solltrajektorienparameter z_k in Abbildung 6.9 sowie den konstanten Anteilen der Reglerparameter $\theta^{[1]}$ in Abbildung 6.7 zu entnehmen. Der Regler mit linearer Approximation des Sollgeschwindigkeitsprofils kann jedoch vorausschauender handeln, da der lineare Anteil der Referenzapproximation als Vorsteuerterm agiert (vgl. (6.12) und [Lun20a, S. 11]). Diese zusätzliche Stellgröße hängt dabei direkt mit der aktuellen Steigung der Sollgeschwindigkeit zusammen, d. h. der aktuellen Sollbeschleunigung, die dem linearen Anteil des Solltrajektorienparameters z_k in Abbildung 6.9 entspricht. Aufgrund dieses zusätzlichen Vorsteuerterms sendet der Regler, der eine lineare Approximation des Sollverlaufs nutzt, die durch u_k gegebenen Beschleunigungssignale früher als der Vergleichsregler. Dies ist in Abbildung 6.10 gezeigt.

Abschließend werden noch die mit den beiden Reglern während der Auswertungsfahrt resultierenden Kosten im Sinne des Gütefunktional (6.2) betrachtet. Da dieses Gütemaß im

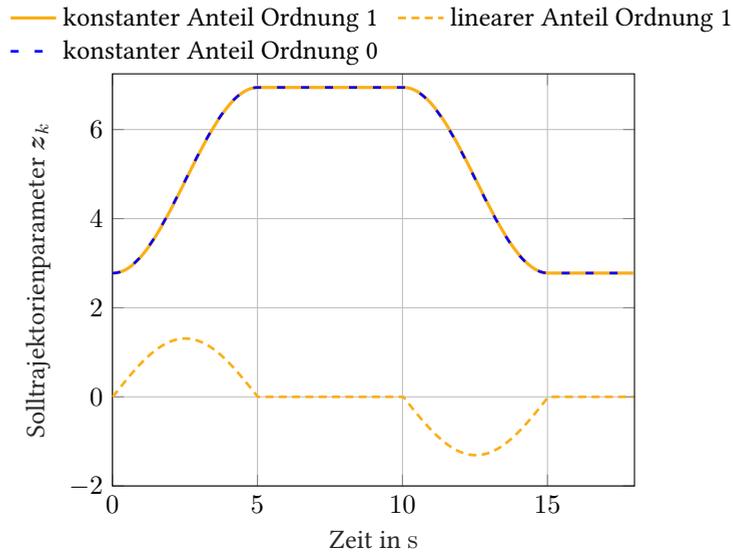


Abbildung 6.9: Solltrajektorienparameter z_k während der Auswertungsfahrt. Der konstante Anteil repräsentiert die aktuelle Sollgeschwindigkeit, der lineare Anteil die aktuelle Sollbeschleunigung.

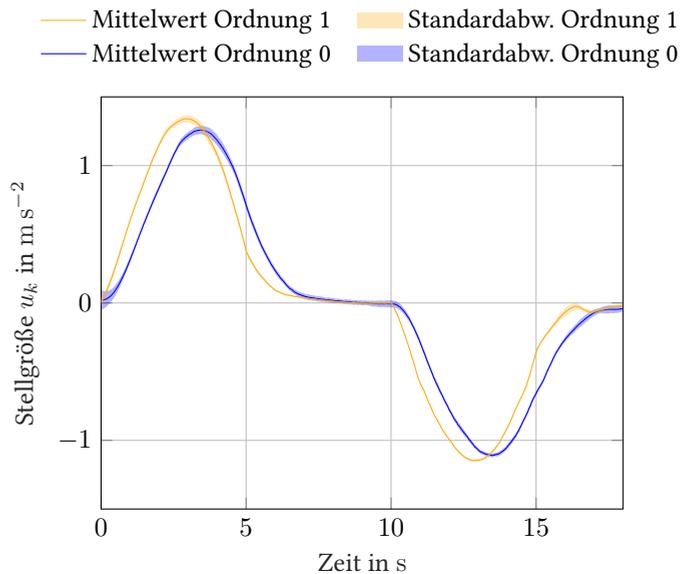


Abbildung 6.10: Stellgröße u_k während des Auswertungsmanövers (gemittelt über 20 Durchgänge). Die Amplituden der Stellgröße sind für die betrachteten Fälle vergleichbar, jedoch reagiert der Regler mit $d = 1$ früher und handelt vorausschauender.

Realversuch aufgrund des unendlichen Zeithorizonts in (6.2) nicht berechnet werden kann, wird als Näherung

$$V^\mu(y_k, z_k) \approx \sum_{\kappa=0}^{100} \gamma^\kappa r\left(y_{k+\kappa}, y_r(z_k^{(\kappa)}), \mu(y_{k+\kappa}, z_k^{(\kappa)})\right) \quad (6.13)$$

genutzt. Diese Näherung, die aufgrund der verwendeten Diskontierung $\gamma = 0,95$ (vgl. Tabelle D.1) zulässig ist, da Kosten in ferner Zukunft nur wenig zur Value Function V^μ beitragen, gibt einen Eindruck über die Güte der beiden Regler. Wie in Abbildung 6.11 ersichtlich ist, reduziert der vorgestellte Regler, der eine lineare Approximation des Sollgeschwindigkeitsverlaufs verwendet, die Kosten (6.13) deutlich gegenüber dem Vergleichsregler. Die besseren Vorausschaugeigenschaften und der daraus resultierende verringerte zeitliche Versatz führen zu geringeren Regelabweichungen und damit letztlich zu geringeren Kosten.

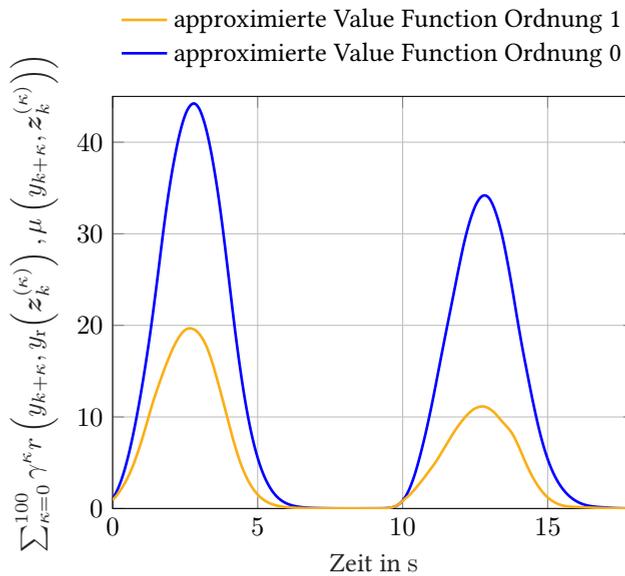


Abbildung 6.11: Approximierte Value Function während des Auswertungsmanövers. Der Regler, der eine lineare Approximation des Sollverlaufs verwendet ($d = 1$), führt dabei zu deutlich geringeren Kosten als der Vergleichsregler.

6.1.4 Diskussion

In diesem Abschnitt wurde, basierend auf der in Abschnitt 3.2 vorgestellten ADP-kompatiblen Referenztrajektorienarstellung, ein lernender Geschwindigkeitsregler für ein Realfahrzeug vorgestellt, der online, d. h. während der Trainingsfahrt, die Schätzparameter der Q-Function und die Reglerparameter adaptiert. Dieser Ansatz zeichnet sich durch die explizite Einbeziehung der lokalen linearen Approximation des Sollgeschwindigkeitsverlaufs aus. Die Konsistenz

der Trainingsdaten mit der Bellman-Gleichung wird hierbei durch die Definition eines erweiterten Zustands \tilde{x}_k sowie der virtuellen Verschiebung der Solltrajektorienparameter z_k erreicht.

Das erfolgreiche Online-Training in einem Realfahrzeug demonstriert die grundsätzliche Anwendbarkeit der in dieser Arbeit entwickelten, neuartigen, ADP-basierten Solltrajektorienregler in realen regelungstechnischen Problemstellungen und liefert somit einen anwendungsorientierten Beitrag zu Forschungsfrage 1 aus Abschnitt 2.4.1. Insbesondere offenbart die explizite Einbeziehung des approximierten Solltrajektorienverlaufs in die Value Function Vorteile. So führt die zusätzliche Verwendung der Änderungsrate der aktuellen Sollgeschwindigkeit zu besseren Vorausschauereigenschaften und signifikant reduzierten Kosten.

Eine effizientere Implementierung der verwendeten Algorithmen sowie ein auf einer externen Recheneinheit durchgeführter Trainingsprozess können künftig die Updaterate Δl des Trainings wesentlich reduzieren und ein schnelleres Training ermöglichen. Ein schnellerer Trainingsvorgang erleichtert schließlich mögliche Erweiterungen, zum Beispiel das Training in anderen Geschwindigkeitsbereichen und Gängen oder die Einbeziehung zusätzlicher Daten, wie beispielsweise das aktuelle Straßenprofil.

Im nächsten Abschnitt wird die in dieser Arbeit vorgestellte ADP-kompatible Referenztrajektorienendarstellung auf ein zweites reales regelungstechnisches Problem angewandt und ein daraus resultierender modellfreier ADP-Regler mit einem modellbasierten Regler verglichen.

6.2 Modellfreier Trajektorienfolgeregler für ein reales Ball-auf-Platte-System

In diesem Abschnitt wird die in Kapitel 3 vorgestellte ADP-kompatible parametrisierte Referenztrajektorienendarstellung verwendet, um einen Solltrajektorienfolgeregler für ein reales Ball-auf-Platte-System, das in Abbildung 6.12 gezeigt ist, zu trainieren¹³¹. Obwohl Ball-auf-Platte-Systeme weitverbreitete Beispielsysteme in der Regelungstechnik darstellen, sind die in der Literatur vorhandenen Regler hierfür entweder vollständig modellbasiert [KIB⁺19], [ABB⁺02], [BFG12], [DHS17], [KCS03] oder modellbasiert mit zusätzlicher Fuzzy-Regelung [MSV08]. Im Gegensatz zu bestehenden Reglern benötigt der im Folgenden verwendete Regelungsansatz kein exaktes Modell des Ball-auf-Platte-Systems, da Messdaten des realen Systems dazu verwendet werden, einen optimierungsbasierten Solltrajektorienfolgeregler zu trainieren. Dieses Vorgehen vermeidet daher eine Modellbildung, die beispielsweise durch konstruktive Eigenschaften, wie einer vertikal versetzten Drehachse der Platte, erschwert wird (vgl. [KIB⁺19]). Ebenso wird eine anschließende manuelle Feinjustage nicht mehr benötigt. Somit stellt die vorliegende Arbeit die erste Anwendung eines modellfreien ADP-basierten Solltrajektorienfolgeregelungsansatzes auf ein reales Ball-auf-Platte-System dar.

¹³¹ Ergebnisse dieses Abschnitts wurden im Rahmen eines Konferenzbeitrags veröffentlicht [KKIH21].



Abbildung 6.12: Ball-auf-Platte-System des Instituts für Regelungs- und Steuerungssysteme am Karlsruher Institut für Technologie, das für die modellfreie, ADP-basierte Solltrajektorienfolgeregelung verwendet wird.

Weiterhin erlauben existierende Regler für Ball-auf-Platte-Systeme bisher entweder überhaupt keine Solltrajektorienvorgabe der Ballposition [KIB⁺19], [ABB⁺02], [BFGB12] oder betrachten nur die aktuelle Abweichung der Ballposition von einer stationären Sollpositionsvorgabe [DHS17], [KCS03], [MSV08]. Letzteres führt jedoch zu einer merklichen Zeitverzögerung im Vergleich zum gewünschten Trajektorienverlauf. Demgegenüber wird in der vorliegenden Arbeit eine ADP-kompatible lokale Approximation des Solltrajektorienverlaufs der Ballposition explizit in den Regler integriert. Diese flexible, aber dennoch kompakte Darstellung ermöglicht prädiktives Verhalten des resultierenden Reglers.

Die Nutzung eines Off-Policy-Ansatzes (vgl. Abschnitt 2.1.4.4) erlaubt zudem die Wiederverwendung der Messdaten und somit eine dateneffiziente Implementierung. Die Verwendung einer On-Policy-Methode würde nicht nur dazu führen, dass zwingend benötigtes Anregungsrauschen (vgl. Kapitel 5) zu einem Offset in der Schätzung der Critic-Gewichte führen könnte,

sondern insbesondere nach jedem Policy-Improvement-Schritt komplett neue Messdaten vom realen System aufgezeichnet werden müssten. Der vorgestellte ADP-basierte Solltrajektorienfolgeregler mit lokaler Approximation des Sollpositionsverlaufs wird schließlich sowohl mit einem ADP-basierten Regler mit stationärer Sollvorgabe als auch mit modellbasierten Reglern, die ebenfalls eine lokale Approximation des Sollpositionsverlaufs bzw. eine konstante Sollvorgabe verwenden, verglichen.

6.2.1 Ball-auf-Platte-System und Problembeschreibung

Nachfolgend wird zunächst das verwendete reale Ball-auf-Platte-System, das ohne genaue Kenntnis eines Systemmodells mithilfe eines ADP-Solltrajektorienfolgeregelungsansatzes geregelt werden soll, vorgestellt und anschließend die Problemstellung definiert.

6.2.1.1 Ball-auf-Platte-System

Das in dieser Arbeit verwendete und in Abbildung 6.12 gezeigte Ball-auf-Platte-System besitzt eine quadratische Platte mit einer Seitenlänge von 1 m und einer Masse von 16,3 kg. Die Platte kann in zwei zueinander orthogonalen Richtungen, die im Folgenden durch X und Y gekennzeichnet sind, geneigt werden. Für beide Plattendimensionen steht hierfür ein eigener Motor zur Verfügung. Die Plattenwinkel ($\alpha^{[X]}$, $\alpha^{[Y]}$) und zugehörigen Winkelgeschwindigkeiten ($\omega^{[X]}$, $\omega^{[Y]}$) werden alle 10 ms gemessen. Ein Ball mit einer Masse von 0,042 kg und einem Radius von 0,02 m befindet sich auf der Platte. Seine Position in einem plattenfesten Koordinatensystem wird kamerabasiert erfasst, sodass alle $\Delta t = 40$ ms die Ballposition ($s^{[X]}$, $s^{[Y]}$) sowie die Ballgeschwindigkeit ($v^{[X]}$, $v^{[Y]}$) zur Verfügung stehen. Für $\iota \in \mathcal{D} = \{X, Y\}$ werden die Systemzustände der jeweiligen Plattendimension zu

$$\mathbf{x}_k^{[\iota]} = \begin{bmatrix} s_k^{[\iota]} & v_k^{[\iota]} & \alpha_k^{[\iota]} & \omega_k^{[\iota]} \end{bmatrix}^\top \quad (6.14)$$

zusammengefasst. Die Eingangsgrößen des Systems sind durch die Motorströme $u_k^{[\iota]} = I_k^{[\iota]}$ gegeben. Die Systemarchitektur ist in Abbildung 6.13 skizziert. Eine detailliertere Beschreibung der Systemarchitektur und der verwendeten Hardware ist in [KIB⁺19], [Kil20] und [Bla18] zu finden¹³².

Da die beiden Plattendimensionen X und Y sich nur geringfügig gegenseitig beeinflussen, ist eine getrennte Regelung in beiden Dimensionen üblich [ABB⁺02], [BFGB12], [KCS03], [KIB⁺19]. Die Dynamik des Ball-auf-Platte-Systems ist zwar in X - und Y -Richtung unterschiedlich, jedoch werden die Regler für beide Plattendimensionen grundsätzlich mit demselben Verfahren trainiert. Daher wird auf den Index ι im Folgenden aus Gründen der Lesbarkeit zumeist verzichtet.

¹³² Im Gegensatz zu [KIB⁺19] und [Bla18] wird in der vorliegenden Arbeit eine schwerere Platte sowie ein anderer Ball verwendet. Dies wurde insbesondere bei der Auslegung des modellbasierten Vergleichsreglers berücksichtigt.

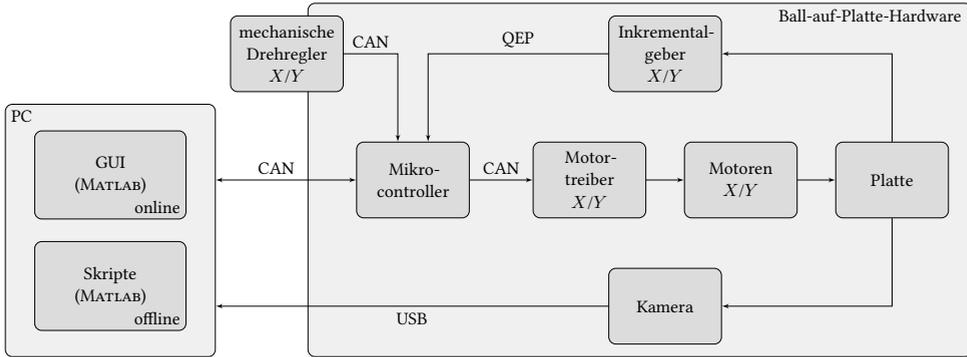


Abbildung 6.13: Architektur des Ball-auf-Platte-Systems. Abbildung nach [Kil20].

6.2.1.2 Problemstellung

Betrachtet werde ein zeitdiskretes, steuerbares System mit der Dynamik

$$\mathbf{x}_{k+1} = \mathbf{F}(\mathbf{x}_k, \mathbf{u}_k), \quad (6.15)$$

wobei $k \in \mathbb{N}_{\geq 0}$ den diskreten Zeitschritt indiziert, $\mathbf{x}_k \in \mathbb{R}^n$ dem Systemzustand entspricht (vgl. (6.14)), $\mathbf{u}_k \in \mathbb{R}^p$ die Stellgröße I_k darstellt, und \mathbf{F} unbekannt ist. Aus Abschnitt 6.2.1.1 folgt für jede Dimension in \mathcal{D} die Systemordnung $n = 4$ sowie eine Stellgrößenanzahl von $p = 1$. In jedem Zeitschritt k sei eine lokale Approximation des Solltrajektorienverlaufs der Ballposition durch

$$s_r\left(\mathbf{z}_k^{(\kappa)}\right) := s_r(\mathbf{z}_k, \kappa) = \mathbf{z}_k^\top \boldsymbol{\rho}(\kappa) \quad (6.16)$$

mit $\kappa \in \mathbb{N}_{\geq 0}$ gegeben, wobei $s_r\left(\mathbf{z}_k^{(\kappa)}\right)$ die Sollposition im Zeitschritt $k + \kappa$ bezeichnet. Der Vektor $\mathbf{z}_k \in \mathbb{R}^{d+1}$ parametrisiert den Sollverlauf und $\boldsymbol{\rho}(\kappa)$ bezeichnet gegebene Basisfunktionen (vgl. (3.15)).

Die folgende Problemstellung formalisiert das Ziel, dass die Ballposition einem vorgegebenen Verlauf folgen soll, während für die Werte der übrigen Systemzustände bei entsprechender Parametrierung eine möglichst geringe Abweichung von null gefordert werden kann. Außerdem soll der Stellaufwand gering gehalten werden.

Problem 6.1

Gegeben seien Basisfunktionen $\rho(\kappa)$ zur ADP-kompatiblen Solltrajektorienapproximation nach Definition 3.1 sowie M Messtupel $\{\mathbf{x}_\kappa, u_\kappa, \mathbf{x}_{\kappa+1}\}$, $\kappa = k, \dots, k + M - 1$. Die Systemdynamik $\mathbf{F}(\mathbf{x}_k, u_k)$ sei unbekannt. Gesucht ist das Regelgesetz $\mu^*(\mathbf{x}_k, \mathbf{z}_k)$, sodass $\forall \mathbf{x}_k, \mathbf{z}_k$ die durch $u_k^* = \mu^*(\mathbf{x}_k, \mathbf{z}_k)$ gegebene Stellgröße das Gütefunktional

$$J_k := \sum_{\kappa=0}^{\infty} \gamma^\kappa \left(\begin{bmatrix} x_{1,k+\kappa} - s_r(\mathbf{z}_k^{(\kappa)}) \\ x_{2,k+\kappa} \\ x_{3,k+\kappa} \\ x_{4,k+\kappa} \end{bmatrix}^\top \mathbf{Q} \begin{bmatrix} x_{1,k+\kappa} - s_r(\mathbf{z}_k^{(\kappa)}) \\ x_{2,k+\kappa} \\ x_{3,k+\kappa} \\ x_{4,k+\kappa} \end{bmatrix} + u_{k+\kappa}^\top R u_{k+\kappa} \right) \\ =: \sum_{\kappa=0}^{\infty} \gamma^\kappa r(\mathbf{x}_{k+\kappa}, u_{k+\kappa}, s_r(\mathbf{z}_k^{(\kappa)})) \quad (6.17)$$

minimiert. Hierbei sei $\gamma \in (0, 1]$ ein Diskontierungsfaktor, \mathbf{Q} positiv semidefinit und R positiv definit.

6.2.2 ADP-Solltrajektorienfolgeregler für ein Ball-auf-Platte-System

Für die optimale Q-Funktion

$$Q^*(\mathbf{x}_k, \mathbf{z}_k, u_k) = r(\mathbf{x}_k, s_r(\mathbf{z}_k), u_k) + \sum_{\kappa=1}^{\infty} \gamma^\kappa r(\mathbf{x}_{k+\kappa}, s_r(\mathbf{z}_k^{(\kappa)}), \mu^*(\mathbf{x}_{k+\kappa}, \mathbf{z}_k^{(\kappa)})) \\ = r(\mathbf{x}_k, s_r(\mathbf{z}_k), u_k) + \gamma Q^*(\mathbf{x}_{k+1}, \mathbf{z}_k^{(1)}, \mu^*(\mathbf{x}_{k+1}, \mathbf{z}_k^{(1)})) \quad (6.18)$$

mit $\mathbf{z}_k^{(\kappa)}$ nach Definition 3.2 folgt gemäß Lemma 3.1, dass

$$u_k^* = \arg \min_{u_k} Q^*(\mathbf{x}_k, \mathbf{z}_k, u_k) \quad (6.19)$$

eine Lösung für Problem 6.1 darstellt. Da die optimale Q-Funktion $Q^*(\mathbf{x}_k, \mathbf{z}_k, u_k)$ a priori unbekannt ist, sei $\hat{Q}^*(\mathbf{x}_k, \mathbf{z}_k, u_k) = \hat{\mathbf{w}}^\top \phi(\mathbf{x}_k, \mathbf{z}_k, u_k)$, mit dem zu schätzenden Gewicht $\hat{\mathbf{w}} \in \mathbb{R}^h$ und $\phi(\cdot) \in \mathbb{R}^h$ (vgl. (3.24)). Mithilfe einer Policy Iteration (vgl. Abschnitt 2.1.4.1) wird in der l -ten Iteration im Policy-Evaluation-Schritt die geschätzte Q-Funktion

$$\hat{Q}^{[l+1]}(\mathbf{x}_k, \mathbf{z}_k, u_k) = \hat{\mathbf{w}}^{[l+1]\top} \phi(\mathbf{x}_k, \mathbf{z}_k, u_k), \quad (6.20)$$

welche die Gesamtkosten des aktuellen Regelgesetzes $\hat{\mu}^{[l]}$ beschreibt, gesucht. Hierzu wird $\hat{\mathbf{w}}^{[l+1]}$ adaptiert, um nach Möglichkeit die Gleichung

$$\hat{Q}^{[l+1]}(\mathbf{x}_k, \mathbf{z}_k, u_k) = r(\mathbf{x}_k, s_r(\mathbf{z}_k), u_k) + \gamma \hat{Q}^{[l+1]}(\mathbf{x}_{k+1}, \mathbf{z}_k^{(1)}, \hat{\mu}^{[l]}(\mathbf{x}_{k+1}, \mathbf{z}_k^{(1)})) \quad (6.21)$$

zu erfüllen. Der anschließende Policy-Improvement-Schritt ist dann durch

$$\hat{\mu}^{[l+1]}(\mathbf{x}_k, \mathbf{z}_k) = \arg \min_{u_k} \hat{Q}^{[l+1]}(\mathbf{x}_k, \mathbf{z}_k, u_k) \quad (6.22)$$

gegeben. In den nächsten Abschnitten wird die verwendete lokale Solltrajektorienapproximation, die Approximation der Q-Function und der Trainingsalgorithmus beschrieben.

6.2.2.1 Quadratische lokale Referenzapproximation

Im Folgenden werde der Solltrajektorienverlauf der Ballposition aus lokaler Perspektive zum Zeitschritt k durch ein quadratisches Polynom ($d = 2$)

$$s_r(\mathbf{z}_k^{(\kappa)}) = \mathbf{z}_k^\top \boldsymbol{\rho}(\kappa) = p_{k,2}(\kappa\Delta t)^2 + p_{k,1}\kappa\Delta t + p_{k,0} \quad (6.23)$$

mit $\boldsymbol{\rho}(\kappa) := [(\kappa\Delta t)^2 \quad \kappa\Delta t \quad 1]^\top$ und dem Parametervektor $\mathbf{z}_k := [p_{k,2} \quad p_{k,1} \quad p_{k,0}]^\top$ approximiert. Die Abtastzeit ist durch Δt gegeben. Die Verschiebungsmatrix $\mathbf{D}(\kappa)$, die benötigt wird, um den propagierten Parameter $\mathbf{z}_k^{(\kappa)}$ nach Definition 3.2 aus \mathbf{z}_k zu berechnen, ergibt sich aus

$$\begin{aligned} s_r(\mathbf{z}_k^{(\kappa+j)}) &= s_r(\mathbf{z}_k^{(\kappa)}, j) \\ &= \mathbf{z}_k^\top \boldsymbol{\rho}(\kappa + j) \\ &= \mathbf{z}_k^\top \begin{bmatrix} ((\kappa + j)\Delta t)^2 \\ (\kappa + j)\Delta t \\ 1 \end{bmatrix} \\ &= \mathbf{z}_k^\top \underbrace{\begin{bmatrix} 1 & 2\kappa\Delta t & (\kappa\Delta t)^2 \\ 0 & 1 & \kappa\Delta t \\ 0 & 0 & 1 \end{bmatrix}}_{=: \mathbf{D}(\kappa)} \boldsymbol{\rho}(j) \\ &= \mathbf{z}_k^{(\kappa)\top} \boldsymbol{\rho}(j), \end{aligned} \quad (6.24)$$

$\forall \kappa, j \in \mathbb{N}_{\geq 0}$. Für einen beliebigen Verlauf der Sollposition $s_{r,\text{soll},k}$ ist in jedem Zeitschritt k ein Parametervektor \mathbf{z}_k gesucht, sodass $s_r(\mathbf{z}_k^{(\kappa)})$, $\kappa \in \mathbb{N}_{\geq 0}$, eine Approximation von $s_{r,\text{soll},k+\kappa}$ darstellt. Während der Laufzeit wird dabei angenommen, dass der Sollpositionsverlauf über einen Vorausschauhorizont von $n_h \in \mathbb{N}_{>0}$ Zeitschritten vorliegt. In jedem Zeitschritt k wird \mathbf{z}_k dann mittels gewichteter Least-Squares-Regression geschätzt. Analog zu (3.71) und (3.72) folgt

$$\mathbf{z}_k^\top = s_{r,\text{soll},k:k+n_h-1} \mathbf{W}_p \boldsymbol{\rho}_{0,n_h-1} (\boldsymbol{\rho}_{0,n_h-1}^\top \mathbf{W}_p \boldsymbol{\rho}_{0,n_h-1})^{-1} \quad (6.25)$$

mit

$$s_{r,\text{soll},k:k+n_h-1} := [s_{r,\text{soll},k} \quad s_{r,\text{soll},k+1} \quad \dots \quad s_{r,\text{soll},k+n_h-1}], \quad (6.26)$$

$$\mathbf{W}_p := \text{diag}(1, \beta, \dots, \beta^{n_h-1}), \quad \beta \leq 1, \quad (6.27)$$

$$\boldsymbol{\rho}_{0:n_h-1} := [\boldsymbol{\rho}(0) \quad \boldsymbol{\rho}(1) \quad \dots \quad \boldsymbol{\rho}(n_h-1)]^\top. \quad (6.28)$$

6.2.2.2 Approximation der Q-Function und Policy Iteration

Die Approximation der Q-Function nach (6.20) wird zu

$$\begin{aligned} \hat{Q}^{[l]}(\mathbf{x}_k, \mathbf{z}_k, u_k) &:= \begin{bmatrix} \mathbf{x}_k \\ u_k \\ \mathbf{z}_k \\ 1 \end{bmatrix}^\top \begin{bmatrix} h_{xx}^{[l]} & h_{xu}^{[l]} & h_{xz}^{[l]} & h_{x1}^{[l]} \\ h_{ux}^{[l]} & h_{uu}^{[l]} & h_{uz}^{[l]} & h_{u1}^{[l]} \\ h_{zx}^{[l]} & h_{zu}^{[l]} & h_{zz}^{[l]} & h_{z1}^{[l]} \\ h_{1x}^{[l]} & h_{1u}^{[l]} & h_{1z}^{[l]} & h_{11}^{[l]} \end{bmatrix} \begin{bmatrix} \mathbf{x}_k \\ u_k \\ \mathbf{z}_k \\ 1 \end{bmatrix} \\ &=: \bar{\mathbf{y}}_k^\top \mathbf{H}^{[l]} \bar{\mathbf{y}}_k =: \hat{\mathbf{w}}^{[l]\top} \phi(\mathbf{x}_k, \mathbf{z}_k, u_k) \end{aligned} \quad (6.29)$$

(vgl. (3.46)) mit $\mathbf{H}^{[l]} = \mathbf{H}^{[l]\top}$ gewählt. Somit gilt $\phi(\mathbf{x}_k, \mathbf{z}_k, u_k) = \bar{\mathbf{y}}_k \otimes_{\mathbb{R}} \bar{\mathbf{y}}_k$ und $\hat{\mathbf{w}}^{[l]}$ entspricht den nicht-redundanten Elementen der symmetrischen Matrix $\mathbf{H}^{[l]}$ ¹³³. Für jede Dimension $\iota \in \mathcal{D}$ sind im Fall einer lokalen Referenzapproximation nach Abschnitt 6.2.2.1 mit $d = 2$ somit $h = 45$ Gewichte zu lernen. Die Wahl einer quadratischen Q-Function ist durch die in [KIB⁺19] veröffentlichte Regelung des betrachteten Ball-auf-Platte-Systems mithilfe eines modellbasierten linear-quadratischen Optimalregelungsansatzes motiviert. Nach Lemma 3.2 weist die Q-Function eines LQ-Problems eine quadratische Struktur auf.

Im Policy-Evaluation-Schritt zur Erfüllung von (6.21) wird der Least-Squares-Temporal-Difference-Q-Learning-Ansatz [LP03] mit der in [LP03, Abschnitt 5.2] beschriebenen Fixpunktbedingung verwendet¹³⁴. Somit werden M Tupel $\{\mathbf{x}_k, u_k, \mathbf{x}_{k+1}, \mathbf{z}_k, \mathbf{z}_k^{(1)}\}$ genutzt, um

$$\hat{\mathbf{w}}^{[l+1]} = \left(\bar{\Phi}^\top \left(\bar{\Phi} - \gamma \bar{\Phi}^{+[l]} \right) \right)^{-1} \bar{\Phi}^\top \mathbf{r} \quad (6.30a)$$

zu bestimmen. Dabei gilt

$$\bar{\Phi} := \begin{bmatrix} \phi_1^\top \\ \vdots \\ \phi_M^\top \end{bmatrix}, \quad \bar{\Phi}^{+[l]} := \begin{bmatrix} \phi_1^{+[l]\top} \\ \vdots \\ \phi_M^{+[l]\top} \end{bmatrix}, \quad \mathbf{r} := \begin{bmatrix} r_1 \\ \vdots \\ r_M \end{bmatrix} \quad (6.30b)$$

¹³³ Elemente von $\hat{\mathbf{w}}^{[l]}$, die zu Nebendiagonalelementen von $\mathbf{H}^{[l]}$ gehören, werden dabei jeweils mit dem Faktor 2 multipliziert, um die redundanten Elemente der symmetrischen Matrix $\mathbf{H}^{[l]}$ zu berücksichtigen.

¹³⁴ Der Bellman-Operator $\mathcal{B}Q(\mathbf{x}_k, \mathbf{z}_k, u_k) := r(\mathbf{x}_k, s_r(\mathbf{z}_k), u_k) + \gamma Q(\mathbf{x}_{k+1}, \mathbf{z}_k^{(1)}, \hat{\mu}^{[l]}(\mathbf{x}_{k+1}, \mathbf{z}_k^{(1)}))$ stellt nach [BBdE10, S. 24] eine Kontraktion mit dem Kontraktionsfaktor γ in der L_∞ -Norm dar, deren Fixpunkt durch die korrekte Q-Function zum betrachteten Regelgesetz $\hat{\mu}^{[l]}(\cdot)$ gegeben ist. Da die Anwendung des Bellman-Operators und die Forderung $\hat{Q}^{[l+1]} \stackrel{!}{=} \mathcal{B}\hat{Q}^{[l+1]}$ jedoch zu einer geschätzten Q-Function führen kann, die möglicherweise nicht mehr in dem durch $\phi(\cdot)$ aufgespannten Raum liegt, wird eine orthogonale Projektion mit der Projektionsmatrix $\bar{\Phi} \left(\bar{\Phi}^\top \bar{\Phi} \right)^{-1} \bar{\Phi}^\top$ [Mey00, S. 430] mit $\bar{\Phi}$ nach (6.30b) vorgenommen. Dies führt zur Forderung $\hat{Q}^{[l+1]} \stackrel{!}{=} \bar{\Phi} \left(\bar{\Phi}^\top \bar{\Phi} \right)^{-1} \bar{\Phi}^\top \mathcal{B}\hat{Q}^{[l+1]}$ (vgl. [LP03, S. 1117]). Eine ausführliche Beschreibung ist in [LP03], [BBdE10] und [Kil20, Abschnitt 2.3.6] zu finden. Für den Fall, dass die Q-Function exakt durch die gewählten Funktionsapproximatoren abgebildet werden kann, führt die verwendete Fixpunktforderung zur selben Lösung wie die in Abschnitt 3.2.3 beschriebene Methode (vgl. [LP03, Abschnitt 5.3]). Da die Fixpunktforderung nach Aussage von [LP03, Abschnitt 5.3] jedoch in Experimenten häufig zu besseren Regelgesetzen führt und für die Anwendung empfohlen wird, findet diese Methode im Folgenden Verwendung.

mit

$$\phi_k := \phi(\mathbf{x}_k, \mathbf{z}_k, u_k), \quad k = 1, \dots, M, \quad (6.30c)$$

$$r_k := r(\mathbf{x}_k, s_{\Gamma}(\mathbf{z}_k), u_k) \quad (6.30d)$$

und

$$\phi_k^{+[l]} := \phi\left(\mathbf{x}_{k+1}, \mathbf{z}_k^{(1)}, \hat{\mu}^{[l]}(\mathbf{x}_{k+1}, \mathbf{z}_k^{(1)})\right). \quad (6.30e)$$

Aufgrund der Off-Policy-Charakteristik können die Messdaten in jeder Iteration wiederverwendet werden. Zudem erfordert die Minimierung in (6.22), dass

$$\frac{\partial \hat{Q}^{[l+1]}}{\partial u_k} = 2 \left(\mathbf{h}_{\text{ux}}^{[l+1]} \mathbf{x}_k + \mathbf{h}_{\text{uz}}^{[l+1]} \mathbf{z}_k + h_{\text{u1}}^{[l+1]} + h_{\text{uu}}^{[l+1]} u_k \right) \stackrel{!}{=} \mathbf{0} \quad (6.31)$$

gilt. Hieraus folgt für den Policy-Improvement-Schritt (6.22) der analytische Ausdruck

$$\hat{\mu}^{[l+1]}(\mathbf{x}_k, \mathbf{z}_k) = - \underbrace{\left(h_{\text{uu}}^{[l+1]} \right)^{-1} \begin{bmatrix} \mathbf{h}_{\text{ux}}^{[l+1]} & \mathbf{h}_{\text{uz}}^{[l+1]} & h_{\text{u1}}^{[l+1]} \end{bmatrix}}_{=: \mathbf{K}^{[l+1]}} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{z}_k \\ 1 \end{bmatrix} \quad (6.32)$$

(vgl. Satz 3.2). Somit hängt der Motorstrom I_k von $\mathbf{x}_k, \mathbf{z}_k$ und einer statischen Offsetkorrektur ab.

Bemerkung 6.3

Die Wahl von $\hat{Q}(\cdot)$ in (6.29) erweitert die in Abschnitt 3.2.4 präsentierte Approximationsstruktur um einen konstanten Term in $\bar{\mathbf{y}}_k$. Hierdurch wird insbesondere auch die Darstellung eines statischen Offsetstroms ermöglicht. Ein solcher statischer Offsetstrom kann beispielsweise eine asymmetrische Masseverteilung der Platte kompensieren. Während dieser Offsetstrom in [KIB⁺ 19] heuristisch ermittelt wurde, wird er in der vorliegenden Arbeit somit automatisiert gelernt.

6.2.2.3 Trainingsablauf

Der zur Adaption eingesetzte LSPI-Algorithmus [LP03] verwendet Trainingsdaten, um iterativ die geschätzte Q-Funktion und somit das Regelgesetz anzupassen. Die während des Trainings genutzten Datentupel bestehen aus zwei Teilen. Einerseits werden am realen System aufgezeichnete Messdaten verwendet, andererseits werden ADP-kompatible Solltrajektorien erzeugt. Diese Trainingsdaten werden anschließend vorverarbeitet und normiert, bevor eine Policy Iteration zum Training durchgeführt wird.

Reale Messdaten

Während der Aufzeichnung von Messdaten am realen System muss sowohl das System ausreichend angeregt werden¹³⁵, als auch ein sicherer Betrieb gewährleistet sein. Inspiriert durch Satz 5.2 und die Simulationsergebnisse aus Abschnitt 5.7.3 hat sich zur Anregung auch im vorliegenden Anwendungsbeispiel eine Überlagerung hinreichend vieler unterschiedlicher Frequenzen, die sich nicht gegenseitig auslöschen, bewährt. Da für die Betrachtung des modellfreien ADP-basierten Ansatzes angenommen wird, dass kein genaues Modell der Systemdynamik vorliegt, ist der Entwurf eines Reglers, der das System geeignet anregt (vgl. beispielsweise Abschnitt 5.7.2) und dabei gleichzeitig einen sicheren Betrieb des Systems gewährleistet, nicht trivial. Daher wird zur Aufzeichnung von Trainingsdaten am realen System eine Anregung des Systems durch einen Menschen genutzt. Dieser kann mithilfe der in Abbildung 6.13 gezeigten mechanischen Drehregler den Plattenwinkel manuell in X - und Y -Richtung beeinflussen. Die am Ball-auf-Platte-System eingebauten Inkrementalgeber sowie eine über der Platte befindliche Kamera zeichnen währenddessen, wie in [KIB⁺19] beschrieben, die Systemzustände $\mathbf{x}_k^{[X]}$ und $\mathbf{x}_k^{[Y]}$ nach (6.14) auf. Zudem wird der applizierte Motorstrom gespeichert. Somit werden Datentupel $\{\mathbf{x}_k, u_k, \mathbf{x}_{k+1}\}$ mit einer Abtastzeit von $\Delta t = 40$ ms erzeugt. An dieser Stelle sei darauf hingewiesen, dass der Mensch bei der Datenaufzeichnung zwar eine ausreichende Anregung sicherstellen muss, also Plattenwinkel und Ballposition variiert werden müssen, jedoch dank der Off-Policy-Charakteristik der anschließend verwendeten ADP-Methode dabei weder optimal noch entsprechend des Initialgewichts $\hat{\mathbf{w}}^{[l]}$ oder einer vorgegebenen Solltrajektorie gehandelt werden muss.

ADP-kompatible Solltrajektorien

Wiederum durch die in Kapitel 5 betrachtete Überlagerung mehrerer Anregungsfrequenzen motiviert, wird der Sollpositionsverlauf des Balls für die Trainingsdaten in Form einer Summe aus Sinus- und Kosinusfunktionen konstruiert. Konkret wurde hierzu $s_{r,\text{soll},k}$ durch Abtastung von

$$0,1 \left(\sin(0,4t) + \sin(t) + 0,5 \sin(1,4t) + \cos(1,6t) + 0,2 \sin(2t) \right) \quad (6.33)$$

mit $\Delta t = 40$ ms erzeugt. Zu jedem Zeitschritt k wird dieser Sollverlauf mittels gewichteter Least-Squares-Approximation nach (6.25) durch ein quadratisches Polynom ($d = 2$) mit $\beta = 0,8$ und $n_h = 10$ in Form des Parameters z_k beschrieben. Anschließend wird z_k mithilfe von $\mathbf{D}(1)$ nach (6.24) propagiert, um $z_k^{(1)}$ im Sinne einer ADP-kompatiblen Solltrajektorien-darstellung zu erhalten.

Vorverarbeitung und Normierung

Die aufgezeichneten Messdaten \mathbf{x}_k , u_k , und \mathbf{x}_{k+1} werden zunächst mithilfe eines gleitenden Mittelwertfilters der Länge 5 geglättet. Die resultierenden Zustands- und Stellgrößenverläufe sind Abbildung 6.14 zu entnehmen. Anschließend werden diese Messdaten zusammen mit

¹³⁵ Ähnlich wie zuvor in Kapitel 5 und der Rangbedingung in (3.29) muss auch hier eine von den Systemzuständen, Solltrajektorienparametern und Basisfunktionen abhängige Matrix vollen Rang aufweisen. Konkret muss die Existenz der in (6.30a) auftretenden Inverse sichergestellt sein.

den Trainingsreferenzparametern z_k und $z_k^{(1)}$ zu Tupeln $\{x_k, u_k, x_{k+1}, z_k, z_k^{(1)}\}$ zusammengefasst. Für den Trainingsvorgang werden $M = 1200$ Datentupel verwendet. Bei einer Abtastzeit von $\Delta t = 40$ ms entspricht dies einer Aufzeichnungsdauer von 48 s. Zudem werden die Zustände und Solltrajektorienparameter mit einem Normierungsfaktor $c_{\text{norm}} = 10$ skaliert ($x_{\text{norm},k} = c_{\text{norm}} x_k$ und $z_{\text{norm},k} = c_{\text{norm}} z_k$)¹³⁶.

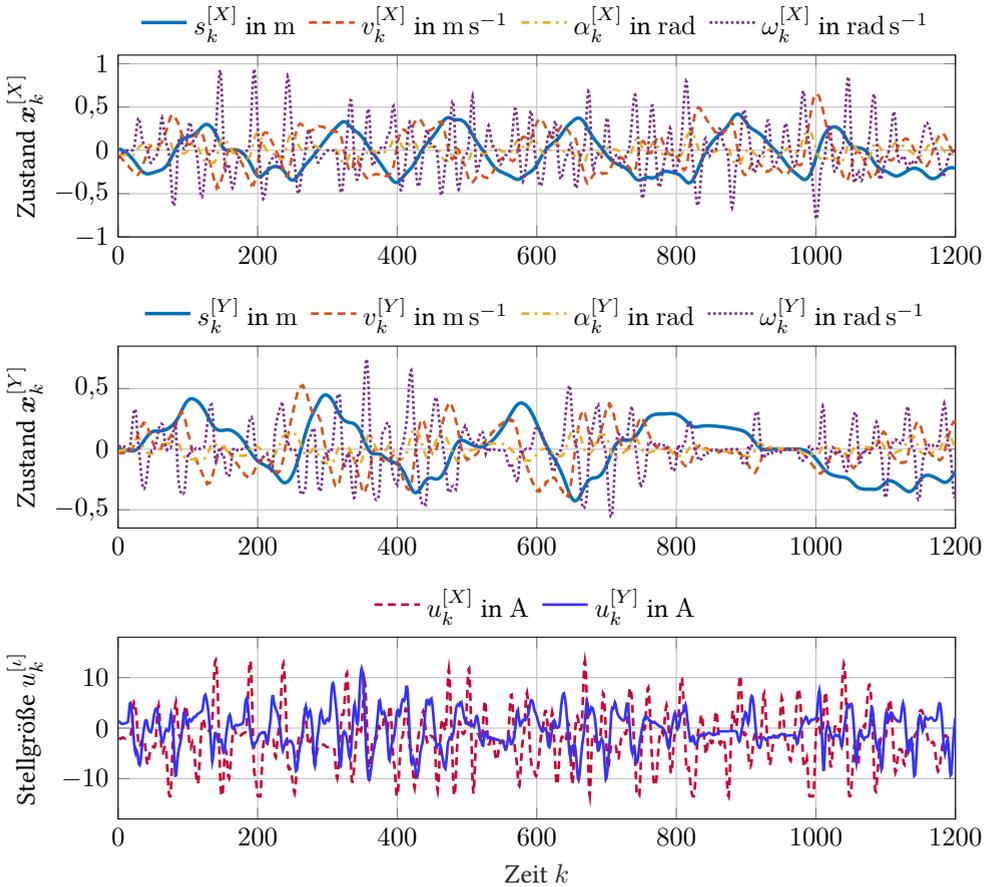


Abbildung 6.14: Mithilfe der mechanischen Drehregler durch einen Menschen aufgezeichnete, geglättete Messdaten des Ball-auf-Platte-Systems. Die Abtastzeit beträgt $\Delta t = 40$ ms.

¹³⁶ Abhängig vom gewählten Verhältnis der Parameter Q und R im Gütefunktional nach (6.17) ist der Wertebereich der Stellgröße u_k zumeist deutlich größer als der Wertebereich der Zustandsgrößen x_k . Die Verwendung des heuristisch ermittelten Skalierungsfaktors c_{norm} wurde daher zugunsten einer verbesserten numerischen Stabilität eingeführt.

Parametrierung und Algorithmus

Primäres Ziel des Reglers ist, dass die Ballposition dem Solltrajektorienverlauf möglichst präzise folgt. Zusätzlich soll die Platte nach Möglichkeit in einer horizontalen Lage gehalten und der Stellaufwand bestraft werden. Im Folgenden werden dazu

$$\mathbf{Q}_{\text{norm}} := \frac{\mathbf{Q}}{c_{\text{norm}}^2} = \text{diag}(800, 0, 400, 0) \quad \text{und} \quad R = 1 \quad (6.34)$$

gewählt. Hierbei stellt \mathbf{Q} die Gewichtungsmatrix nach (6.17) dar und \mathbf{Q}_{norm} entspricht der Gewichtung der normierten Größen $\mathbf{x}_{\text{norm},k}$ und $\mathbf{z}_{\text{norm},k}$. Somit wird insbesondere eine Abweichung der Ballposition s_k von dem durch \mathbf{z}_k parametrisierten Sollpositionsverlauf $s_r(\mathbf{z}_k^{(\kappa)})$ nach (6.23) sowie die Abweichung des Plattenwinkels α_k von der Horizontalen $\alpha = 0$ bestraft. Zudem werde der Diskontierungsfaktor zu $\gamma = 0,9$ gewählt. Für die erste Iteration werden alle Gewichte $\hat{\mathbf{w}}^{[0]}$ zu 1 initialisiert¹³⁷.

Mithilfe der in Abschnitt 6.2.2.2 beschriebenen Q-Function-Approximation sowie des LSPI-Algorithmus wird in jeder Iteration l eine Adaption des Gewichtsvektors $\hat{\mathbf{w}}^{[l]}$ vorgenommen¹³⁸. Der Algorithmus endet, sobald die zu

$$\left\| \hat{\mathbf{w}}^{[l]} - \hat{\mathbf{w}}^{[l-1]} \right\|_2 \leq e_{\hat{\mathbf{w}}} = 10^{-6} \quad (6.35)$$

gewählte Abbruchbedingung erfüllt ist. Aus dem letzten Policy-Improvement-Schritt resultiert nach einer Rücknormierung schließlich die Reglermatrix \mathbf{K} (vgl. (6.32)) und somit das Regelgesetz

$$\hat{\boldsymbol{\mu}}(\mathbf{x}_k, \mathbf{z}_k) = - \begin{bmatrix} \mathbf{K}_x & \mathbf{K}_z & K_{\text{off}} \end{bmatrix} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{z}_k \\ 1 \end{bmatrix}. \quad (6.36)$$

Diese Schritte sind im Ablaufdiagramm in Abbildung 6.15 veranschaulicht.

6.2.3 Ergebnisse

Der modellfreie, ADP-basierte Regler \mathbf{K}_{ADP} wird zu Vergleichszwecken mit einem modellbasierten Optimalregler $\mathbf{K}_{\text{Modell}}$ verglichen. Da die Regler für die beiden Plattendimensionen X und Y auf gleiche Weise trainiert bzw. berechnet werden, wird zunächst nur eine Dimension betrachtet, bevor schließlich in Abschnitt 6.2.3.4 eine gleichzeitige Regelung beider Plattendimensionen erfolgt.

¹³⁷ Zwar entspricht dies einem instabilen initialen Regelgesetz, aufgrund der Diskontierung $\gamma = 0,9$ konvergiert der erste Policy-Evaluation-Schritt in diesem Fall aber dennoch gegen eine endliche Lösung und die Policy Iteration letztlich gegen einen stabilisierenden Regler. Unter Nutzung von Vorwissen über die Systemdynamik könnte $\hat{\mathbf{w}}^{[0]}$ alternativ so initialisiert werden, dass $\mathbf{K}_x = \begin{bmatrix} 10 & 10 & 100 & 10 \end{bmatrix}$ (vgl. (6.36)) gilt. Dies entspricht einem stabilisierenden, jedoch suboptimalen initialen Regelgesetz.

¹³⁸ Die Komplexität jeder Iteration wird hierbei durch den Policy-Evaluation-Schritt mit $\mathcal{O}(h^3 + h^2M)$ dominiert.

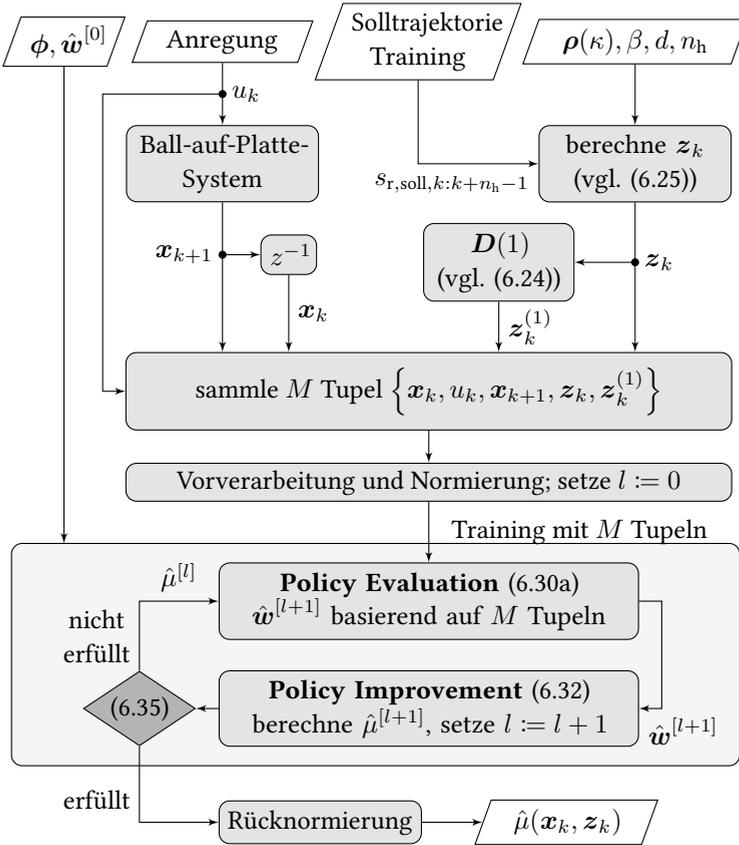


Abbildung 6.15: Ablaufschema des Trainingsvorgangs für den modellfreien ADP-Solltrajektorienfolgeregler für das Ball-auf-Platte-System. Hieraus resultiert der geschätzte Optimalregler $\hat{\mu}(x_k, z_k)$, welcher die durch z_k beschriebene lokale Approximation des Solltrajektorienverlaufs explizit berücksichtigt.

Der ADP-basierte Regler wird wie in Abschnitt 6.2.2.3 beschrieben trainiert. Die Konvergenz von $\hat{w}^{[l]}$ ist in Abbildung 6.16 zu sehen. Die gelernte Reglermatrix ergibt sich schließlich zu

$$\mathbf{K}_{\text{ADP}}^{[Y]} = \underbrace{[64,8 \quad 32,3 \quad 145,3 \quad 16,2]}_{\mathbf{K}_x} \quad \underbrace{[-27,9 \quad -36,9 \quad -60,7]}_{\mathbf{K}_z} \quad \underbrace{[-0,1]}_{\mathbf{K}_{\text{off}}}. \quad (6.37)$$

Die modellbasierte Vergleichslösung, die das in [KIB⁺19] gegebene Systemmodell mit entsprechend angepassten Parametern der verwendeten Platte und des verwendeten Balls nutzt, wird nach Satz 3.1 berechnet. Daraus ergibt sich die Reglermatrix

$$\mathbf{K}_{\text{Modell}}^{[Y]} = \underbrace{[75,5 \quad 55,9 \quad 213,3 \quad 33,0]}_{\mathbf{K}_x} \quad \underbrace{[-41,1 \quad -55,9 \quad -75,1]}_{\mathbf{K}_z} \quad (6.38)$$

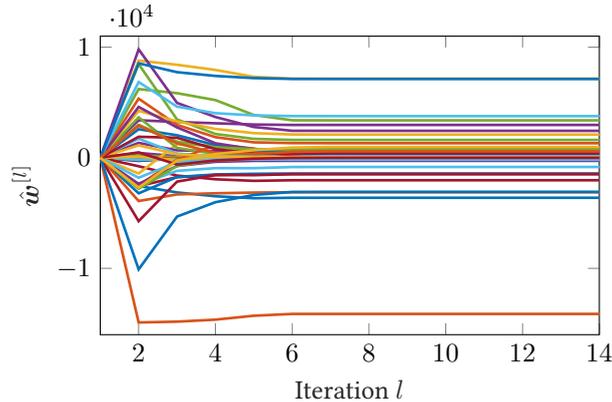


Abbildung 6.16: Gewichtsvektor $\hat{w}^{[l]}$ über die Iterationen l des LSPI-Algorithmus bis zur Erfüllung der Abbruchbedingung nach (6.35).

des modellbasierten Vergleichsreglers. Im Folgenden wird zunächst eine konstante Sollpositionsvorgabe, d. h. $d = 0$, und anschließend die Verwendung einer Approximation des Solltrajektorienverlaufs für Polynome vom Grad $d = 2$ betrachtet.

6.2.3.1 Konstante Sollvorgabe: geglätteter Rechteckverlauf

Der gelernte, modellfreie ADP-Regler wird anhand eines geglätteten Rechteckverlaufs¹³⁹ als Sollpositionsvorgabe $s_{r,\text{soll}}$ des Balls mit dem berechneten, modellbasierten Regler verglichen. Hierzu wird zunächst eine konstante Sollvorgabe ($d = 0$) verwendet. In Abbildung 6.17 ist die mittlere Ballposition über 11 Wiederholungen für den gelernten Regler in Blau und für den berechneten, modellbasierten Regler in Gelb gegeben. Die jeweilige Standardabweichung ist transparent dargestellt. Der gesamte Systemzustand \mathbf{x}_k sowie die Stellgröße u_k sind in Abbildung D.1 in Anhang D.4 visualisiert.

Bei beiden Reglern liegt eine merkliche zeitliche Verzögerung der Ballposition gegenüber der Sollposition vor. Dies ist darin begründet, dass den Reglern lediglich die aktuelle Sollposition im Zeitschritt k , nicht jedoch deren zukünftiger Verlauf zur Verfügung steht, das System träge ist und die Stellgrößen im Gütefunktional bestraft werden. Der gelernte Regler reagiert dabei etwas schneller und regelt die Ballposition präziser, was sich, wie Abbildung 6.17 zu entnehmen ist, auch in geringeren akkumulierten Einschrittkosten

$$\sum_{\kappa=0}^k r(\mathbf{x}_{\kappa}, u_{\kappa}, s_{\mathbf{r}}(\mathbf{z}_{\kappa})) \quad (6.39)$$

widerspiegelt.

¹³⁹ Genauer gesagt wird der Sollpositionsverlauf aus einem um 0,15 m verschobenen sinusförmigen Signal mit einer Amplitude von 0,15 m und einer Periodendauer von 147 Zeitschritten, das im Maximum für 200 Zeitschritte konstant gehalten wird, erzeugt.

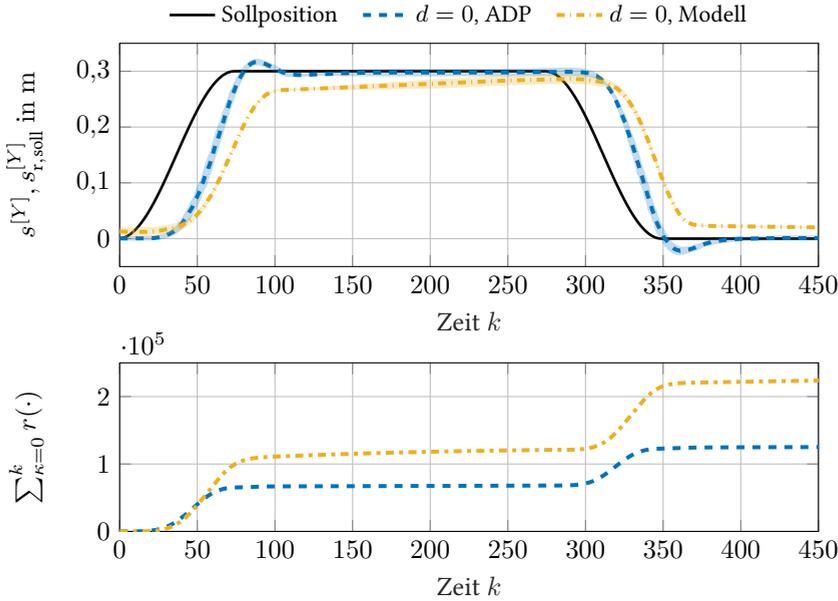


Abbildung 6.17: Vorgabe einer stationären Sollposition ($d = 0$) für den gelernten, ADP-basierten Regler (blau) sowie den modellbasierten Vergleichsregler (gelb). *Oben:* Mittlere Ballposition und Standardabweichung über 11 Wiederholungen. *Unten:* Mittlere akkumulierte Einschnittkosten.

6.2.3.2 Approximation des Sollverlaufs: geglätteter Rechteckverlauf

Für $d = 2$ zeigt Abbildung 6.18 den Vergleich eines gelernten, ADP-basierten Solltrajektorienfolgereglers (rot) mit einem modellbasierten Solltrajektorienfolgeregler (grün). Beide Regler folgen dem Sollpositionsverlauf deutlich besser als der Regler mit konstanter Sollvorgabe, der zu Vergleichszwecken in Blau dargestellt ist¹⁴⁰. Aufgrund der lokalen Approximation der Solltrajektorie nutzen die Regler mit $d = 2$ Information über den zukünftigen Verlauf der Sollposition. Dies ermöglicht signifikant geringere akkumulierte Einschnittkosten, wie Abbildung 6.18 zu entnehmen ist. Ähnlich wie bei den in Abschnitt 6.2.3.1 betrachteten Reglern mit konstanter Sollvorgabe weist auch im Fall $d = 2$ der gelernte, ADP-basierte Regler geringere Gesamtkosten im Vergleich zum modellbasierten Regler auf.

6.2.3.3 Approximation des Sollverlaufs: Validierungstrajektorie

In diesem Abschnitt wird die Vorgabe eines Sollpositionsverlaufs betrachtet, der sich aus einer Summation von Sinussignalen sowie aus Sprüngen und Rampen zusammensetzt. Hierdurch soll insbesondere die in Abschnitt 2.4.1 beschriebene und in der Simulation in Abschnitt 3.2.5

¹⁴⁰ Der gesamte Systemzustand \mathbf{x}_k und die Stellgröße u_k unter Verwendung der drei hier betrachteten Regler sind in Abbildung D.2 in Anhang D.4 gegeben.

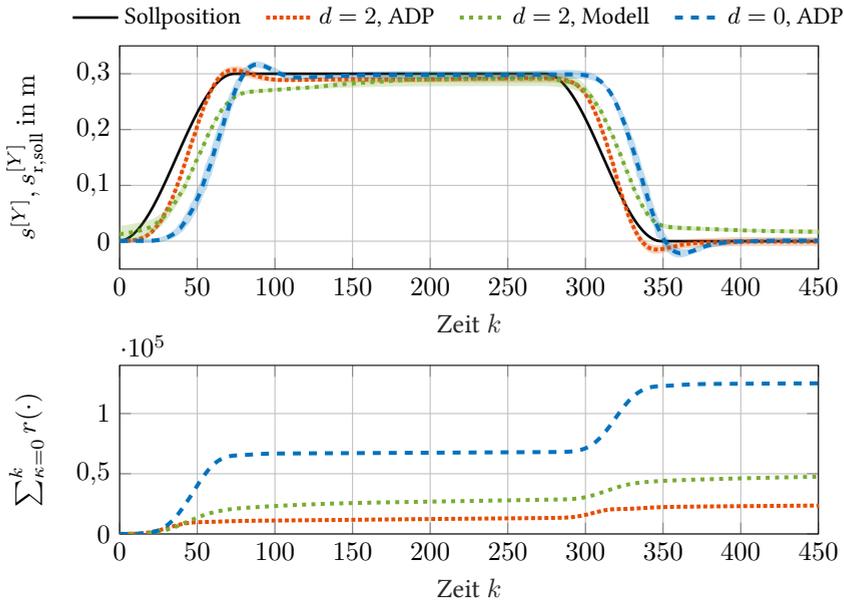


Abbildung 6.18: Vorgabe eines polynomiellen Sollverlaufs mit $d = 2$ für den gelernten, ADP-basierten Regler (rot) sowie den modellbasierten Vergleichsregler (grün) im Vergleich zur Vorgabe einer stationären Sollposition ($d = 0$) für den gelernten, ADP-basierten Regler (blau). *Oben:* Mittlere Ballposition und Standardabweichung über 18, 12, bzw. 11 Wiederholungen. *Unten:* Mittlere akkumulierte Einschnittkosten.

bereits gezeigte Flexibilität des neuartigen ADP-basierten Solltrajektorienfolgereglers verdeutlicht werden, die auch in der realen Anwendung Gültigkeit besitzt. Der gelernte Regler kann unterschiedlichen Solltrajektorienvorgaben erfolgreich folgen, ohne dass ein erneuter Trainingsvorgang benötigt wird. In Abbildung 6.19 ist neben der Solltrajektorie der resultierende Verlauf für den gelernten ADP-Solltrajektorienfolgeregler ($d = 2$) sowie für den ADP-Sollzustandsregler ($d = 0$) gezeigt¹⁴¹. Auch hier offenbaren sich die Vorteile des ADP-Solltrajektorienfolgereglers, der den Sollverlauf explizit einbezieht, im Gegensatz zur konstanten Sollzustandsvorgabe. Ersterer kann, aufgrund der lokalen Solltrajektorienapproximation, dem Sollverlauf besser folgen, wodurch insbesondere deutlich geringere akkumulierte Kosten entstehen.

¹⁴¹ Der gesamte Systemzustand \mathbf{x}_k und die Stellgröße u_k sind in Abbildung D.3 in Anhang D.4 gezeigt.

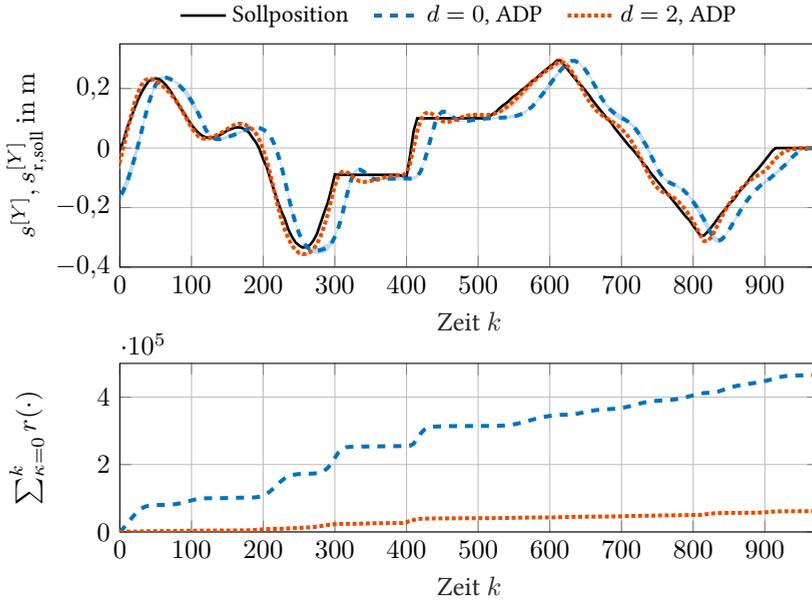


Abbildung 6.19: Vergleich der gelernten Solltrajektorienfolgeregler bei Vorgabe eines polynomiellen Sollverlaufs ($d = 2$, rot) und einer stationären Sollposition ($d = 0$, blau) für eine Validierungstrajektorie. *Oben:* Mittlere Ballposition und Standardabweichung über 4 Wiederholungen. *Unten:* Mittlere akkumulierte Einschrittkosten.

6.2.3.4 Offsetkorrektur und zweidimensionale Regelung mit Approximation des Sollverlaufs

Der Regler, der die X -Dimension des Ball-auf-Platte-Systems regelt, wird mit derselben Methodik und denselben Parametern wie der Regler für die Y -Dimension trainiert. Daraus ergibt sich die durch

$$\mathbf{K}_{ADP}^{[X]} = \underbrace{\begin{bmatrix} 65,3 & 37,0 & 135,1 & 18,6 \end{bmatrix}}_{\mathbf{K}_x} \underbrace{\begin{bmatrix} -28,8 & -38,2 & -61,2 \end{bmatrix}}_{\mathbf{K}_z} \underbrace{\begin{bmatrix} 2,2 \end{bmatrix}}_{\mathbf{K}_{off}} \quad (6.40)$$

gegebene gelernte Reglermatrix. Aufgrund einer nahezu symmetrischen Masseverteilung der Platte in Y -Richtung findet in $\mathbf{K}_{ADP}^{[Y]}$ (6.37) nahezu keine Offsetkorrektur statt. Demgegenüber findet in X -Richtung aufgrund von $K_{off}^{[X]} = 2,2$ in (6.40) eine statische Offsetkorrektur in Form eines konstanten Ausgleichstroms von $-2,2$ A statt, um ein konstruktions- oder fertigungsbedingtes Ungleichgewicht auszugleichen. Im Fall des modellbasierten Vergleichsreglers muss solch ein Ausgleichstrom in der Regel heuristisch ermittelt werden (vgl. [KIB⁺19]), da potenzielle Ungleichgewichte und Fertigungstoleranzen meist nicht präzise im Systemmodell berücksichtigt sind. Wie wichtig jedoch die Verwendung eines solchen Ausgleichstroms ist, zeigt Abbildung 6.20¹⁴². Bei Verwendung einer Q-Function-Approximation, die keine

¹⁴² Der gesamte Systemzustand \mathbf{x}_k sowie die Stellgröße u_k sind in Abbildung D.4 in Anhang D.4 gezeigt.

Offsetkorrektur ermöglicht, kommt es zu einem asymmetrischen Verhalten der Ballposition und somit zu signifikanten Abweichungen von der Sollposition. Im Vergleich dazu führt der gelernte Offsetterm zu einem symmetrischen Verhalten.

Nach Satz 3.1 ergibt sich die Reglermatrix

$$\mathbf{K}_{\text{Modell}}^{[X]} = \underbrace{\begin{bmatrix} 75,9 & 56,8 & 217,0 & 34,8 \end{bmatrix}}_{\mathbf{K}_x} \underbrace{\begin{bmatrix} -45,3 & -56,8 & -75,4 \end{bmatrix}}_{\mathbf{K}_z} \quad (6.41)$$

des modellbasierten Vergleichsreglers, für den ebenfalls ein (in diesem Fall heuristisch ermittelter) Ausgleichsstrom von $-2,2$ A verwendet wird. Die gleichzeitige Vorgabe eines Solltrajektorienverlaufs für die Position des Balls in beiden Plattendimensionen ist schließlich sowohl für den ADP-basierten Regler als auch für den Vergleichsregler in Abbildung 6.21 gezeigt. Die zugehörigen Zustands- und Stellgrößenverläufe sind in Abbildung D.5 und Abbildung D.6 (Anhang D.4) gegeben.

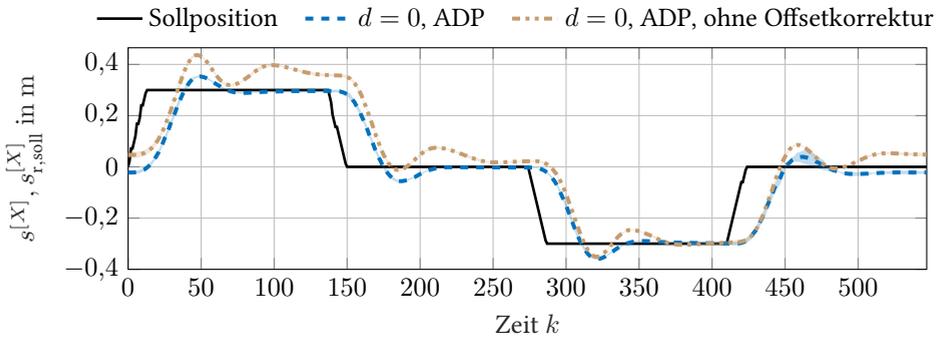


Abbildung 6.20: Vergleich eines gelernten Reglers mit (blau) und ohne (braun) lernbare Offsetkorrektur bei Vorgabe einer stationären Sollposition ($d = 0$).

6.2.4 Diskussion

In diesem Unterkapitel wurde ein modellfreier, ADP-basierter Solltrajektorienfolgereger für ein reales Ball-auf-Platte-System vorgestellt. Mit weniger als einer Minute aufgezeichneter Messdaten konnte ein Regler trainiert werden, der gegebenen Solltrajektorienverläufen folgen kann und dabei hinsichtlich des zugrunde liegenden Gütefunktionalen einem modellbasierten Vergleichsregler überlegen ist. Letzteres resultiert aus einer etwas schnelleren Sprungantwort sowie einem geringeren statischen Fehler: Während der modellbasierte Vergleichsregler Abweichungen vom Systemmodell nicht berücksichtigen kann, verwendet der ADP-Ansatz reale Messdaten und kann so, im Rahmen der Approximationsfähigkeit der verwendeten Basis-

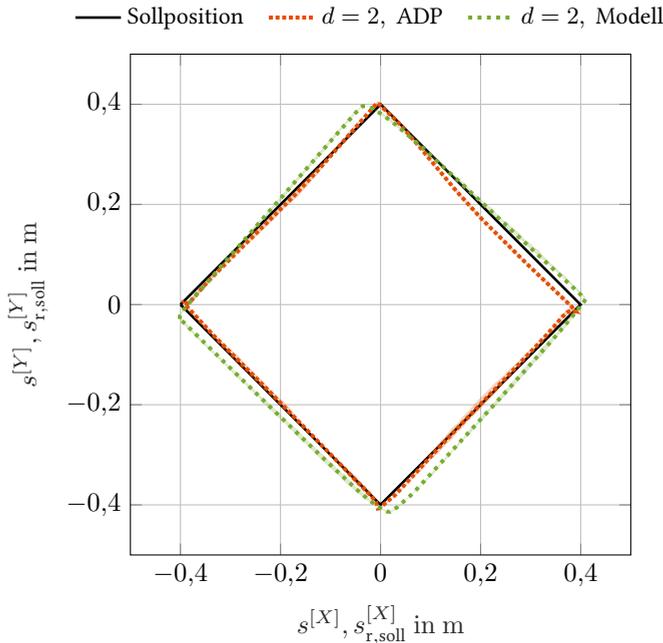


Abbildung 6.21: Gleichzeitige Vorgabe eines Solltrajektorienverlaufs in beiden Plattendimensionen und resultierende Ballposition für den gelernten Regler und den modellbasierten Vergleichsregler mit jeweils $d = 2$. Gezeigt sind die mittlere Ballposition sowie die Standardabweichung über jeweils 3 Wiederholungen.

funktionen, das Systemverhalten bestmöglich im Sinne der zu minimierenden Gütefunktion berücksichtigen¹⁴³.

Die Verwendung einer ADP-kompatiblen lokalen Approximation des Sollzustandsverlaufs anstelle einer stationären Sollzustandsvorgabe ermöglicht ein vorausschauendes Verhalten des Reglers, einen verringerten zeitlichen Versatz des Zustands zum Sollzustand und letztlich reduzierte Kosten.

Die experimentellen Ergebnisse zeigen letztlich auch anhand des zweiten realen Anwendungsbeispiels eine grundsätzliche Anwendbarkeit der in dieser Arbeit vorgestellten ADP-kompatiblen Solltrajektorienfolgeregelungsmethoden auf reale, regelungstechnische Problemstellungen. Aufgrund der messdatenbasierten Adaption der Gewichte der Q-Funktion-Approximation und somit des Regelgesetzes ist weder eine aufwendige exakte Modellbildung noch eine manuelle Feinabstimmung notwendig.

¹⁴³ Die Basisfunktionen des ADP-Ansatzes sind aus Gründen der Fairness gegenüber dem Vergleichsregler so gewählt, dass die resultierenden Regelgesetze dieselbe Form aufweisen. Dies ist in (6.37) und (6.38) zu erkennen, wobei anzumerken sei, dass auch im Fall des modellbasierten Reglers ein Ausgleichsstrom, der jedoch heuristisch ermittelt wurde, verwendet wird.

Schließlich soll noch die Wahl der Basisfunktionen zur Approximation der Q-Function diskutiert werden. Der vorgestellte Ansatz kann ohne Weiteres um zusätzliche Basisfunktionen erweitert werden, um beispielsweise auch Nichtlinearitäten im Systemverhalten abbilden zu können. Dies würde jedoch zu einer größeren Anzahl h an Gewichten und somit einer höheren Komplexität der Policy Iteration (vgl. Abschnitt 6.2.2.3) führen. Neben einem höheren Rechenaufwand wären zudem vor allem mehr Trainingstupel M und somit mehr reale Messdaten benötigt. Erfahrungsgemäß erschwert dies auch eine geeignete Anregung und kann in Kombination mit vorhandenem Messrauschen die Trainingsergebnisse verschlechtern. Somit ist eine angemessene Wahl der Basisfunktionen entscheidend für einen dateneffizienten ADP-Algorithmus. Die Wahl der Basisfunktionen sollte mit Bedacht erfolgen und bietet die Chance, vorhandenes Vorwissen über geeignete Reglerstrukturen oder prinzipielles Systemverhalten zu integrieren. Am Beispiel des betrachteten Ball-auf-Platte-Systems wurde auf diese Weise Vorwissen über geeignete Basisfunktionen eingebracht und ein erlernbarer statischer Ausgleichsstrom verwendet. Die Nutzung dieses Vorwissens auf der einen Seite und die Adaptierbarkeit der Q-Function-Gewichte auf der anderen Seite haben letztlich einen dateneffizienten adaptiven Solltrajektorienfolgeregler ermöglicht.

Zusammenfassend kann konstatiert werden, dass Kapitel 6 die in Abschnitt 2.4.1 formulierte Forschungsfrage 1 aus einer anwendungsorientierten Perspektive betrachtet. Der neuartige, ADP-basierte Trajektorienfolgeregler, der vielfältige und flexible Solltrajektorienvorgaben erlaubt, kann anhand realer Messdaten erfolgreich trainiert werden. Zudem weist die Einbeziehung des approximierten Solltrajektorienverlaufs auch in der realen Anwendung Vorteile gegenüber einer stationären Sollzustandsvorgabe auf.

7 Zusammenfassung

Obwohl selbstlernende, ADP-basierte Regelungsansätze vermehrt in den Fokus der aktuellen Forschung gerückt sind, ist die Integration flexibler Solltrajektorien in diese modellfreien Methoden bislang nur unzureichend gelöst. Zahlreiche regelungstechnische Problemstellungen erfordern jedoch, dass Systemgrößen vorgegebenen und flexiblen Solltrajektorienverläufen optimal im Sinne eines gegebenen Gütefunktionalen folgen. Die Analyse bestehender ADP-Ansätze aus der Literatur offenbart, dass diese hauptsächlich stationäre Sollzustandsvorgaben oder globale Vorgaben durch Exosysteme betrachten. Ersteren fehlt die Berücksichtigung des weiteren Solltrajektorienverlaufs, letzteren die Möglichkeit, den Sollverlauf flexibel von außen vorzugeben. Eine zweite bislang in der ADP-Literatur nur unzureichend analysierte Problematik betrifft eine angemessene Systemanregung, um Konvergenz zu gewährleisten. Wenngleich ADP-basierte Ansätze stets eine geeignete Anregung erfordern, finden sich hierzu in der Literatur bislang nur wenige theoretische Analysen.

Die vorliegende Arbeit liefert daher Beiträge zu diesen bislang ungelösten wissenschaftlichen Fragestellungen. Insbesondere wurden für adaptive Optimalregler geeignete Solltrajektorienrepräsentationen präsentiert, analysiert und in realen Anwendungsbeispielen umgesetzt. Weiterhin wurden Konvergenzeigenschaften von ADP-Reglern untersucht und Bedingungen an den Systemzustand präsentiert, die eine ausreichende Anregung gewährleisten.

Zunächst wurden Konzepte zur effizienten Einbettung flexibler, zur Laufzeit vorgegebener Solltrajektorienverläufe in den ADP-Mechanismus entwickelt sowie die Existenz und Stabilität der Lösung der zugehörigen Optimierungsprobleme analysiert. Um Solltrajektorienverläufe in einen modellfreien ADP-Formalismus integrieren zu können, wurde in der vorliegenden Arbeit erstmals der Begriff der ADP-kompatiblen Solltrajektorienrepräsentation in zeitdiskreter und zeitkontinuierlicher Form definiert. Als hierfür zentrale Eigenschaft, nicht nur des Systemzustands, sondern insbesondere auch der Beschreibung des Solltrajektorienverlaufs, wurde die Einhaltung der Markov-Bedingung identifiziert. Diese ermöglicht eine adäquate und zeitinvariante Kostenrepräsentation in Form einer erlernbaren Value- oder Q-Function, anhand derer Regelgesetze verbessert werden können. Mithilfe der vorgestellten ADP-kompatiblen Solltrajektorienrepräsentationen wurden modellfreie, adaptive optimale Trajektorienfolgeregler entworfen. Aufgrund der expliziten Abhängigkeit der Value- bzw. Q-Function und somit auch des damit verknüpften Regelgesetzes vom aktuellen Solltrajektorienverlauf ergeben sich zwei wesentliche Vorteile: Erstens behält ein gelerntes Regelgesetz auch bei sich änderndem Solltrajektorienverlauf seine Gültigkeit. Zweitens weisen die gelernten Regler vorausschauendes Verhalten auf und sind somit Vergleichsreglern aus der Literatur mit stationärer Sollzustandsvorgabe oder ohne flexible Solltrajektorienrepräsentation hinsichtlich des Kostenfunktionalen überlegen.

Neben diesen Beiträgen zu ADP-kompatiblen Solltrajektorienfolgereglern wurde in der vorliegenden Arbeit die Fragestellung einer geeigneten Systemanregung im ADP-Kontext verfolgt. Hierbei wurden am Beispiel eines zeitkontinuierlichen Nicht-Nullsummen-Differenzialspiels erstmalig allgemeingültige hinreichende Frequenzbedingungen an den Systemzustand hergeleitet. Diese Bedingungen gewährleisten die Erfüllung der PE-Eigenschaft für die häufig verwendete Klasse polynomieller Basisfunktionen zur Funktionsapproximation. Die PE-Eigenschaft ist essenziell für Konvergenzgarantien von ADP-Methoden gegen die optimale Lösung. Die vorgestellten Bedingungen beinhalten zudem Freiheitsgrade, die in der Praxis zur Berücksichtigung anwendungsspezifischer Anforderungen genutzt werden können. Die präsentierten Simulationsergebnisse bestätigen die theoretischen Aussagen und zeigen zudem, dass eine Anregung in Form einer Überlagerung harmonischer Schwingungen zu wesentlich schnellerer Konvergenz der Policy Evaluation führen kann als die Verwendung weißen Gaußschen Rauschens. Hieraus kann schließlich die Handlungsempfehlung abgeleitet werden, Systeme, die mithilfe ADP-basierter Methoden geregelt werden sollen, nicht durch klassische Gaußsche Rauschprozesse, sondern vielmehr durch eine Überlagerung harmonischer Schwingungen oder anderer hauptsächlich niederfrequenter Signale anzuregen. Dadurch können mögliche Tiefpasscharakteristiken des Systems berücksichtigt und eine effiziente Anregung erreicht werden.

Schließlich offenbaren die beiden realen Anwendungsbeispiele die Umsetzbarkeit sowie potenzielle Vorteile der vorgestellten adaptiven optimalen Trajektorienfolgeregler. Die vorliegende Arbeit liefert die erstmalige Anwendung ADP-kompatibler Solltrajektorienfolgeregler, die einen zeitvarianten, von außen vorgebbaren, flexiblen Solltrajektorienverlauf explizit in die Q-Funktion integrieren und mithilfe realer Messdaten trainiert werden. So konnte ein selbstlernender Geschwindigkeitsregler ohne ein Modell der Longitudinaldynamik in einem Realfahrzeug umgesetzt und mithilfe eines Actor-Critic-Ansatzes online, d. h. während der Fahrt, trainiert werden. Die Verwendung der lokalen Approximation des Sollgeschwindigkeitsprofils gemäß des in dieser Arbeit vorgestellten ADP-kompatiblen Mechanismus lieferte hierbei einen Regler, der einer konstanten Sollgeschwindigkeitsvorgabe bezüglich des Kostenfunktionalis überlegen ist. Als zweites Anwendungsbeispiel wurde erstmalig ein ADP-basierter Solltrajektorienfolgeregler an einem realen Ball-auf-Platte-System angewandt. Ohne Kenntnis eines exakten Systemmodells können aus aufgezeichneten Messdaten des Systems die Q-Funktion des optimalen Trajektorienfolgeregelungsproblems und das zugehörige Regelgesetz erlernt werden. Der vorgestellte selbstlernende Regler übertrifft in experimentellen Untersuchungen den modellbasierten Vergleichsregler und erfordert zudem keine manuelle Feinabstimmung. Auch in diesem Anwendungsbeispiel zeigt die neuartige, ADP-kompatible, lokale Solltrajektorienapproximation Vorteile gegenüber einer stationären Sollzustandsvorgabe. Im Fall der neu vorgestellten Methode folgt der Systemzustand der Solltrajektorie präziser, zudem entstehen signifikant geringere akkumulierte Kosten.

Abschließend lässt sich zusammenfassen, dass die ADP-kompatible Darstellung des lokalen Solltrajektorienverlaufs eine effiziente Approximation der mit der Solltrajektorie verbundenen Gesamtkosten und somit eine erfolgreiche Adaption der Reglergewichte ermöglicht. Die grundsätzliche Anwendbarkeit der entwickelten adaptiven optimalen Solltrajektorienfolge-

regler auf reale regelungstechnische Problemstellungen wurde gezeigt. Bei geeigneter Wahl der Basisfunktionen und angemessener Systemanregung können Solltrajektorienfolgeregler erlernt werden, die ohne aufwendige Modellbildung und ohne vorherige Kenntnis der Systemparameter modellbasierte Vergleichsregler hinsichtlich der Performanz sogar übertreffen können. Dies ist insbesondere darauf zurückzuführen, dass die vorgestellten lernbasierten Ansätze im Gegensatz zu rein modellbasierten Reglerentwurfverfahren Modellierungsungenauigkeiten und Fertigungstoleranzen im Rahmen der gewählten Funktionsapproximatoren ausgleichen können. Die in dieser Arbeit präsentierten Beiträge schließen somit aktuelle Forschungslücken im Kontext selbstlernender Optimalregler und motivieren eine stärkere interdisziplinäre Verflechtung von Regelungstechnik und Maschinellern Lernen.

A Anhang zu Kapitel 3

A.1 Beweisskizze zu Satz 3.4

Beweis:

Die Grundidee des Beweises basiert auf der Verwendung der dynamischen Programmierung und der Herleitung der Lösung der Q-Function durch Rückwärtsinduktion. Für $i = 1, \dots, d$ sei die i -te Teilmatrix einer Matrix $\mathbf{\Pi} \in \mathbb{R}^{m \times nd}$ durch

$$\mathbf{\Pi}[i] := \begin{bmatrix} \mathbf{\Pi}(1, (i-1)n+1) & \cdots & \mathbf{\Pi}(1, ni) \\ \vdots & \ddots & \vdots \\ \mathbf{\Pi}(m, (i-1)n+1) & \cdots & \mathbf{\Pi}(m, ni) \end{bmatrix} \quad (\text{A.1})$$

mit $\mathbf{\Pi}[i] \in \mathbb{R}^{m \times n}$ definiert. Durch $\iota \in \mathbb{N}_{\geq 0}$ sei zunächst ein Platzhalter gegeben, der später durch η (bzw. $\eta + 1$ im Induktionsschritt) ersetzt wird, wobei $\eta = K - \kappa$ die verbliebenen Zeitschritte auf dem Horizont der Länge K bezeichne. Des Weiteren sei m ein Index mit $m \in \mathbb{N} : m > 1$. Im Folgenden werden einige Kurzschreibweisen definiert. Sei

$$\mathbf{X}_\iota^0 := \mathbf{X}^0 := [\mathbf{I}_n \quad -\mathbf{I}_n], \quad (\text{A.2})$$

$$\mathbf{X}_\iota^1 := \sqrt{\gamma} \left(-\mathbf{X}^0[1] \mathbf{B} \mathbf{G}_\iota [F_\iota \quad K_\iota^{\iota-1} \quad \cdots \quad K_\iota^0] + [\mathbf{X}^0[1] \mathbf{A} \quad \mathbf{X}^0[2] \quad \mathbf{0} \quad \cdots \quad \mathbf{0}] \right) \quad (\text{A.3})$$

und

$$\begin{aligned} \mathbf{X}_\iota^m := & \sqrt{\gamma} \left(-\mathbf{X}_\iota^{m-1}[1] \mathbf{B} \mathbf{G}_\iota [F_\iota \quad K_\iota^{\iota-1} \quad \cdots \quad K_\iota^0] \right. \\ & \left. + [\mathbf{X}_\iota^{m-1}[1] \mathbf{A} \quad \mathbf{0} \quad \mathbf{X}_\iota^{m-1}[2] \quad \cdots \quad \mathbf{X}_\iota^{m-1}[\iota-1]] \right) \end{aligned} \quad (\text{A.4})$$

sowie

$$\mathbf{U}_\iota^1 := -\mathbf{G}_\iota [F_\iota \quad K_\iota^{\iota-1} \quad \cdots \quad K_\iota^0] \quad (\text{A.5})$$

und

$$\begin{aligned} \mathbf{U}_\iota^m := & \sqrt{\gamma} \left(-\mathbf{U}_\iota^{m-1}[1] \mathbf{B} \mathbf{G}_\iota [F_\iota \quad K_\iota^{\iota-1} \quad \cdots \quad K_\iota^0] \right. \\ & \left. + [\mathbf{U}_\iota^{m-1}[1] \mathbf{A} \quad \mathbf{0} \quad \mathbf{U}_\iota^{m-1}[2] \quad \cdots \quad \mathbf{U}_\iota^{m-1}[\iota-1]] \right) \end{aligned} \quad (\text{A.6})$$

mit

$$M_\ell := \gamma \mathbf{B}^\top \left(\sum_{i=0}^{\ell-2} (\mathbf{X}_\ell^i[1])^\top \mathbf{Q} \mathbf{X}_\ell^i[1] + \sum_{i=1}^{\ell-2} (\mathbf{U}_\ell^i[1])^\top \mathbf{R} \mathbf{U}_\ell^i[1] \right), \quad (\text{A.7})$$

$$\mathbf{F}_\ell := M_\ell \mathbf{A}, \quad (\text{A.8})$$

$$\mathbf{G}_\ell^{-1} := M_\ell \mathbf{B} + \mathbf{R}, \quad (\text{A.9})$$

$$\mathbf{K}_\ell^j := \begin{cases} \gamma \mathbf{B}^\top (\mathbf{X}^0[1])^\top \mathbf{Q} \mathbf{X}^0[2], & \text{für } j = \ell - 1, \\ \gamma \mathbf{B}^\top \left(\sum_{i=1}^{\ell-2} (\mathbf{X}_\ell^i[1])^\top \mathbf{Q} \mathbf{X}_\ell^i[\ell - j] + \sum_{i=1}^{\ell-2} (\mathbf{U}_\ell^i[1])^\top \mathbf{R} \mathbf{U}_\ell^i[\ell - j] \right), & \text{für } j < \ell - 1, \end{cases} \quad (\text{A.10})$$

wobei $j \in \mathbb{N}_{\geq 0}$. Sei weiterhin

$$\boldsymbol{\rho}_0^\kappa := [\mathbf{x}_{k_\kappa}^\top \quad \mathbf{x}_{r, k_\kappa}^\top] (\mathbf{X}^0)^\top, \quad (\text{A.11})$$

$$\boldsymbol{\rho}_1^\kappa := [\mathbf{x}_{k_\kappa}^\top \quad \mathbf{u}_{k_\kappa}^\top \quad \mathbf{x}_{r, k_\kappa+1}^\top] \begin{bmatrix} (\mathbf{X}^0[1] \mathbf{A})^\top \\ (\mathbf{X}^0[1] \mathbf{B})^\top \\ (\mathbf{X}^0[2])^\top \end{bmatrix}, \quad (\text{A.12})$$

$$\boldsymbol{\mu}_i^\kappa := [\mathbf{x}_{k_\kappa}^\top \quad \mathbf{u}_{k_\kappa}^\top \quad \mathbf{x}_{r, k_\kappa+2}^\top \quad \cdots \quad \mathbf{x}_{r, k+K}^\top] \begin{bmatrix} (\mathbf{U}_\eta^{\eta-i}[1] \mathbf{A})^\top \\ (\mathbf{U}_\eta^{\eta-i}[1] \mathbf{B})^\top \\ (\mathbf{U}_\eta^{\eta-i}[2])^\top \\ \vdots \\ (\mathbf{U}_\eta^{\eta-i}[\eta])^\top \end{bmatrix}, \quad (\text{A.13})$$

$$\boldsymbol{\chi}_i^\kappa := [\mathbf{x}_{k_\kappa}^\top \quad \mathbf{u}_{k_\kappa}^\top \quad \mathbf{x}_{r, k_\kappa+2}^\top \quad \cdots \quad \mathbf{x}_{r, k+K}^\top] \begin{bmatrix} (\mathbf{X}_\eta^{\eta-i}[1] \mathbf{A})^\top \\ (\mathbf{X}_\eta^{\eta-i}[1] \mathbf{B})^\top \\ (\mathbf{X}_\eta^{\eta-i}[2])^\top \\ \vdots \\ (\mathbf{X}_\eta^{\eta-i}[\eta])^\top \end{bmatrix} \quad (\text{A.14})$$

mit $k_\kappa = k + K - \eta = k + \kappa$ und $i \in \mathbb{N}$.

Nachfolgend wird mittels Rückwärtsinduktion gezeigt, dass die zu dem System (3.80) und dem Gütefunktional (3.82) gehörende Q-Funktion ${}^K Q_{k_\kappa}$ nach Definition 3.6 durch

$$\begin{aligned} {}^K Q_{k_\kappa} = & \frac{1}{2} \left(\boldsymbol{\rho}_0^\kappa \mathbf{Q} (\boldsymbol{\rho}_0^\kappa)^\top + \mathbf{u}_{k+\kappa}^\top \mathbf{R} \mathbf{u}_{k+\kappa} + \gamma \boldsymbol{\rho}_1^\kappa \mathbf{Q} (\boldsymbol{\rho}_1^\kappa)^\top \right. \\ & \left. + \gamma \sum_{i=1}^{K-\kappa-2} \left(\boldsymbol{\chi}_i^\kappa \mathbf{Q} (\boldsymbol{\chi}_i^\kappa)^\top + \boldsymbol{\mu}_i^\kappa \mathbf{R} (\boldsymbol{\mu}_i^\kappa)^\top \right) \right) \end{aligned} \quad (\text{A.15})$$

gegeben ist. Ausgehend von ${}^K Q_K$ (vgl. (3.92)) folgt aus Definition 3.6 und

$$\left. \frac{\partial {}^K Q_K}{\partial \mathbf{u}_{k+K}} \right|_{\mathbf{u}_{k+K}^*} = \mathbf{0}, \quad \left. \frac{\partial^2 {}^K Q_K}{\partial \mathbf{u}_{k+K}^2} \right|_{\mathbf{u}_{k+K}^*} = \mathbf{R} \succ \mathbf{0} \quad (\text{A.16})$$

direkt $\mathbf{u}_{k+K}^* = \mathbf{0}$. Durch rückwärtige Iteration über die Zeit und Verwendung von (3.91) sowie der Systemdynamik (3.80) kann mit $\eta = K - \kappa$ gezeigt werden, dass (A.15) für $\eta = 0, 1, 2$, d. h. $\kappa = K, K - 1, K - 2$, gilt. Des Weiteren minimiert

$$\mathbf{u}_{k+\kappa}^* = -\mathbf{G}_\eta \left(\mathbf{F}_\eta \mathbf{x}_{k+\kappa} + \sum_{j=0}^{\eta-1} \mathbf{K}_\eta^j \mathbf{x}_{r,k+K-j} \right) \quad (\text{A.17})$$

den Ausdruck (A.15), weil aufgrund von $\mathbf{R} \succ \mathbf{0}$ und $\mathbf{Q} \succeq \mathbf{0}$

$$\left. \frac{\partial^K Q_\kappa}{\partial \mathbf{u}_{k+\kappa}} \right|_{\mathbf{u}_{k+\kappa}^*} = \mathbf{0}, \quad \text{und} \quad \left. \frac{\partial^2 Q_\kappa}{\partial \mathbf{u}_{k+\kappa}^2} \right|_{\mathbf{u}_{k+\kappa}^*} \succ \mathbf{0} \quad (\text{A.18})$$

sichergestellt ist.

Die Induktionsbehauptung ${}^K Q_{\kappa-1}$ (vgl. (A.15) mit $\kappa \rightarrow \kappa - 1$) wird dann im Induktionsschritt bewiesen. Hierzu wird ${}^K Q_{\kappa-1}$ durch (3.91) ausgedrückt und $\mathbf{u}_{k+\kappa}^*$ aus (A.17) verwendet. Damit folgt

$$\begin{aligned} {}^K Q_{\kappa-1} &= \frac{1}{2} \begin{bmatrix} \mathbf{x}_{k+\kappa-1} \\ \mathbf{x}_{r,k+\kappa-1} \end{bmatrix}^\top (\mathbf{X}^0)^\top \mathbf{Q} \mathbf{X}^0 \begin{bmatrix} \mathbf{x}_{k+\kappa-1} \\ \mathbf{x}_{r,k+\kappa-1} \end{bmatrix} + \frac{1}{2} \mathbf{u}_{k+\kappa-1}^\top \mathbf{R} \mathbf{u}_{k+\kappa-1} \\ &+ \frac{1}{2} \gamma \begin{bmatrix} \mathbf{x}_{k+\kappa} \\ \mathbf{x}_{r,k+\kappa} \end{bmatrix}^\top (\mathbf{X}^0)^\top \mathbf{Q} \mathbf{X}^0 \begin{bmatrix} \mathbf{x}_{k+\kappa} \\ \mathbf{x}_{r,k+\kappa} \end{bmatrix} \\ &+ \frac{1}{2} \gamma \bar{\mathbf{z}}_k^\top \sum_{i=1}^{\eta-1} \left((\mathbf{X}_{\eta+1}^i)^\top \mathbf{Q} \mathbf{X}_{\eta+1}^i + (\mathbf{U}_{\eta+1}^i)^\top \mathbf{R} \mathbf{U}_{\eta+1}^i \right) \bar{\mathbf{z}}_k \end{aligned} \quad (\text{A.19})$$

mit $\bar{\mathbf{z}}_k^\top := [\mathbf{x}_{k+\kappa}^\top \quad \mathbf{x}_{r,k+\kappa+1}^\top \quad \dots \quad \mathbf{x}_{r,k+K}^\top]$. Einsetzen von $\mathbf{x}_{k+\kappa} = \mathbf{A} \mathbf{x}_{k+\kappa-1} + \mathbf{B} \mathbf{u}_{k+\kappa-1}$ (vgl. (3.80)) in (A.19) liefert

$$\begin{aligned} {}^K Q_{\kappa-1} &= \frac{1}{2} \left(\boldsymbol{\rho}_0^{\kappa-1} \mathbf{Q} (\boldsymbol{\rho}_0^{\kappa-1})^\top + \mathbf{u}_{k+\kappa-1}^\top \mathbf{R} \mathbf{u}_{k+\kappa-1} + \gamma \boldsymbol{\rho}_1^{\kappa-1} \mathbf{Q} (\boldsymbol{\rho}_1^{\kappa-1})^\top \right. \\ &\quad \left. + \gamma \sum_{i=1}^{K-(\kappa-1)-2} \left(\boldsymbol{\chi}_i^{\kappa-1} \mathbf{Q} (\boldsymbol{\chi}_i^{\kappa-1})^\top + \boldsymbol{\mu}_i^{\kappa-1} \mathbf{R} (\boldsymbol{\mu}_i^{\kappa-1})^\top \right) \right) \end{aligned} \quad (\text{A.20})$$

und somit die Induktionsbehauptung ((A.15) mit $\kappa \rightarrow \kappa - 1$). Daher gilt (A.15).

Somit ist die analytische Lösung von ${}^K Q_0$ quadratisch bezüglich $\boldsymbol{\rho}_0^\kappa$, \mathbf{u}_k , $\boldsymbol{\rho}_1^\kappa$, $\boldsymbol{\chi}_i^\kappa$ und $\boldsymbol{\mu}_i^\kappa$. Da nach (A.11)–(A.14) jede dieser Komponenten linear bezüglich \mathbf{x}_k , \mathbf{u}_k und $\mathbf{x}_{r,k+1}, \dots, \mathbf{x}_{r,k+n_h}$ ist, folgt Satz 3.4 direkt für $\kappa = 0$ und $K \geq n_h$. Die genauen Werte von \mathbf{H}_K können aus dem im Beweis gegebenen Schema beispielsweise mithilfe von MATLAB berechnet werden. \square

A.2 Beweis zu Lemma 3.8

Beweis:

Nach Lemma 3.7 minimiert die durch $-\mathbf{K}(\mathbf{Z}^{[l]}) [\mathbf{x}_{k+1}^\top \quad \mathbf{x}_{r,k+2}^\top \quad \cdots \quad \mathbf{x}_{r,k+n_h}^\top \quad \mathbf{0}^\top]$ gegebene Stellgröße den Ausdruck $\tilde{\mathbf{y}}_{k+1}^\top \mathbf{Z}^{[l]} \tilde{\mathbf{y}}_{k+1}, \forall l > 0$. Aufgrund von

$$\mathbf{M}(-\mathbf{K}(\mathbf{Z}^{[l]})) \tilde{\mathbf{y}}_k = \tilde{\mathbf{y}}_{k+1} \Big|_{\mathbf{u}_{k+1} = -\mathbf{K}(\mathbf{Z}^{[l]}) [\mathbf{x}_{k+1}^\top \quad \mathbf{x}_{r,k+2}^\top \quad \cdots \quad \mathbf{x}_{r,k+n_h}^\top \quad \mathbf{0}^\top]} \quad (\text{A.21})$$

gilt daher

$$\begin{aligned} & \tilde{\mathbf{y}}_k^\top \mathbf{M}(\mathbf{W}^{[l]})^\top \mathbf{Z}^{[l]} \mathbf{M}(\mathbf{W}^{[l]}) \tilde{\mathbf{y}}_k \\ & \geq \tilde{\mathbf{y}}_k^\top \mathbf{M}(-\mathbf{K}(\mathbf{Z}^{[l]}))^\top \mathbf{Z}^{[l]} \mathbf{M}(-\mathbf{K}(\mathbf{Z}^{[l]})) \tilde{\mathbf{y}}_k. \end{aligned} \quad (\text{A.22})$$

Daraus folgt

$$\mathbf{M}(\mathbf{W}^{[l]})^\top \mathbf{Z}^{[l]} \mathbf{M}(\mathbf{W}^{[l]}) - \mathbf{M}(-\mathbf{K}(\mathbf{Z}^{[l]}))^\top \mathbf{Z}^{[l]} \mathbf{M}(-\mathbf{K}(\mathbf{Z}^{[l]})) \succeq \mathbf{0} \quad (\text{A.23})$$

und somit für die nach (3.129) definierte Folge

$$\mathbf{Z}^{[l+1]} = F(\mathbf{Z}^{[l]}, \mathbf{W}^{[l]}) \succeq F(\mathbf{Z}^{[l]}, -\mathbf{K}(\mathbf{Z}^{[l]})) =: \hat{\mathbf{Z}}^{[l+1]}. \quad (\text{A.24})$$

Daher gilt auch

$$F(\mathbf{H}^{[l]}, -\mathbf{K}(\mathbf{Z}^{[l]})) \succeq F(\mathbf{H}^{[l]}, -\mathbf{K}(\mathbf{H}^{[l]})) = \mathbf{H}^{[l+1]}. \quad (\text{A.25})$$

Mit dem Induktionsanfang $\mathbf{0} \preceq \mathbf{H}^{[0]} \preceq \mathbf{Z}^{[0]}$ und der Induktionsbehauptung $\mathbf{H}^{[l]} \preceq \mathbf{Z}^{[l]}$ ergibt sich

$$\begin{aligned} F(\mathbf{H}^{[l]}, -\mathbf{K}(\mathbf{Z}^{[l]})) &= \mathbf{G} + \gamma \mathbf{M}(-\mathbf{K}(\mathbf{Z}^{[l]}))^\top \mathbf{H}^{[l]} \mathbf{M}(-\mathbf{K}(\mathbf{Z}^{[l]})) \\ &\succeq \mathbf{G} + \gamma \mathbf{M}(-\mathbf{K}(\mathbf{Z}^{[l]}))^\top \mathbf{Z}^{[l]} \mathbf{M}(-\mathbf{K}(\mathbf{Z}^{[l]})) \\ &= F(\mathbf{Z}^{[l]}, -\mathbf{K}(\mathbf{Z}^{[l]})) = \hat{\mathbf{Z}}^{[l+1]}. \end{aligned} \quad (\text{A.26})$$

Zusammen mit (A.25) folgt daraus $\mathbf{H}^{[l+1]} \preceq \hat{\mathbf{Z}}^{[l+1]}$ und unter Einbeziehung von (A.24) ergibt sich schließlich

$$\mathbf{0} \preceq \mathbf{H}^{[l+1]} \preceq \mathbf{Z}^{[l+1]}. \quad (\text{A.27})$$

□

A.3 Beweis zu Lemma 3.9

Beweis:

Der Beweis stellt eine Erweiterung des Beweises von [Lan97, Lemma B.1.2] auf den Solltrajektorienfolgeregelungsfall dar. Sei

$$\mathbf{Z}^{[l+1]} = F\left(\mathbf{Z}^{[l]}, -\tilde{\mathbf{K}}\right) \quad (\text{A.28})$$

mit $\mathbf{Z}^{[0]} = \mathbf{H}^{[0]}$ gegeben, wobei $\tilde{\mathbf{K}}$ so gewählt ist, dass alle Eigenwerte von $(\mathbf{A} - \mathbf{B}\tilde{\mathbf{K}}_x)$ innerhalb des Einheitskreises liegen. Die Existenz von $\tilde{\mathbf{K}}$ ist aufgrund der Steuerbarkeit von (\mathbf{A}, \mathbf{B}) gewährleistet. Mit $\mathbf{W}^{[l]} = -\tilde{\mathbf{K}}$ folgt aus Lemma 3.8, dass

$$\mathbf{0} \preceq \mathbf{H}^{[l]} \preceq \mathbf{Z}^{[l]} \quad (\text{A.29})$$

gilt. Aus

$$\begin{aligned} \mathbf{Z}^{[l+1]} - \mathbf{Z}^{[l]} &= F\left(\mathbf{Z}^{[l]}, -\tilde{\mathbf{K}}\right) - F\left(\mathbf{Z}^{[l-1]}, -\tilde{\mathbf{K}}\right) \\ &= \gamma \mathbf{M}\left(-\tilde{\mathbf{K}}\right)^\top \left(\mathbf{Z}^{[l]} - \mathbf{Z}^{[l-1]}\right) \mathbf{M}\left(-\tilde{\mathbf{K}}\right) \end{aligned} \quad (\text{A.30})$$

folgt¹⁴⁴

$$\text{vec}\left(\mathbf{Z}^{[l+1]} - \mathbf{Z}^{[l]}\right) = \underbrace{\gamma \mathbf{M}\left(-\tilde{\mathbf{K}}\right)^\top \otimes \mathbf{M}\left(-\tilde{\mathbf{K}}\right)}_{=: \mathbf{E}} \text{vec}\left(\mathbf{Z}^{[l]} - \mathbf{Z}^{[l-1]}\right) \quad (\text{A.31})$$

und somit

$$\text{vec}\left(\mathbf{Z}^{[l]} - \mathbf{Z}^{[l-1]}\right) = \mathbf{E}^{l-1} \text{vec}\left(\mathbf{Z}^{[1]} - \mathbf{Z}^{[0]}\right). \quad (\text{A.32})$$

Falls nun alle Eigenwerte von $\sqrt{\gamma} \mathbf{M}\left(-\tilde{\mathbf{K}}\right)$ im Inneren des Einheitskreises liegen, so gilt dies auch für die Eigenwerte von \mathbf{E} . Aufgrund der in (3.121) gegebenen Struktur folgt, dass mindestens $(n_h + 1)n$ Eigenwerte von $\mathbf{M}\left(-\tilde{\mathbf{K}}\right)$ im Koordinatenursprung liegen. Alle übrigen Eigenwerte, d. h. die Eigenwerte von

$$\bar{\mathbf{M}} := \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ -\tilde{\mathbf{K}}_x \mathbf{A} & -\tilde{\mathbf{K}}_x \mathbf{B} \end{bmatrix}, \quad (\text{A.33})$$

werden analog zum Vorgehen in [Lan97, Lemma B.1.2] nachfolgend analysiert. Sei $\|\cdot\|_2$ die Spektralnrm einer Matrix bzw. die Euklidische Norm eines Vektors. Dann gilt

$$\begin{aligned} \lim_{l \rightarrow \infty} \left\| \bar{\mathbf{M}}^l \right\|_2 &= \lim_{l \rightarrow \infty} \left\| \begin{bmatrix} \mathbf{I}_n \\ -\tilde{\mathbf{K}}_x \end{bmatrix} \left(\mathbf{A} - \mathbf{B}\tilde{\mathbf{K}}_x\right)^{l-1} [\mathbf{A} \ \mathbf{B}] \right\|_2 \\ &\leq \lim_{l \rightarrow \infty} \left\| \begin{bmatrix} \mathbf{I}_n \\ -\tilde{\mathbf{K}}_x \end{bmatrix} \right\|_2 \left\| [\mathbf{A} \ \mathbf{B}] \right\|_2 \left\| \left(\mathbf{A} - \mathbf{B}\tilde{\mathbf{K}}_x\right)^{l-1} \right\|_2 = 0. \end{aligned}$$

¹⁴⁴ Der Operator $\text{vec}(\cdot)$ nimmt eine vertikale Konkatenation der Spalten einer Matrix vor und \otimes bezeichnet das Kronecker-Produkt.

Aus $\lim_{l \rightarrow \infty} \bar{M}^l = \mathbf{0}$ folgt, dass alle Eigenwerte von \bar{M} im Inneren des Einheitskreises liegen. Folglich sind alle Eigenwerte von E ebenfalls innerhalb des Einheitskreises und es gilt $e := \|E\|_2 < 1$. Aus

$$\begin{aligned} \text{vec}(\mathbf{Z}^{[j]}) &= \text{vec}(\mathbf{Z}^{[0]}) + \sum_{l=1}^j \text{vec}(\mathbf{Z}^{[l]} - \mathbf{Z}^{[l-1]}) \\ &\stackrel{\text{(A.32)}}{=} \text{vec}(\mathbf{Z}^{[0]}) + \sum_{l=1}^j E^{l-1} \text{vec}(\mathbf{Z}^{[1]} - \mathbf{Z}^{[0]}) \end{aligned} \quad (\text{A.34})$$

folgt

$$\left\| \text{vec}(\mathbf{Z}^{[j]}) \right\|_2 \leq \left\| \text{vec}(\mathbf{Z}^{[0]}) \right\|_2 + \sum_{l=0}^{\infty} \|E\|_2^l \left\| \text{vec}(\mathbf{Z}^{[1]} - \mathbf{Z}^{[0]}) \right\|_2 =: e_0, \quad (\text{A.35})$$

wobei die obere Schranke e_0 unabhängig von j ist. Da $\left\| \text{vec}(\mathbf{Z}^{[j]}) \right\|_2$ durch e_0 nach oben beschränkt ist, existiert e_1 , sodass $\left\| \mathbf{Z}^{[j]} \right\|_2 \leq e_1, \forall j$. Mit $\mathbf{Y} := e_1 \mathbf{I}_{\text{Dim}(\mathbf{H})}$ ergibt sich schließlich

$$\mathbf{0} \preceq \mathbf{H}^{[l]} \preceq \mathbf{Z}^{[l]} \preceq \left\| \mathbf{Z}^{[l]} \right\|_2 \mathbf{I}_{\text{Dim}(\mathbf{H})} \preceq e_1 \mathbf{I}_{\text{Dim}(\mathbf{H})} = \mathbf{Y} \quad (\text{A.36})$$

und somit die Aussage von Lemma 3.9. \square

A.4 Ergänzungen zum linearen Einspurmodell

Die physikalischen Parameter sowie das zugehörige zeitkontinuierliche Modell des in (3.143) gegebenen linearen Einspurmodells sind in [Fla16, Anhang B] zu finden. Die Stellgröße u_k entspricht einem auf das Lenkrad aufgebrachten Drehmoment. Der Systemzustand ist durch

$$\mathbf{x}_k = [\beta_k \quad \psi_{r,k} \quad \psi_k \quad y_k \quad \delta_{\text{LR},v,k} \quad \delta_{\text{LR},k}]^T \quad (\text{A.37})$$

gegeben, wobei β_k der Schwimmwinkel, ψ_k der Gierwinkel, $\psi_{r,k}$ die Gierrate, y_k die laterale Abweichung vom Koordinatenursprung, $\delta_{\text{LR},k}$ der Lenkradwinkel und $\delta_{\text{LR},v,k}$ die Winkelgeschwindigkeit des Lenkrads ist. Die geometrischen Zusammenhänge sind Abbildung A.1 zu entnehmen¹⁴⁵, hierbei stellt v die als konstant angenommene Geschwindigkeit des Fahrzeugs dar. Im vorliegenden Beispiel sei $v = 20 \text{ m s}^{-1}$. Zudem stellt $\delta_s = 0,0625 \delta_{\text{LR}}$ den Lenkwinkel und dessen Zusammenhang zum Lenkradwinkel δ_{LR} dar.

¹⁴⁵ Für eine ausführliche Diskussion des Modells sei auf [Fla16, Anhang B] verwiesen.

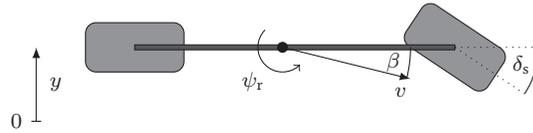


Abbildung A.1: Geometrische Zusammenhänge des linearen Einspurmodells.

A.5 ADP-Solltrajektorienregler für das lineare Einspurmodell mit $\gamma = 1$

Für das Beispiel des linearen Einspurmodells (3.143) mit den in (3.144) gegebenen Gütemaßparametern und $\gamma = 1$ sind die resultierenden Trajektorienverläufe in Abbildung A.2 und die Gewichtsfehler in Abbildung A.3 gegeben.

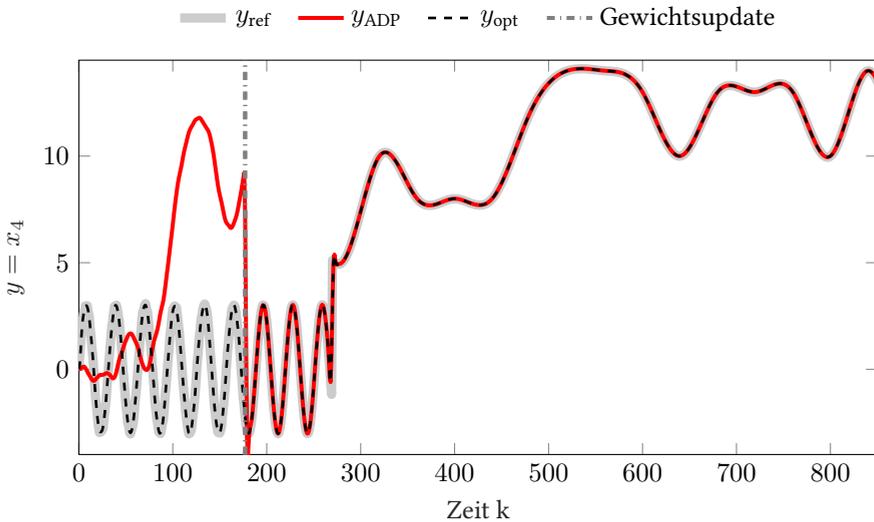


Abbildung A.2: Ergebnis der ADP-Solltrajektorienregelung für System 2 (lineares Einspurmodell sechster Ordnung) für $\gamma = 1$.

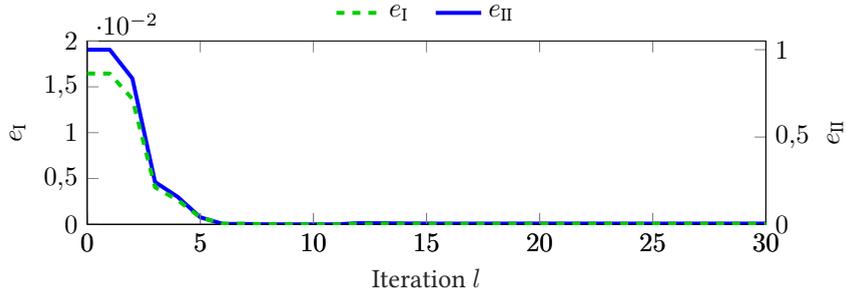


Abbildung A.3: Gewichtsfehlerverlauf während der ersten 30 Iterationen des Lernvorgangs für System 2 (lineares Einspurmodell sechster Ordnung) für $\gamma = 1$. Hierbei stellt e_I (3.148) den Mittelwert und e_{II} (3.149) das Maximum des elementweisen absoluten Fehlers von \hat{w} , jeweils durch $\max_j \{w^*\}_j$ normiert, dar.

B Anhang zu Kapitel 4

B.1 Beweis zu Lemma 4.1

Beweis:

Betrachtet werde $x_{r,i}(t)$ nach (4.7). Sei $J_{i,r}$ die Anzahl verschiedener reeller Eigenwerte sowie $J_{i,c}$ die Anzahl verschiedener konjugiert-komplexer Eigenwertpaare, d. h. $J_i = J_{i,r} + 2J_{i,c}$. Dann folgt mit Annahme 4.1

$$x_{r,i}(t) = \sum_{j=1}^{J_{i,r}} \left(\sum_{k=0}^{\nu_{ij}-1} c_{ijk} t^k \right) e^{\lambda_{ij} t} \quad (\text{B.1})$$

$$+ \sum_{j=J_{i,r}+1}^{J_{i,r}+J_{i,c}} \left(\sum_{k=0}^{\nu_{ij}-1} c_{ijk} t^k \right) e^{\lambda_{ij} t} + \sum_{j=J_{i,r}+J_{i,c}+1}^{J_i} \left(\sum_{k=0}^{\nu_{ij}-1} c_{ijk}^* t^k \right) e^{\lambda_{ij}^* t} \quad (\text{B.2})$$

$$= \text{Re} \left\{ \sum_{j=1}^{J_{i,r}} \left(\sum_{k=0}^{\nu_{ij}-1} c_{ijk} t^k \right) e^{\lambda_{ij} t} \right\} + 2 \text{Re} \left\{ \sum_{j=J_{i,r}+1}^{J_{i,r}+J_{i,c}} \left(\sum_{k=0}^{\nu_{ij}-1} c_{ijk} t^k \right) e^{\lambda_{ij} t} \right\} \quad (\text{B.3})$$

und somit Reellwertigkeit von $x_{r,i}(t)$ und $\mathbf{x}_r(t)$. \square

B.2 Beweis zu Lemma 4.2

Beweis:

Der Beweis erfolgt mittels vollständiger Induktion. Für den Induktionsanfang ergibt sich

$$x_{r,ij}^{(0)}(t) = e^{\lambda_{ij} t} \binom{0}{0} \lambda_{ij}^0 \sum_{k=0}^{\nu_{ij}-1} \frac{k!}{k!} c_{ijk} t^j = e^{\lambda_{ij} t} \sum_{k=0}^{\nu_{ij}-1} c_{ijk} t^j = x_{r,ij}(t). \quad (\text{B.4})$$

Im Induktionsschritt folgt aus der zeitlichen Ableitung des Ausdrucks $x_{r,ij}^{(l)}(t)$ in (4.23)

$$x_{r,ij}^{(l+1)}(t) = \sum_{m=0}^l e^{\lambda_{ij} t} \binom{l}{m} \lambda_{ij}^{l+1-m} \sum_{k=m}^{\nu_{ij}-1} \frac{k!}{(k-m)!} c_{ijk} t^{k-m}$$

$$+ \sum_{m=0}^l e^{\lambda_{ij}t} \binom{l}{m} \lambda_{ij}^{l-m} \sum_{k=m+1}^{\nu_{ij}-1} \frac{k!}{(k-m)!} (k-m) c_{ijk} t^{k-(m+1)} \quad (\text{B.5})$$

$$\begin{aligned} &= e^{\lambda_{ij}t} \binom{l}{0} \lambda_{ij}^{l+1-0} \sum_{k=0}^{\nu_{ij}-1} \frac{k!}{(k-0)!} c_{ijk} t^{k-0} \\ &+ \sum_{m=1}^l e^{\lambda_{ij}t} \binom{l}{m} \lambda_{ij}^{l+1-m} \sum_{k=m}^{\nu_{ij}-1} \frac{k!}{(k-m)!} c_{ijk} t^{k-m} \\ &+ \sum_{m=0}^{l-1} e^{\lambda_{ij}t} \binom{l}{m} \lambda_{ij}^{l-m} \sum_{k=m+1}^{\nu_{ij}-1} \frac{k!}{(k-(m+1))!} c_{ijk} t^{k-(m+1)} \\ &+ e^{\lambda_{ij}t} \binom{l}{l} \lambda_{ij}^{l-l} \sum_{k=l+1}^{\nu_{ij}-1} \frac{k!}{(k-(l+1))!} c_{ijk} t^{k-(l+1)} \quad (\text{B.6}) \end{aligned}$$

$$\begin{aligned} &= e^{\lambda_{ij}t} \binom{l+1}{0} \lambda_{ij}^{l+1-0} \sum_{k=0}^{\nu_{ij}-1} \frac{k!}{(k-0)!} c_{ijk} t^{k-0} \\ &+ \sum_{m=1}^l e^{\lambda_{ij}t} \binom{l}{m} \lambda_{ij}^{l+1-m} \sum_{k=m}^{\nu_{ij}-1} \frac{k!}{(k-m)!} c_{ijk} t^{k-m} \\ &+ \sum_{m=1}^l e^{\lambda_{ij}t} \binom{l}{m-1} \lambda_{ij}^{l+1-m} \sum_{k=m}^{\nu_{ij}-1} \frac{k!}{(k-m)!} c_{ijk} t^{k-m} \dots \\ &\dots + e^{\lambda_{ij}t} \binom{l+1}{l+1} \lambda_{ij}^{l+1-(l+1)} \sum_{k=l+1}^{\nu_{ij}-1} \frac{k!}{(k-(l+1))!} c_{ijk} t^{k-(l+1)} \quad (\text{B.7}) \end{aligned}$$

$$\begin{aligned} &= e^{\lambda_{ij}t} \binom{l+1}{0} \lambda_{ij}^{l+1-0} \sum_{k=0}^{\nu_{ij}-1} \frac{k!}{(k-0)!} c_{ijk} t^{k-0} \\ &+ \sum_{m=1}^l e^{\lambda_{ij}t} \left[\binom{l}{m-1} + \binom{l}{m} \right] \lambda_{ij}^{l+1-m} \sum_{k=m}^{\nu_{ij}-1} \frac{k!}{(k-m)!} c_{ijk} t^{k-m} \\ &+ e^{\lambda_{ij}t} \binom{l+1}{l+1} \lambda_{ij}^{l+1-(l+1)} \sum_{k=l+1}^{\nu_{ij}-1} \frac{k!}{(k-(l+1))!} c_{ijk} t^{k-(l+1)} \quad (\text{B.8}) \end{aligned}$$

$$= \sum_{m=0}^{l+1} e^{\lambda_{ij}t} \binom{l+1}{m} \lambda_{ij}^{l+1-m} \sum_{k=m}^{\nu_{ij}-1} \frac{k!}{(k-m)!} c_{ijk} t^{k-m}, \quad (\text{B.9})$$

wobei im letzten Schritt die Beziehung $\binom{l}{m-1} + \binom{l}{m} = \binom{l+1}{m}$ verwendet wird. Der Ausdruck (B.9) entspricht der Induktionsbehauptung, d. h. (4.23) mit $l+1$. \square

B.3 Zustandstransformation für reellwertige $\zeta(t)$ und D

Für jedes konjugiert-komplexe Polpaar λ_{ij} und $\lambda_{ij}^* = \lambda_{ij}^*$ lassen sich nach [Büh20, S. 31 f.] durch eine bijektive Zustandstransformation stets reellwertige Ersatzzustände $\zeta_1(t)$ und $\zeta_2(t)$ mit reeller Dynamik definieren, die den Sollzustandsverlauf nicht beeinflussen. Durch Einsetzen von $\lambda_{ij}^* = \lambda_{ij}^*$ in (4.9) ist ersichtlich, dass $D_{ij^*} = D_{ij}^*$ gilt (vgl. (4.14)). Zudem folgt aus (4.7) unter Annahme 4.1 mit (4.12) $\zeta_{ij^*}(t) = \zeta_{ij}^*(t)$.

Definiere

$$\zeta_1(t) := \zeta_{ij}(t) + \zeta_{ij^*}(t) = 2 \operatorname{Re}\{\zeta_{ij}(t)\} \quad (\text{B.10})$$

und

$$\zeta_2(t) := j(\zeta_{ij}(t) - \zeta_{ij^*}(t)) = 2 \operatorname{Im}\{\zeta_{ij}(t)\}. \quad (\text{B.11})$$

Da durch $\zeta_{ij}(t) = \frac{\zeta_1(t) - j\zeta_2(t)}{2}$ und $\zeta_{ij^*}(t) = \frac{\zeta_1(t) + j\zeta_2(t)}{2}$ Umkehrfunktionen existieren, gewährleistet eine Ersetzung von $\zeta_{ij}(t)$ und $\zeta_{ij^*}(t)$ durch $\zeta_1(t)$ und $\zeta_2(t)$ aufgrund von Proposition 4.2 weiterhin eine eindeutige Repräsentation der Parameter c . Die reellwertigen Zustände $\zeta_1(t)$ und $\zeta_2(t)$ sind durch die reelle Dynamik

$$\frac{d}{dt} \begin{bmatrix} \zeta_1(t) \\ \zeta_2(t) \end{bmatrix} = \begin{bmatrix} \dot{\zeta}_{ij}(t) + \dot{\zeta}_{ij^*}(t) \\ j(\dot{\zeta}_{ij}(t) - \dot{\zeta}_{ij^*}(t)) \end{bmatrix} = \begin{bmatrix} D_{ij}\zeta_{ij}(t) + D_{ij}^*\zeta_{ij}^*(t) \\ j(D_{ij}\zeta_{ij}(t) - D_{ij}^*\zeta_{ij}^*(t)) \end{bmatrix} \quad (\text{B.12})$$

$$= \begin{bmatrix} D_{ij} \left(\frac{\zeta_1(t) - j\zeta_2(t)}{2} \right) + D_{ij}^* \left(\frac{\zeta_1(t) + j\zeta_2(t)}{2} \right) \\ j \left(D_{ij} \left(\frac{\zeta_1(t) - j\zeta_2(t)}{2} \right) - D_{ij}^* \left(\frac{\zeta_1(t) + j\zeta_2(t)}{2} \right) \right) \end{bmatrix} \quad (\text{B.13})$$

$$= \begin{bmatrix} \frac{D_{ij} + D_{ij}^*}{2} & j \frac{D_{ij}^* - D_{ij}}{2} \\ j \frac{D_{ij} - D_{ij}^*}{2} & \frac{D_{ij} + D_{ij}^*}{2} \end{bmatrix} \begin{bmatrix} \zeta_1(t) \\ \zeta_2(t) \end{bmatrix} \quad (\text{B.14})$$

verknüpft. Schließlich müssen noch s_{ij} und s_{ij^*} (vgl. (4.18)) durch $s_1 = s_{ij}$ und $s_2 = 0$ ersetzt werden, damit

$$s_1^T \zeta_1(t) + s_2^T \zeta_2(t) = \zeta_{ij}(t) + \zeta_{ij}^*(t) = s_{ij}^T \zeta_{ij}(t) + s_{ij^*}^T \zeta_{ij^*}(t) \quad (\text{B.15})$$

gilt und somit $x_{r,i}(t)$ nicht verändert wird.

B.4 Äquivalenz zwischen Problem 4.1 und Problem 4.2

Lemma B.1 (vgl. [Büh20, S. 33 f.])

Problem 4.1 und Problem 4.2 sind äquivalent.

Beweis:

Unter Annahme eines linearen Regelgesetzes der Form $\boldsymbol{\mu}(\tilde{\boldsymbol{x}}(\tau)) = -\tilde{\boldsymbol{K}}\tilde{\boldsymbol{x}}(\tau)$ ergibt sich der Verlauf des erweiterten Zustands (vgl. (4.33)) nach [Föl16, S. 277] zu

$$\tilde{\boldsymbol{x}}(\tau) = e^{(\tilde{\boldsymbol{A}} - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}})(\tau-t)} \tilde{\boldsymbol{x}}(t). \quad (\text{B.16})$$

Einsetzen in das Gütefunktional (4.34) ergibt

$$\int_t^\infty e^{-\gamma(\tau-t)} \left(\tilde{\boldsymbol{x}}^\top(\tau) \tilde{\boldsymbol{Q}} \tilde{\boldsymbol{x}}(\tau) + \boldsymbol{\mu}^\top(\boldsymbol{x}(\tau), \zeta(\tau)) \boldsymbol{R} \boldsymbol{\mu}(\boldsymbol{x}(\tau), \zeta(\tau)) \right) d\tau \quad (\text{B.17})$$

$$= \tilde{\boldsymbol{x}}^\top(t) \int_t^\infty e^{-\frac{\gamma}{2}(\tau-t)} e^{(\tilde{\boldsymbol{A}} - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}})^\top(\tau-t)} \left(\tilde{\boldsymbol{Q}} + \tilde{\boldsymbol{K}}^\top \boldsymbol{R} \tilde{\boldsymbol{K}} \right) e^{(\tilde{\boldsymbol{A}} - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}})(\tau-t)} e^{-\frac{\gamma}{2}(\tau-t)} d\tau \tilde{\boldsymbol{x}}(t) \quad (\text{B.18})$$

$$= \tilde{\boldsymbol{x}}^\top(t) \int_t^\infty e^{((\tilde{\boldsymbol{A}} - \frac{\gamma}{2}) - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}})^\top(\tau-t)} \left(\tilde{\boldsymbol{Q}} + \tilde{\boldsymbol{K}}^\top \boldsymbol{R} \tilde{\boldsymbol{K}} \right) e^{((\tilde{\boldsymbol{A}} - \frac{\gamma}{2}) - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}})(\tau-t)} d\tau \tilde{\boldsymbol{x}}(t) \quad (\text{B.19})$$

$$= \int_t^\infty \left(e^{((\tilde{\boldsymbol{A}} - \frac{\gamma}{2}) - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}})^\top(\tau-t)} \tilde{\boldsymbol{x}}(t) \right)^\top \left(\tilde{\boldsymbol{Q}} + \tilde{\boldsymbol{K}}^\top \boldsymbol{R} \tilde{\boldsymbol{K}} \right) e^{((\tilde{\boldsymbol{A}} - \frac{\gamma}{2}) - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}})(\tau-t)} \tilde{\boldsymbol{x}}(t) d\tau. \quad (\text{B.20})$$

In (B.18) wurden dabei die Zusammenhänge

$$e^{-\frac{\gamma}{2}(\tau-t)} e^{(\tilde{\boldsymbol{A}} - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}})^\top(\tau-t)} = e^{(\tilde{\boldsymbol{A}} - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}} - \frac{\gamma}{2}\boldsymbol{I})^\top(\tau-t)} \quad (\text{B.21})$$

und

$$e^{(\tilde{\boldsymbol{A}} - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}})(\tau-t)} e^{-\frac{\gamma}{2}(\tau-t)} = e^{(\tilde{\boldsymbol{A}} - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}} - \frac{\gamma}{2}\boldsymbol{I})(\tau-t)} \quad (\text{B.22})$$

genutzt. Der Ausdruck $e^{((\tilde{\boldsymbol{A}} - \frac{\gamma}{2}) - \tilde{\boldsymbol{B}}\tilde{\boldsymbol{K}})^\top(\tau-t)} \tilde{\boldsymbol{x}}(t)$ in (B.20) entspricht dem Verlauf des erweiterten Zustands für die Systemdynamik $\dot{\tilde{\boldsymbol{x}}}(\tau) = \left(\tilde{\boldsymbol{A}} - \frac{\gamma}{2}\boldsymbol{I} \right) \tilde{\boldsymbol{x}}(\tau) + \tilde{\boldsymbol{B}}\boldsymbol{\mu}(\boldsymbol{x}(\tau), \zeta(\tau))$ für den durch $\boldsymbol{\mu}(\tilde{\boldsymbol{x}}(\tau)) = -\tilde{\boldsymbol{K}}\tilde{\boldsymbol{x}}(\tau)$ geschlossenen Regelkreis. Somit folgt Äquivalenz zwischen Problem 4.1 und Problem 4.2. Zudem rechtfertigt die LQ-Form von Problem 4.2 die Annahme eines linearen Regelgesetzes (vgl. [Büh20, S. 33 f.], [ML14a, Beweis zu Lemma 2]). \square

B.5 Beweis zu Lemma 4.3

Beweis:

Es gilt

$$\text{Rang}(\boldsymbol{M}) = \text{Dim}(\text{Bild}(\boldsymbol{M})) = n_1 + n_2 \quad (\text{B.23})$$

$$\Leftrightarrow \forall \boldsymbol{z} = \begin{bmatrix} \boldsymbol{z}_1 \\ \boldsymbol{z}_2 \end{bmatrix} \in \mathbb{R}^{n_1+n_2} \exists \boldsymbol{v} = \begin{bmatrix} \boldsymbol{v}_1 \\ \boldsymbol{v}_2 \end{bmatrix} : \boldsymbol{M}\boldsymbol{v} = \boldsymbol{z}. \quad (\text{B.24})$$

Sei ein solches beliebiges z gegeben. Aus

$$\begin{bmatrix} M_1 & M_3 \\ \mathbf{0} & M_2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \quad (\text{B.25})$$

folgt

$$M_2 v_2 = z_2. \quad (\text{B.26})$$

Aufgrund von (4.53) existiert v_2 . Des Weiteren gilt

$$M_1 v_1 = z_1 - M_3 v_2. \quad (\text{B.27})$$

Wegen (4.52) existiert auch v_1 . Daher existiert v , es folgt (B.24) und somit gilt (4.55). \square

B.6 Beweis zu Lemma 4.4

Beweis:

Mit dem Hautus-Kriterium für Stabilisierbarkeit [ZD98, Theorem 3.2] folgt, dass (\tilde{A}', \tilde{B}) genau dann stabilisierbar ist, wenn $[\tilde{A}' - \lambda I \quad \tilde{B}]$ für alle nicht-negativen Eigenwerte λ von \tilde{A}' Maximalrang aufweist. Es gilt

$$\begin{aligned} \text{Rang} \left[\tilde{A}' - \lambda I \quad \tilde{B} \right] &= \text{Rang} \begin{bmatrix} A - \lambda I & \mathbf{0} & B \\ \mathbf{0} & (D - \frac{\gamma}{2} I) - \lambda I & \mathbf{0} \end{bmatrix} \\ &= \text{Rang} \left[\begin{array}{cc|c} A - \lambda I & B & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} & (D - \frac{\gamma}{2} I) - \lambda I \end{array} \right]. \end{aligned} \quad (\text{B.28})$$

Da (A, B) aufgrund von Annahme 4.2 stabilisierbar ist, hat nach dem Hautus-Kriterium $[A - \lambda I \quad B]$ vollen Rang. Zudem ist $(D - \frac{\gamma}{2} I) - \lambda I$ regulär, da $D - \frac{\gamma}{2} I$ für $\gamma > \gamma_{\min}$ ausschließlich negative Eigenwerte besitzt. Nach Lemma 4.3 weist der Ausdruck in (B.28) somit vollen Rang auf. Somit sind alle instabilen Eigenwerte steuerbar, woraus Stabilisierbarkeit von (\tilde{A}', \tilde{B}) resultiert. \square

B.7 Beweis zu Lemma 4.6

Beweis:

Mit dem Hautus-Kriterium für Detektierbarkeit [ZD98, Theorem 3.4] folgt, dass $(\tilde{A}', \sqrt{\tilde{Q}})$

genau dann detektierbar ist, wenn $\begin{bmatrix} \tilde{\mathbf{A}}' - \lambda \mathbf{I} \\ \sqrt{\tilde{\mathbf{Q}}} \end{bmatrix}$ für alle nicht-negativen Eigenwerte λ von $\tilde{\mathbf{A}}'$ vollen Rang aufweist. Aus

$$\tilde{\mathbf{Q}} = [\mathbf{I} \quad -\mathbf{S}]^\top \mathbf{Q} [\mathbf{I} \quad -\mathbf{S}] = [\mathbf{I} \quad -\mathbf{S}]^\top \sqrt{\mathbf{Q}}^\top \sqrt{\mathbf{Q}} [\mathbf{I} \quad -\mathbf{S}] \quad (\text{B.29})$$

folgt mit der Definition von $\tilde{\mathbf{Q}}$ nach (4.36)

$$\sqrt{\tilde{\mathbf{Q}}} = \sqrt{\mathbf{Q}} [\mathbf{I} \quad -\mathbf{S}]. \quad (\text{B.30})$$

Daher gilt

$$\begin{aligned} \text{Rang} \begin{bmatrix} \tilde{\mathbf{A}}' - \lambda \mathbf{I} \\ \sqrt{\tilde{\mathbf{Q}}} \end{bmatrix} &= \text{Rang} \begin{bmatrix} \mathbf{A} - \lambda \mathbf{I} & \mathbf{0} \\ \mathbf{0} & (\mathbf{D} - \frac{\gamma}{2} \mathbf{I}) - \lambda \mathbf{I} \\ \sqrt{\mathbf{Q}} & -\sqrt{\mathbf{Q}} \mathbf{S} \end{bmatrix} \\ &= \text{Rang} \begin{bmatrix} \mathbf{A} - \lambda \mathbf{I} & \mathbf{0} \\ \sqrt{\mathbf{Q}} & -\sqrt{\mathbf{Q}} \mathbf{S} \\ \mathbf{0} & (\mathbf{D} - \frac{\gamma}{2} \mathbf{I}) - \lambda \mathbf{I} \end{bmatrix}^\top. \end{aligned} \quad (\text{B.31})$$

Unter Annahme 4.2 ist $(\mathbf{A}, \sqrt{\mathbf{Q}})$ detektierbar, daher besitzt $\begin{bmatrix} \mathbf{A} - \lambda \mathbf{I} \\ \sqrt{\mathbf{Q}} \end{bmatrix}$ nach dem Hautuskriterium vollen Rang. Zudem ist $(\mathbf{D} - \frac{\gamma}{2} \mathbf{I}) - \lambda \mathbf{I}$ regulär, da $\mathbf{D} - \frac{\gamma}{2} \mathbf{I}$ für $\gamma > \gamma_{\min}$ ausschließlich negative Eigenwerte besitzt. Mit Lemma 4.5 folgt somit, dass (B.31) Maximalrang aufweist. Daher sind alle nicht-negativen Eigenwerte beobachtbar und $(\tilde{\mathbf{A}}', \sqrt{\tilde{\mathbf{Q}}})$ ist detektierbar. \square

B.8 Value Iteration nach Bian und Jiang [BJ16a]

Einsetzen von (4.73) und (4.74) in (4.72) ergibt

$$\frac{\partial V^{\mu_s}(\tilde{\mathbf{x}})}{\partial s} = \phi^\top(\tilde{\mathbf{x}}) \frac{d\hat{\mathbf{w}}_\phi(s)}{ds} = \min_{\boldsymbol{\mu}} (r(\tilde{\mathbf{x}}, \boldsymbol{\mu}) + \psi^\top(\tilde{\mathbf{x}}, \boldsymbol{\mu}) \hat{\mathbf{w}}_\psi(s)). \quad (\text{B.32})$$

Analog zu [BJ16a] wird ein Datensatz aus $M \geq \frac{\tilde{n}(\tilde{n}+1)}{2} + \tilde{n}p$ Tupeln $(\tilde{\mathbf{x}}(t_j), \Theta_j, \tilde{\mathbf{x}}(t_j + T_{\text{IRL}}))$, $j = 1, \dots, M$, verwendet, wobei

$$\Theta_j := \int_{t_j}^{t_j + T_{\text{IRL}}} \psi(\tilde{\mathbf{x}}, \boldsymbol{\mu}) d\tau \quad (\text{B.33})$$

definiert ist. Mithilfe der Messdatenmatrizen

$$\mathbf{D}_\phi := \frac{1}{M} \sum_{j=1}^M \phi(\tilde{\mathbf{x}}_j) \phi^\top(\tilde{\mathbf{x}}_j), \quad (\text{B.34})$$

$$\mathbf{D}_\psi := \frac{1}{M} \sum_{j=1}^M \Theta_j \Theta_j^\top \quad (\text{B.35})$$

folgt, da die skalare Gleichung (B.32) für jedes der M Tupel gilt, nach Multiplikation von $\phi(\tilde{\mathbf{x}}_j)$ von links und Summation über alle M Datentupel

$$\sum_{j=1}^M \phi(\tilde{\mathbf{x}}_j) \phi^\top(\tilde{\mathbf{x}}_j) \frac{d\hat{\mathbf{w}}_\phi(s)}{ds} = \sum_{j=1}^M \phi(\tilde{\mathbf{x}}_j) \min_{\boldsymbol{\mu}} (r(\tilde{\mathbf{x}}_j, \boldsymbol{\mu}) + \boldsymbol{\psi}^\top(\tilde{\mathbf{x}}_j, \boldsymbol{\mu}) \hat{\mathbf{w}}_\psi(s)). \quad (\text{B.36})$$

Falls \mathbf{D}_ϕ invertierbar ist, resultiert somit die Adaptionvorschrift

$$\frac{d\hat{\mathbf{w}}_\phi(s)}{ds} = \frac{1}{M} \mathbf{D}_\phi^{-1} \sum_{j=1}^M \phi(\tilde{\mathbf{x}}_j) \min_{\boldsymbol{\mu}} (r(\tilde{\mathbf{x}}_j, \boldsymbol{\mu}) + \boldsymbol{\psi}^\top(\tilde{\mathbf{x}}_j, \boldsymbol{\mu}) \hat{\mathbf{w}}_\psi(s)) \quad (\text{B.37})$$

der Gewichte $\hat{\mathbf{w}}_\phi(s)$. Aus $\hat{\mathbf{w}}_\phi(s)$ lässt sich, falls \mathbf{D}_ψ invertierbar ist, wegen¹⁴⁶

$$\Theta_j^\top \hat{\mathbf{w}}_\psi(s) = \left(\phi(\tilde{\mathbf{x}}(t_j + T_{\text{IRL}})) - \phi(\tilde{\mathbf{x}}(t_j)) \right)^\top \hat{\mathbf{w}}_\phi(s) \quad (\text{B.38})$$

direkt

$$\hat{\mathbf{w}}_\psi(s) = \Theta \hat{\mathbf{w}}_\phi(s) \quad (\text{B.39})$$

mit

$$\Theta := \frac{1}{M} \mathbf{D}_\psi^{-1} \sum_{j=1}^M \left(\phi(\tilde{\mathbf{x}}(t_j + T_{\text{IRL}})) - \phi(\tilde{\mathbf{x}}(t_j)) \right)^\top \quad (\text{B.40})$$

berechnen. Die Konvergenz dieser Value Iteration beginnend mit einer positiv definiten initialen Value Function wird in [BJ16a, Theorem 2] untersucht.

¹⁴⁶ Gleichung (B.38) folgt direkt aus $V^{\boldsymbol{\mu}_s}(\tilde{\mathbf{x}}(t_j + T_{\text{IRL}})) - V^{\boldsymbol{\mu}_s}(\tilde{\mathbf{x}}(t_j)) = \int_{t_j}^{t_j + T_{\text{IRL}}} \nabla_{\tilde{\mathbf{x}}} V^{\boldsymbol{\mu}_s}(\tilde{\mathbf{x}}) \frac{d\tilde{\mathbf{x}}}{d\tau} d\tau$.

C Anhang zu Kapitel 5

C.1 Beweisskizze zu Lemma 5.2

Beweis:

Einsetzen eines Elements $x_{\text{PE},o}$ von \mathbf{x}_{PE} (5.30) in einen der Faktoren $x_o^{f_{o,\bar{h}}}$ von $\phi_{\bar{h}}(\mathbf{x})$ (5.25) und Anwendung des Multinomialtheorems [Spi68, S. 4] liefert

$$\begin{aligned} x_{\text{PE},o}^{f_{o,\bar{h}}} &= \left(\sum_{j=1}^m g_{j,o} \sin(\omega_j t) + \sum_{j=1}^m \bar{g}_{j,o} \cos(\omega_j t) \right)^{f_{o,\bar{h}}} \\ &= \sum_{r_1 + \dots + r_{2m} = f_{o,\bar{h}}} \frac{f_{o,\bar{h}}!}{r_1! \dots r_{2m}!} g_{1,o}^{r_1} \sin^{r_1}(\omega_1 t) \dots \bar{g}_{m,o}^{r_{2m}} \cos^{r_{2m}}(\omega_m t). \end{aligned} \quad (\text{C.1})$$

Unter der beispielhaften Annahme, dass r_1 geradzahlig ist, ergibt sich aus der Potenzfunktion¹⁴⁷ von $\sin^{r_1}(\omega_1 t)$ [Spi68, S. 17] aus (C.1)

$$\begin{aligned} x_{\text{PE},o}^{f_{o,\bar{h}}} &= \sum_{r_1 + \dots + r_{2m} = f_{o,\bar{h}}} \frac{f_{o,\bar{h}}!}{r_1! \dots r_{2m}!} g_{1,o}^{r_1} \dots \bar{g}_{m,o}^{r_{2m}} \\ &\quad \left(\frac{1}{2^{r_1}} \binom{r_1}{\frac{r_1}{2}} + \frac{1}{2^{r_1-1}} \sum_{s=0}^{\frac{r_1}{2}-1} (-1)^{\frac{r_1}{2}-s} \binom{r_1}{s} \cos((r_1 - 2s)\omega_1 t) \right) \\ &\quad \dots \cos^{r_{2m}}(\omega_m t). \end{aligned} \quad (\text{C.2})$$

Da für den ungeradzahigen Fall sowie die Kosinusfunktionen ähnliche Potenzfunktionen existieren (vgl. [Spi68, S. 17]), kann (C.2) unter zusätzlicher Anwendung der Produktregeln¹⁴⁸ für $\sin(\cdot) \sin(\cdot)$, $\sin(\cdot) \cos(\cdot)$ und $\cos(\cdot) \cos(\cdot)$ zu

$$x_{\text{PE},o}^{f_{o,\bar{h}}} = \sum_{l=1}^{L_{\bar{h}}^{(o)}} a_{l,\bar{h}}^{(o)} \sin \left(\sum_{j=1}^m b_{j,l,\bar{h}}^{(o)} \omega_j t \right) + \sum_{k=1}^{K_{\bar{h}}^{(o)}} c_{k,\bar{h}}^{(o)} \cos \left(\sum_{j=1}^m d_{j,k,\bar{h}}^{(o)} \omega_j t \right) + e_{\bar{h}}^{(o)} \quad (\text{C.3})$$

¹⁴⁷ Auch bekannt als *power reduction formula*.

¹⁴⁸ Bekannt als *product-to-sum identities*.

umgeformt werden. Die Indizes \bar{h} und $\langle o \rangle$, mit denen die Parameter $a_{l,\bar{h}}^{\langle o \rangle}, c_{k,\bar{h}}^{\langle o \rangle} \in \mathbb{R}_{\neq 0}$, $e_{\bar{h}}^{\langle o \rangle}, b_{j,l,\bar{h}}^{\langle o \rangle}, d_{j,k,\bar{h}}^{\langle o \rangle} \in \mathbb{R}$ in (C.3) gekennzeichnet sind, indizieren die Abhängigkeit von den Elementen $\phi_{\bar{h}}(\mathbf{x}_{\text{PE}})$ bzw. dem Faktor $x_{\text{PE},o}^{f_{o,\bar{h}}}$. Um $a_{l,\bar{h}}^{\langle o \rangle}, c_{k,\bar{h}}^{\langle o \rangle} \neq 0$ annehmen zu können, werden die oberen Summengrenzen $L_{\bar{h}}^{\langle o \rangle}, K_{\bar{h}}^{\langle o \rangle} \in \mathbb{N}_{\geq 0}$ mit denselben Indizes versehen.

Da (C.3) $\forall o \in \{1, \dots, n\}$ gilt, folgt für jedes Element $\phi_{\bar{h}}(\mathbf{x}_{\text{PE}})$ von $\phi(\mathbf{x}_{\text{PE}})$

$$\begin{aligned} \phi_{\bar{h}}(\mathbf{x}_{\text{PE}}) &= \bar{f}_{\bar{h}} \prod_{o=1}^n x_{\text{PE},o}^{f_{o,\bar{h}}} \\ &= \sum_{l=1}^{L_{\bar{h}}} a_{l,\bar{h}} \sin \left(\sum_{j=1}^m b_{j,l,\bar{h}} \omega_j t \right) + \sum_{k=1}^{K_{\bar{h}}} c_{k,\bar{h}} \cos \left(\sum_{j=1}^m d_{j,k,\bar{h}} \omega_j t \right) + e_{\bar{h}}. \end{aligned} \quad (\text{C.4})$$

Die Umformung in (C.4) ergibt sich aus der Ausmultiplikation des Produktes der Summen (C.3) und erneuter Anwendung der Produktregeln für $\sin(\cdot) \sin(\cdot)$, $\sin(\cdot) \cos(\cdot)$ und $\cos(\cdot) \cos(\cdot)$. Aufgrund von $f_{o,\bar{h}} \neq 0$ für mindestens ein o (vgl. Annahme 5.3) sowie $m \geq 1$ und $g_{j,o} \neq 0$ oder $\bar{g}_{j,o} \neq 0$ für mindestens ein j (vgl. Annahme 5.4) folgt, dass die Summenobergrenzen $L_{\bar{h}} + K_{\bar{h}} \geq 1, \forall \bar{h} \in \{1, \dots, h\}$, erfüllen. Mit der Kurzschreibweise $\underline{\omega}_{l,\bar{h}} := \sum_{j=1}^m b_{j,l,\bar{h}} \omega_j$ und $\underline{\omega}_{k,\bar{h}} := \sum_{j=1}^m d_{j,k,\bar{h}} \omega_j$ folgt aus (C.4) schließlich (5.31). \square

C.2 Beweis zu Lemma 5.9

Beweis:

Einsetzen von (5.53) in (5.25) liefert $\forall \bar{h} \in \{1, \dots, h\}$

$$\phi_{\bar{h}}(\bar{\mathbf{x}}_{\text{PE}}) = \bar{f}_{\bar{h}} \prod_{j=1}^n \bar{x}_{\text{PE},j}^{f_{j,\bar{h}}} = \left(\prod_{j=1}^n \nu_j^{f_{j,\bar{h}}} \right) \left(\bar{f}_{\bar{h}} \prod_{j=1}^n x_{\text{PE},j}^{f_{j,\bar{h}}} \right) = \left(\prod_{j=1}^n \nu_j^{f_{j,\bar{h}}} \right) \phi_{\bar{h}}(\mathbf{x}_{\text{PE}}). \quad (\text{C.5})$$

Nach (C.5) beeinflusst eine Skalierung von \mathbf{x}_{PE} wie in (5.53) lediglich die Koeffizienten $a_{l,\bar{h}}$, $c_{k,\bar{h}}$ und $e_{\bar{h}}$ in (5.31). Da die exakten Werte dieser Koeffizienten jedoch irrelevant für die Berechnung der Menge Ω sowie die Erfüllung von (5.35) sind, gewährleistet jeder Vektor $\omega \in \Omega$, dass $\bar{\mathbf{x}}_{\text{PE}}$ SR bezüglich $\dot{\phi}(\cdot)$ ist, wenn Ω basierend auf Satz 5.2 mit \mathbf{x}_{PE} berechnet wurde. \square

C.3 M_i und $\mathbf{v}_{\text{freq}}(t)$ für das verwendete Beispielsystem

Im Folgenden wird gezeigt, dass das in Abschnitt 5.7 betrachtete Beispielsystem die durch (5.79) gegebene Bedingung in Annahme 5.6 sowie $\text{Rang}(M_i) = h_i$ erfüllt. Da \tilde{x}_1 , mit $\tilde{\mathbf{x}}$ wie

in (5.78), einem flachen Ausgang $h(\tilde{\mathbf{x}}) = \tilde{x}_1$ des Systems (5.85) entspricht, folgt für die Wahl von $y_f|_{\tilde{\mathbf{x}}_{\text{PE}}} = h_f(\mathbf{z}_{\text{B}}|_{\tilde{\mathbf{x}}_{\text{PE}}}) = h_f(\mathbf{t}(\tilde{\mathbf{x}}_{\text{PE}})) = \tilde{x}_{\text{PE},1}$ mit (5.66) und (5.88)

$$\tilde{x}_1 = \tilde{x}_{\text{PE},1} = \nu_1 (\sin(\omega_{c,1}t) + \sin(\omega_{c,2}t)). \quad (\text{C.6})$$

Nach (5.72) gilt $z_{\text{B},2} = \dot{z}_{\text{B},1}$ und aus (5.74) folgt $z_{\text{B},1} = t_i(\tilde{\mathbf{x}}) = h(\tilde{\mathbf{x}}) = \tilde{x}_1$. Daraus ergibt sich

$$z_{\text{B},2} = \dot{z}_{\text{B},1} = \frac{d}{dt} \tilde{x}_1 = \nu_1 (\omega_{c,1} \cos(\omega_{c,1}t) + \omega_{c,2} \cos(\omega_{c,2}t)). \quad (\text{C.7})$$

Aufgrund von

$$z_{\text{B},2} = L_{\mathbf{f}_g} h(\tilde{\mathbf{x}}) = -2\tilde{x}_1 + \tilde{x}_2 \quad (\text{C.8})$$

(vgl. (5.74)) folgt mit (C.7)

$$\begin{aligned} \tilde{x}_2 &= 2\tilde{x}_1 + z_{\text{B},2} = 2\tilde{x}_{\text{PE},1} + z_{\text{B},2} \\ &= \nu_1 (2 \sin(\omega_{c,1}t) + 2 \sin(\omega_{c,2}t) + \omega_{c,1} \cos(\omega_{c,1}t) + \omega_{c,2} \cos(\omega_{c,2}t)). \end{aligned} \quad (\text{C.9})$$

Wie (C.6) und (C.9) verdeutlichen, gilt für das gewählte Beispiel $\tilde{\mathbf{x}} \neq \tilde{\mathbf{x}}_{\text{PE}}$. Um zu zeigen, dass (5.79) dennoch gilt, wird $\tilde{\boldsymbol{\sigma}}_i$ berechnet. Es folgt für $i \in \mathcal{N}$

$$\begin{aligned} \tilde{\boldsymbol{\sigma}}_i &= \dot{\boldsymbol{\phi}}_i(\tilde{\mathbf{x}}) \\ &= \nu_1^2 \underbrace{\begin{bmatrix} b_1 & 0 & b_2 & 0 & b_3 & 0 & b_4 & 0 \\ b_5 & a_1 & b_6 & a_2 & b_7 & a_3 & b_8 & a_4 \\ b_9 & a_5 & b_{10} & a_6 & b_{11} & a_7 & b_{12} & a_8 \end{bmatrix}}_{=: M_i} \underbrace{\begin{bmatrix} \sin((\omega_{c,1} - \omega_{c,2})t) \\ \cos((\omega_{c,1} - \omega_{c,2})t) \\ \sin((\omega_{c,1} + \omega_{c,2})t) \\ \cos((\omega_{c,1} + \omega_{c,2})t) \\ \sin(2\omega_{c,1}t) \\ \cos(2\omega_{c,1}t) \\ \sin(2\omega_{c,2}t) \\ \cos(2\omega_{c,2}t) \end{bmatrix}}_{=: \frac{d\mathbf{v}_{\text{freq}}(t)}{dt}} \end{aligned} \quad (\text{C.10})$$

mit

$$\begin{aligned} a_1 &= -\frac{(\omega_{c,1} - \omega_{c,2})^2}{2}, & a_2 &= \frac{(\omega_{c,1} + \omega_{c,2})^2}{2}, \\ a_3 &= \omega_{c,1}^2, & a_4 &= \omega_{c,2}^2, \\ a_5 &= -2(\omega_{c,1} - \omega_{c,2})^2, & a_6 &= 2(\omega_{c,1} + \omega_{c,2})^2, \\ a_7 &= 4\omega_{c,1}^2, & a_8 &= 4\omega_{c,2}^2, \\ b_1 &= \omega_{c,2} - \omega_{c,1}, & b_2 &= \omega_{c,1} + \omega_{c,2}, \\ b_3 &= \omega_{c,1}, & b_4 &= \omega_{c,2}, \\ b_5 &= 2(\omega_{c,2} - \omega_{c,1}), & b_6 &= 2(\omega_{c,1} + \omega_{c,2}), \\ b_7 &= 2\omega_{c,1}, & b_8 &= 2\omega_{c,2}, \\ b_9 &= (\omega_{c,1}\omega_{c,2} + 4)(\omega_{c,2} - \omega_{c,1}), & b_{10} &= (4 - \omega_{c,1}\omega_{c,2})(\omega_{c,1} + \omega_{c,2}), \\ b_{11} &= \omega_{c,1}(4 - \omega_{c,1}^2), & b_{12} &= \omega_{c,2}(4 - \omega_{c,2}^2). \end{aligned} \quad (\text{C.11})$$

Weiterhin gilt $\text{Rang}(M_i) = h_i = 3, \forall \omega_c \in \Omega$. Schließlich muss noch gezeigt werden, dass die Frequenzen gleichartiger trigonometrischer Funktionen in $v_{\text{freq}}(t)$ unterschiedlich sind. Dies ist aufgrund von $\omega_{c,1} \neq 0, \omega_{c,2} \neq 0, \omega_{c,1} \neq \omega_{c,2}, \omega_{c,1} \neq -\omega_{c,2}, \omega_{c,1} \neq 3\omega_{c,2}, \omega_{c,1} \neq -3\omega_{c,2}, \omega_{c,2} \neq 3\omega_{c,1}$ und $\omega_{c,2} \neq -3\omega_{c,1}$ für beliebige $\omega_c \in \Omega$ sichergestellt (vgl. die durch C_z beschriebenen Bedingungen in Tabelle 5.1).

D Anhang zu Kapitel 6

D.1 Exkurs: Reinforcement Learning mit unvollständiger Zustandsinformation

Das Ignorieren nicht gemessener Zustände stellt eine Verletzung der Markov-Annahme dar. Der Lernprozess leidet dann unter hoher Varianz und würde langsamer und ungenauer, da der RL-Agent die Value Function im Allgemeinen rein basierend auf der Ausgangsgröße nicht exakt schätzen kann [PRH19], [Wer13, Abschnitt 1.3.1]. Der Umgang mit nur teilweise bekannten Zustandsgrößen stellt somit eine besondere Herausforderung bei lernbasierten Verfahren dar [Rec19], [DAMH19].

Eine Alternative dazu, die nicht messbaren Zustände der Aktuatordynamik zu ignorieren, stellt die Rekonstruktion fehlender Information aus vergangenen Stellgrößen oder Messungen vergangener Ausgangsgrößen dar. Da im betrachteten Anwendungsbeispiel die Längsdynamik jedoch nicht exakt bekannt ist, erweist sich ein modellbasierter Beobachterentwurf (vgl. beispielsweise [CDG93]) als ungeeignet. Grundsätzlich könnte hier die Verwendung eines rekurrenten neuronalen Netzes für die Approximation der Value Function Abhilfe schaffen [HS15]. Die Nutzung linearer rekurrenter Zellen führt bei deren Anwendung in einem realen Fahrzeug jedoch zu sehr langsamer Konvergenz der Actor-Gewichte [PKRH20]. Zudem erschweren rekurrente neuronale Netze die Verwendung von *Experience Replay* [MKS⁺13], [MKS⁺15], einem Mechanismus, der Daten mehrfach verwendet und somit zu höherer Dateneffizienz führt (vgl. auch Anhang D.2).

Alternativ kann der Zustandsvektor um vergangene Ausgangs- oder Stellgrößen erweitert werden [PMK99], [MKS⁺15], um daraus die fehlende Zustandsinformation ohne Verwendung eines Systemmodells zu rekonstruieren. Diese Idee wird auch in der vorliegenden Arbeit in Form eines FIR-Filters zur Schätzung eines Hilfszustands genutzt (vgl. [PRH19], [PKRH20]). Eine ausführlichere Diskussion bezüglich des Umgangs mit nur teilweise messbaren Zuständen ist in den Arbeiten von Puccetti et al. [PRH19] sowie Werbos [Wer13, Abschnitt 1.3.1] zu finden.

D.2 Actor-Critic-Grundlagen

In diesem Abschnitt wird eine kurze Einführung in die in Abschnitt 6.1 verwendeten Actor-Critic-Mechanismen gegeben. Eine ausführlichere Übersicht ist beispielsweise im Grundlagenwerk von Sutton und Barto [SB18] zu finden.

Im Folgenden seien \mathbf{x}_k der (gegebenenfalls erweiterte) Zustandsvektor, wie er dem RL-Agenten zur Verfügung gestellt wird¹⁴⁹, \mathbf{u}_k die Stellgröße und $r(\mathbf{x}_k, \mathbf{u}_k)$ die resultierenden Einschrittkosten. Als *Actor* wird das durch θ parametrisierte Regelgesetz $\mu_\theta(\mathbf{x}_k)$ bezeichnet, das zum Ziel hat, die Value Function

$$V^{\mu_\theta}(\mathbf{x}_k) = \sum_{\kappa=0}^{\infty} \gamma^\kappa r(\mathbf{x}_{k+\kappa}, \mu_\theta(\mathbf{x}_{k+\kappa})) \quad (\text{D.1})$$

zu minimieren. Der *Critic* hingegen sei ein Funktionsapproximator $Q_w^{\mu_\theta}$ mit dem Parameter w , der die Q-Function

$$Q^{\mu_\theta}(\mathbf{x}_k, \mathbf{u}_k) = r(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_w^{\mu_\theta}(\mathbf{x}_{k+1}, \mu_\theta(\mathbf{x}_{k+1})) \quad (\text{D.2})$$

approximiert. Aufgrund der Verwendung der Approximation $Q_w^{\mu_\theta}$ anstelle von Q^{μ_θ} ergibt sich der TD-Fehler

$$\delta = r(\mathbf{x}_k, \mathbf{u}_k) - Q_w^{\mu_\theta}(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_w^{\mu_\theta}(\mathbf{x}_{k+1}, \mu_\theta(\mathbf{x}_{k+1})) \quad (\text{D.3})$$

aus (D.2). Mit $L := \frac{1}{2}\delta^2$ und dem Gradienten¹⁵⁰

$$\frac{\partial L}{\partial w} = -\delta \frac{\partial Q_w^{\mu_\theta}(\mathbf{x}_k, \mathbf{u}_k)}{\partial w} \quad (\text{D.4})$$

kann der Parametervektor w des Critics mittels Gradientenabstieg oder anderer Optimierungsalgorithmen adaptiert werden.

Die Anpassung des Actor-Gewichts θ hat hingegen die Minimierung von $Q_w^{\mu_\theta}$ zum Ziel. Hierfür kann der sogenannte *Deterministic Policy Gradient* (DPG) [SLH⁺14]

$$\frac{\partial V^{\mu_\theta}}{\partial \theta} = \left. \frac{\partial Q_w^{\mu_\theta}(\mathbf{x}_k, \mathbf{u}_k)}{\partial \mathbf{u}_k} \right|_{\mathbf{u}_k = \mu_\theta(\mathbf{x}_k)} \frac{\partial \mu_\theta(\mathbf{x}_k)}{\partial \theta} \quad (\text{D.5})$$

verwendet werden, um einen Gradientenabstieg oder andere gradientenbasierte Optimierungsmethoden anzuwenden und den Reglerparameter θ zu adaptieren.

Während grundsätzlich beliebige Stellgrößen \mathbf{u}_k auf das System angewandt werden können, approximiert $Q_w^{\mu_\theta}$, welches nach (D.4) adaptiert wird, die Q-Function, die mit dem durch μ_θ gegebenen Regelgesetz assoziiert ist. Somit liegt eine Off-Policy-Methode vor (vgl. Abschnitt 2.1.4.4) und Anregungsrauschen führt nicht zu einem Offset der Critic-Schätzung. Ein weiterer, wesentlicher Vorteil der Off-Policy-Charakteristik ist, dass die zum Training verwendeten Datentupel verändert werden können, solange sie konsistent bezüglich der Bellman-Gleichung (D.2) und der zugrunde liegenden Systemdynamik bleiben. Diese Eigenschaft wird in Abschnitt 6.1.2.3 genutzt, indem zusätzliche Anregungssignale auf die Referenzparameter \mathbf{z}_k addiert werden. Schließlich kann bei dem verwendeten Off-Policy-Actor-Critic-Verfahren

¹⁴⁹ Im Falle eines erweiterten Zustandsvektors wird im Folgenden somit \mathbf{x}_k durch $\tilde{\mathbf{x}}_k$ ersetzt.

¹⁵⁰ In dieser Gradientenberechnung wird der Zielwert $Q_w^{\mu_\theta}(\mathbf{x}_{k+1}, \mu_\theta(\mathbf{x}_{k+1}))$ üblicherweise als konstant behandelt, siehe beispielsweise [SLH⁺14, Abschnitt 4.2].

der im RL häufig genutzte Mechanismus des *Experience Replay* [MKS⁺15] angewandt werden. Anstelle der Verwendung einzelner Datentupel wird in jedem Trainingsschritt ein sogenannter *Batch* $\{X_0, U_0, R_0, X_1\}$ aus M_B Datentupeln $\{\mathbf{x}_k, \mathbf{u}_k, r_k, \mathbf{x}_{k+1}\}$, die zufällig aus gespeicherten Daten gezogen werden, genutzt. Die Batch-Größe kann dabei prinzipiell für Actor und Critic unterschiedlich sein. Die Verwendung von Experience Replay kann zu einer verringerten Varianz des Trainingsergebnisses führen [MKS⁺15].

D.3 Wahl der Hyperparameter für das Online-Training

In Tabelle D.1 sind die Hyperparameter für das Online-Training der in Abschnitt 6.1 vorgestellten modellfreien, adaptiven Längsregelung eines realen Fahrzeugs inklusive Vorsteuerung gegeben.

Hyperparameter	Wert
Größe M des Datentupelspeichers	5000
Batch-Größe M_B für Actor und Critic	200
Anzahl h_{FIR} vergangener Stellgrößen für die Zustandsrekonstruktion	35
Anzahl der Ausgänge des FIR-Filters	1
Anzahl n_h zur Approximation des Sollgeschwindigkeitsverlaufs in (6.8)	2
Amplitude des Anregungssignals auf u_k	1 m s^{-2}
Haltezeit des Anregungssignals	2 s
Standardabweichung des Anregungsrauschens auf z_k	1
Parameter Q_y des Gütefunktional	1
Parameter R des Gütefunktional	0,1
Diskontierungsfaktor γ	0,95
Maximale Norm für die Adaption des Actors	0,01
Maximale Norm für die Adaption des Critics	10
Abtastzeit des Systems und Reglers Δt	0,02 s
Updaterate Δl des Trainingsvorgangs	0,6 s

Tabelle D.1: Wahl der Hyperparameter für das Online-Training.

D.4 Ergänzende Messdaten des realen Ball-auf-Platte-Systems

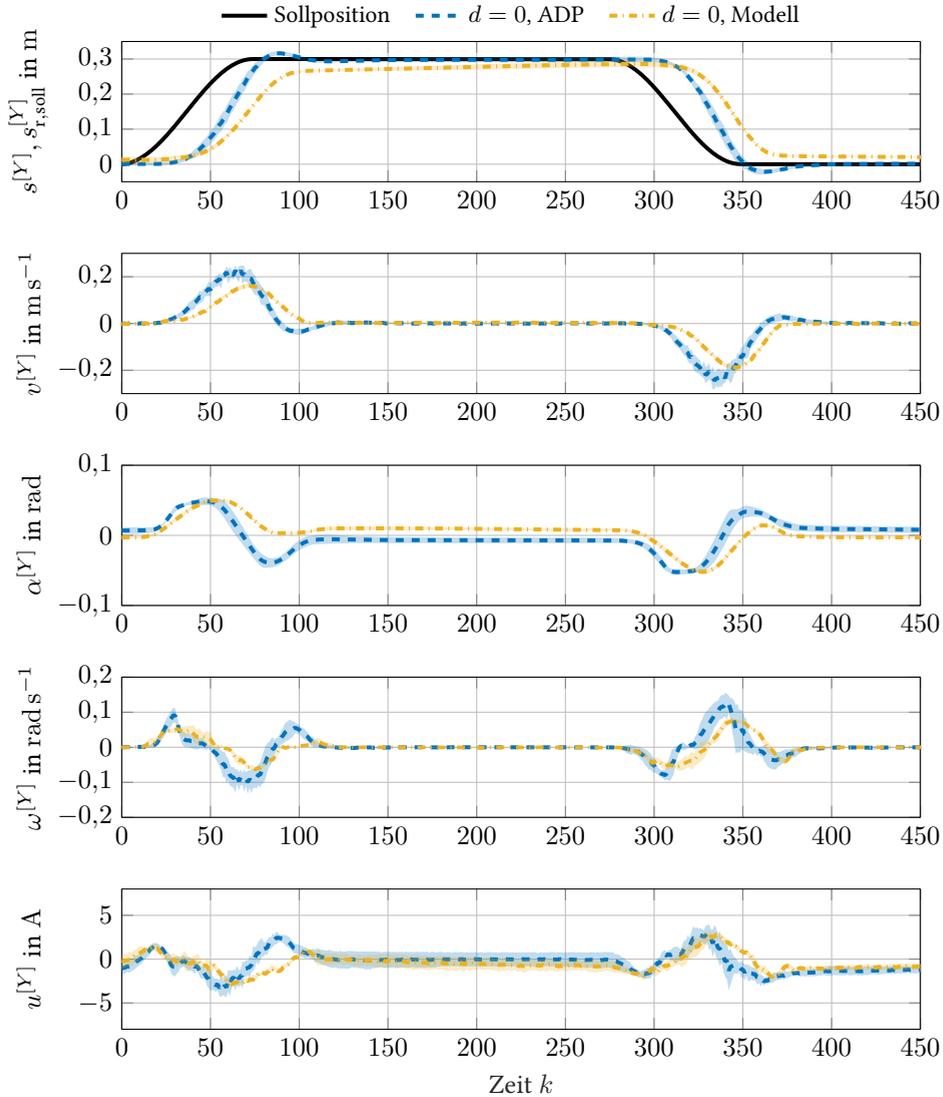


Abbildung D.1: Zustände und Stellgröße bei Vorgabe einer stationären Sollposition ($d = 0$) für den gelernten, ADP-basierten Regler (blau) sowie den modellbasierten Vergleichsregler (gelb).

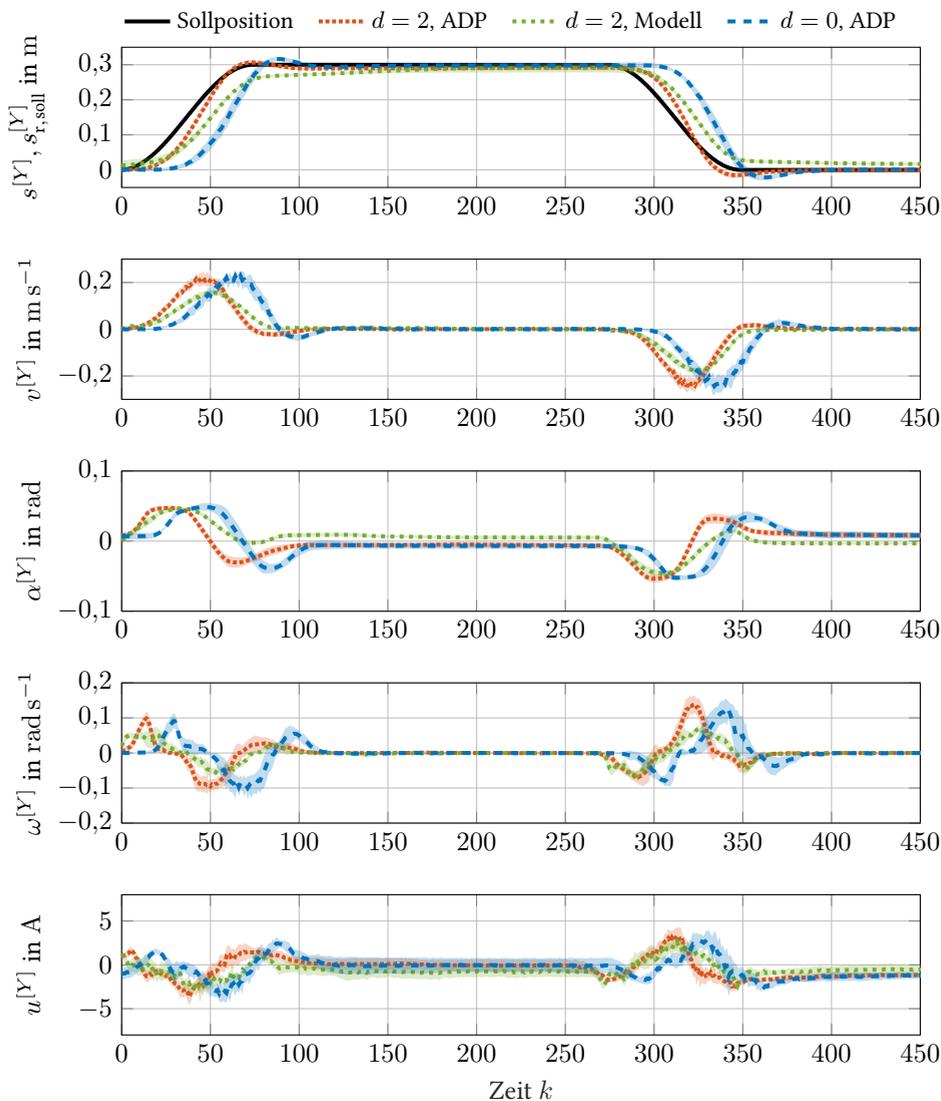


Abbildung D.2: Zustände und Stellgröße bei Vorgabe eines polynomiellen Sollverlaufs ($d = 2$) für den gelernten, ADP-basierten Regler (rot) sowie den modellbasierten Vergleichsregler (grün) im Vergleich zur Vorgabe einer stationären Sollposition ($d = 0$) für den gelernten, ADP-basierten Regler (blau).

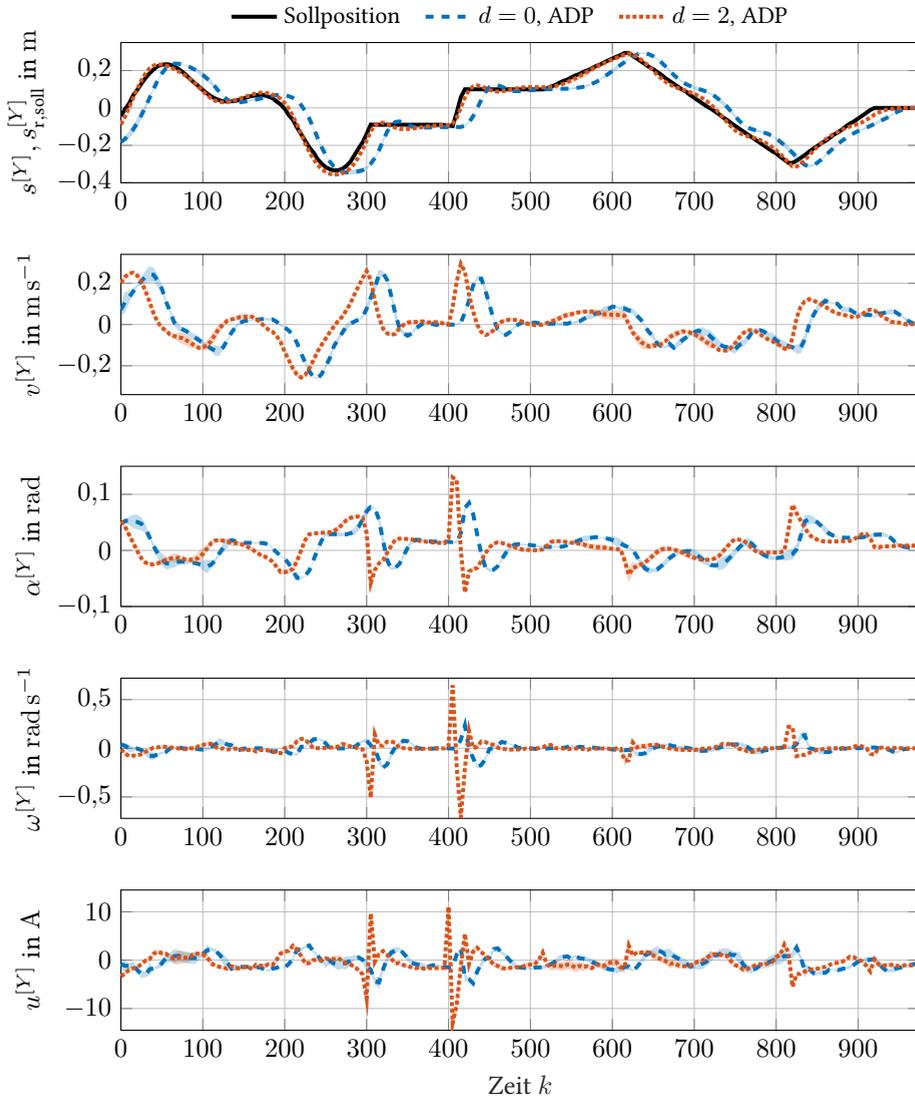


Abbildung D.3: Vergleich der gelernten Solltrajektorienfolgeregler bei Vorgabe eines quadratischen Sollverlaufs ($d = 2$, rot) und einer stationären Sollposition ($d = 0$, blau) für eine Validierungstrajektorie.

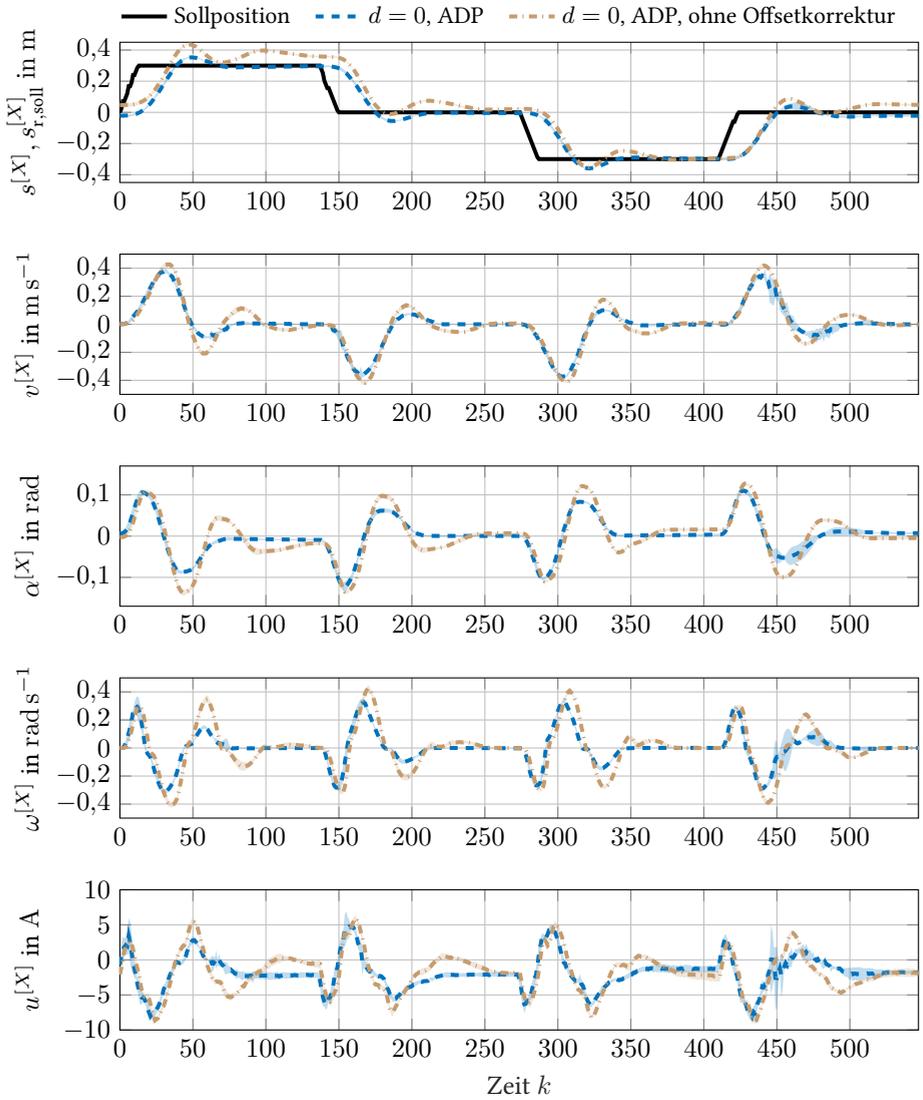


Abbildung D.4: Vergleich eines gelernten Reglers mit (blau) und ohne (braun) lernbare Offsetkorrektur bei Vorgabe einer stationären Sollposition ($d = 0$).

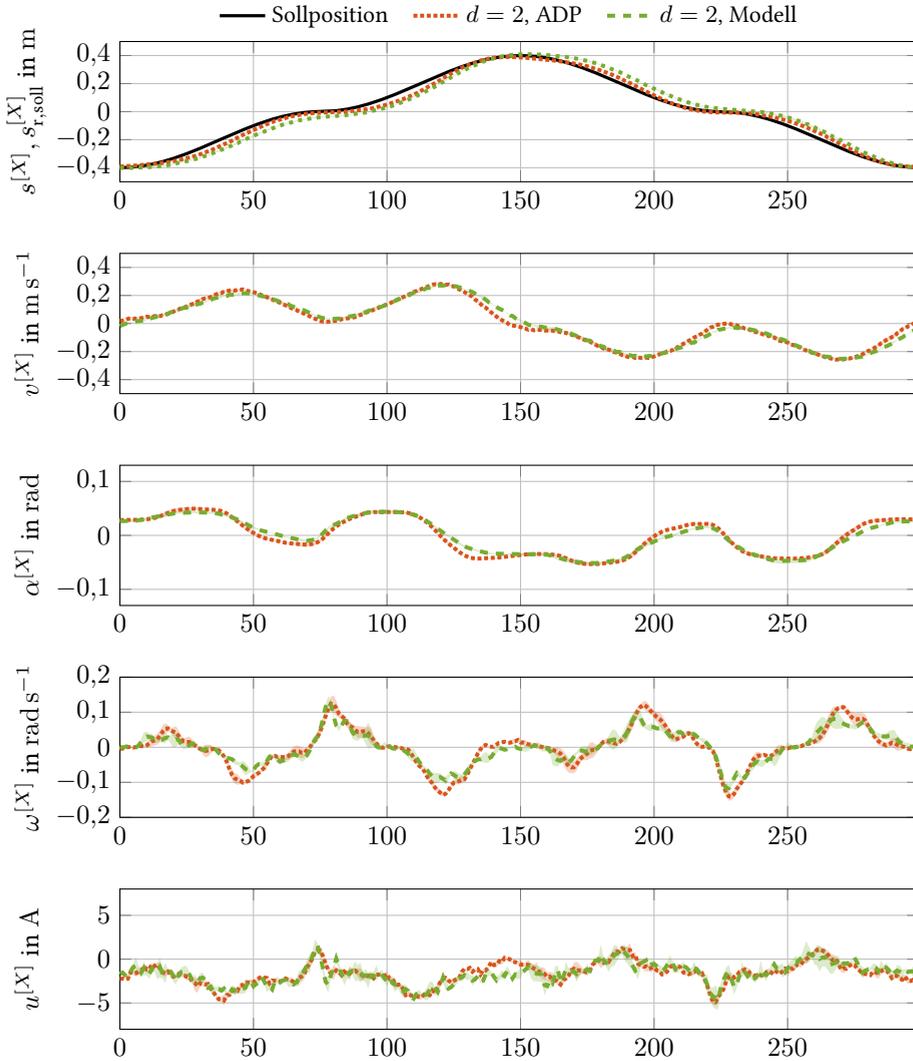


Abbildung D.5: Gleichzeitige Vorgabe eines Sollpositionsverlaufs in beiden Plattendimensionen. Gezeigt sind die resultierenden Zustände und Stellgrößen in X -Richtung für den gelernten Regler und den modellbasierten Vergleichsregler mit jeweils $d = 2$.

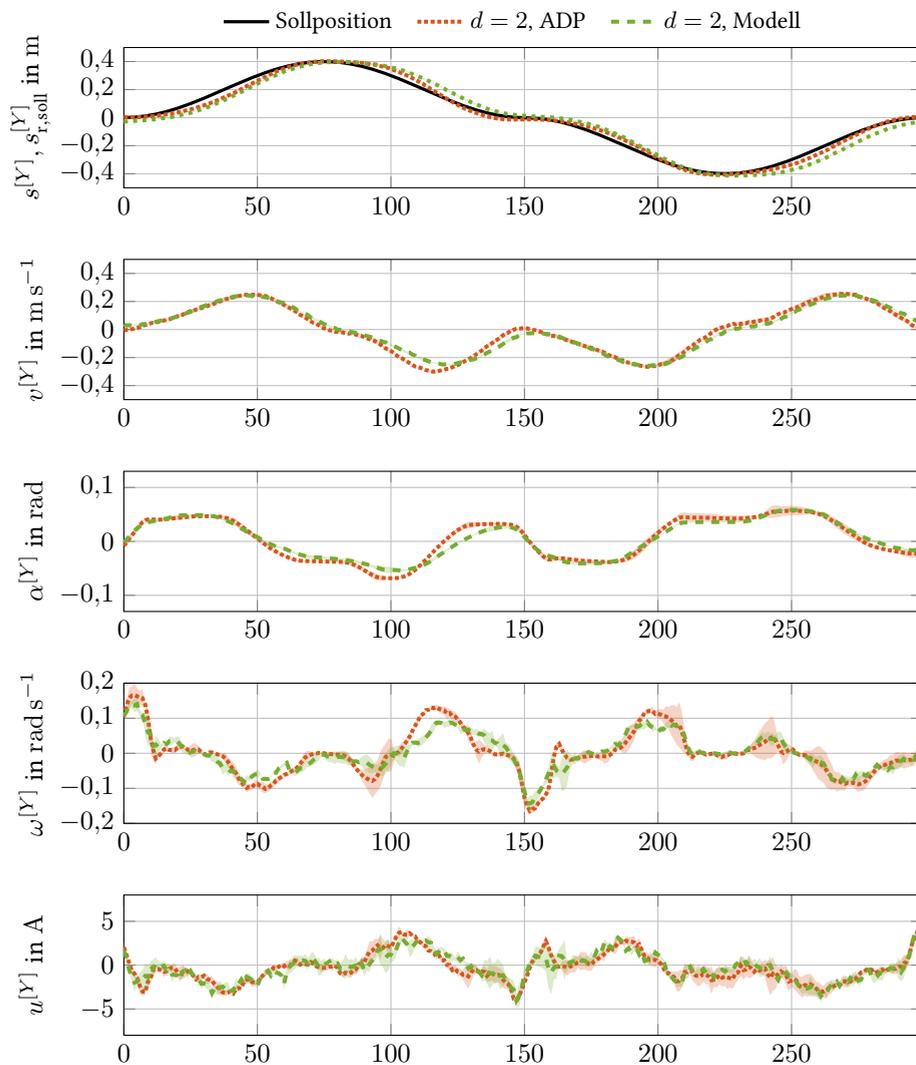


Abbildung D.6: Gleichzeitige Vorgabe eines Sollpositionsverlaufs in beiden Plattendimensionen. Gezeigt sind die resultierenden Zustände und Stellgrößen in Y -Richtung für den gelernten Regler und den modellbasierten Vergleichsregler mit jeweils $d = 2$.

Literaturverzeichnis

Öffentlich zugängliche Quellen

- [ÅB66] ÅSTRÖM, K.-J.; BOHLIN, T.: Numerical Identification of Linear Dynamic Systems from Normal Operating Records. In: *IFAC Proceedings Volumes 2* (1966), Nr. 2, S. 96–111
- [ABB⁺02] AWTAR, S.; BERNARD, C.; BOKLUND, N.; MASTER, A.; UEDA, D.; CRAIG, K.: Mechatronic Design of a Ball-on-Plate Balancing System. In: *Mechatronics 12* (2002), Nr. 2, S. 217–228
- [ABC⁺20] ANDRYCHOWICZ, M.; BAKER, B.; CHOCIEJ, M.; JÓZEFOWICZ, R.; MCGREW, B.; PACHOCKI, J.; PETRON, A.; PLAPPERT, M.; POWELL, G.; RAY, A.; SCHNEIDER, J.; SIDOR, S.; TOBIN, J.; WELINDER, P.; WENG, L.; ZAREMBA, W.: Learning Dexterous In-Hand Manipulation. In: *The International Journal of Robotics Research 39* (2020), Nr. 1, S. 3–20
- [Ada18] ADAMY, J.: *Nichtlineare Systeme und Regelungen*. 3., aktualisierte Auflage. Berlin: Springer Vieweg, 2018
- [AE10] AMANN, H.; ESCHER, J.: *Analysis I*. 3. Auflage. Basel: Birkhäuser, 2010
- [AG06] ADETOLA, V.; GUAY, M.: Excitation Signal Design for Parameter Convergence in Adaptive Control of Linearizable Systems. In: *Proceedings of the 45th IEEE Conference on Decision and Control*, 2006, S. 447–452
- [AG08] ADETOLA, V.; GUAY, M.: Finite-Time Parameter Estimation in Adaptive Control of Nonlinear Systems. In: *IEEE Transactions on Automatic Control 53* (2008), Nr. 3, S. 807–811
- [AKL05] ABU-KHALAF, M.; LEWIS, F. L.: Nearly Optimal Control Laws for Nonlinear Systems with Saturating Actuators Using a Neural Network HJB Approach. In: *Automatica 41* (2005), Nr. 5, S. 779–791
- [AL84] ARNOLD, W. F.; LAUB, A. J.: Generalized Eigenproblem Algorithms and Software for Algebraic Riccati Equations. In: *Proceedings of the IEEE 72* (1984), Nr. 12, S. 1746–1754
- [AM89] ANDERSON, B. D.; MOORE, J. B.: *Optimal Control: Linear Quadratic Methods*. Englewood Cliffs: Prentice Hall, 1989 (Prentice Hall information and system sciences series)

- [Apo67] APOSTOL, T. M.: *Calculus*. 2. Auflage. New York: Wiley, 1967
- [ATLAK07] AL-TAMIMI, A.; LEWIS, F. L.; ABU-KHALAF, M.: Model-Free Q-Learning Designs for Linear Discrete-Time Zero-Sum Games with Application to H-Infinity Control. In: *Automatica* 43 (2007), Nr. 3, S. 473–481
- [ATLAK08] AL-TAMIMI, A.; LEWIS, F. L.; ABU-KHALAF, M.: Discrete-Time Nonlinear HJB Solution Using Approximate Dynamic Programming: Convergence Proof. In: *IEEE Transactions on Systems, Man, and Cybernetics. Part B: Cybernetics* 38 (2008), Nr. 4, S. 943–949
- [ÅW95] ÅSTRÖM, K. J.; WITTENMARK, B.: *Adaptive Control*. 2. Auflage. Reading, Mass.: Addison-Wesley, 1995
- [Axl15] AXLER, S. J.: *Linear Algebra Done Right*. 3. Auflage. Cham, Heidelberg, New York, Dordrecht, London: Springer, 2015
- [BA18] BERNHARD, S.; ADAMY, J.: LQ Optimal Tracking with Unbounded Cost for Unknown Dynamics: An Adaptive Dynamic Programming Approach. In: *European Control Conference*. 2018
- [BBdE10] BUŞONIU, L.; BABUŞKA, R.; DE SCHUTTER, B.; ERNST, D.: *Reinforcement Learning and Dynamic Programming Using Function Approximators*. Boca Raton, London, New York: CRC press, 2010
- [BBT⁺18] BUŞONIU, L.; BRUIN, T. de; TOLIĆ, D.; KOBER, J.; PALUNKO, I.: Reinforcement Learning for Control: Performance, Stability, and Deep Approximators. In: *Annual Reviews in Control* 46 (2018), S. 8–28
- [Bel57a] BELLMAN, R.: *Dynamic Programming*. 1. Auflage. Princeton: Princeton University Press, 1957
- [Bel57b] BELLMAN, R.: A Markovian Decision Process. In: *Journal of Mathematics and Mechanics* 6 (1957), Nr. 5, S. 679–684
- [Ben89] BENTHAM, J.: *An Introduction to the Principles of Morals and Legislation*. London: T. Payne and Son, 1789
- [Ber95] BERTSEKAS, D. P.: *Dynamic Programming and Optimal Control*. Belmont, Mass.: Athena Scientific, 1995
- [Ber17] BERTSEKAS, D. P.: Value and Policy Iterations in Optimal Control and Adaptive Dynamic Programming. In: *IEEE Transactions on Neural Networks and Learning Systems* 28 (2017), Nr. 3, S. 500–509
- [Ber20a] BERNHARD, S.: *Optimale Folgeregelung über unendliche Horizonte und optimale Output Regulation für quadratische, über- und unteraktuierte Systeme*. Darmstadt, Technische Universität Darmstadt, Dissertation, 2020

- [Ber20b] BERTSEKAS, D.: Multiagent Value Iteration Algorithms in Dynamic Programming and Reinforcement Learning. In: *Results in Control and Optimization 1* (2020), S. 1–10
- [BFGB12] BRAESCU, F. C.; FERARIU, L.; GILCA, R.; BORDIANU, V.: Ball on Plate Balancing System for Multi-Discipline Educational Purposes. In: *16th International Conference on System Theory, Control and Computing (ICSTCC)*, 2012, S. 1–6
- [BHWM12] BURG, K.; HAF, H.; WILLE, F.; MEISTER, A.: *Höhere Mathematik für Ingenieure: Band II: Lineare Algebra*. 7. Auflage. Wiesbaden: Vieweg+Teubner Verlag, 2012
- [Bia17] BIAN, T.: *Biologically Inspired Adaptive Optimal Control and Learning*. New York, New York University, Dissertation, 2017
- [Bit84] BITMEAD, R. R.: Persistence of Excitation Conditions and the Convergence of Adaptive Schemes. In: *IEEE Transactions on Information Theory* 30 (1984), Nr. 2, S. 183–191
- [BJ16a] BIAN, T.; JIANG, Z.-P.: Value Iteration, Adaptive Dynamic Programming, and Optimal Control of Nonlinear Systems. In: *55th IEEE Conference on Decision and Control (CDC)*, 2016, S. 3375–3380
- [BJ16b] BIAN, T.; JIANG, Z.-P.: Value Iteration and Adaptive Dynamic Programming for Data-Driven Adaptive Optimal Control Design. In: *Automatica* 71 (2016), S. 348–360
- [BJD10] BHASIN, S.; JOHNSON, M.; DIXON, W. E.: A Model-Free Robust Policy Iteration Algorithm for Optimal Control of Nonlinear Systems. In: *49th IEEE Conference on Decision and Control (CDC)*, 2010, S. 3060–3065
- [BK18] BÜCHEL, M.; KNOLL, A.: Deep Reinforcement Learning for Predictive Longitudinal Control of Automated Vehicles. In: *21st International Conference on Intelligent Transportation Systems (ITSC)*, 2018, S. 2391–2397
- [BKJ⁺13] BHASIN, S.; KAMALAPURKAR, R.; JOHNSON, M.; VAMVOUDAKIS, K. G.; LEWIS, F. L.; DIXON, W. E.: A Novel Actor-Critic-Identifier Architecture for Approximate Optimal Control of Uncertain Nonlinear Systems. In: *Automatica* 49 (2013), Nr. 1, S. 82–92
- [BO99] BAŞAR, T.; OLSDER, G. J.: *Dynamic Noncooperative Game Theory*. 2. Auflage. Society for Industrial and Applied Mathematics, 1999
- [Bro92] BROOKS, R. A.: Artificial Life and Real Robots. In: *Proceedings of the First European Conference on Artificial Life*, MIT Press, 1992, S. 3–10
- [BS83] BOYD, S.; SASTRY, S.: On Parameter Convergence in Adaptive Control. In: *Systems & Control Letters* 3 (1983), Nr. 6, S. 311–319
- [BS86] BOYD, S.; SASTRY, S. S.: Necessary and Sufficient Conditions for Parameter Convergence in Adaptive Control. In: *Automatica* 22 (1986), Nr. 6, S. 629–639

- [BSA83] BARTO, A. G.; SUTTON, R. S.; ANDERSON, C. W.: Neuronlike Adaptive Elements that Can Solve Difficult Learning Control Problems. In: *IEEE Transactions on Systems, Man, and Cybernetics* 13 (1983), Nr. 5, S. 834–846
- [BSMM08] BRONSTEIN, I. N.; SEMENDJAEV, K. A.; MUSIOL, G.; MÜHLIG, H.: *Taschenbuch der Mathematik*. 7., vollständig überarbeitete und ergänzte Auflage. Frankfurt am Main: Verlag Harri Deutsch, 2008
- [BSMM13] BRONSTEIN, I. N.; SEMENDJAEV, K. A.; MUSIOL, G.; MÜHLIG, H.: *Taschenbuch der Mathematik*. 9. Auflage. Haan-Gruiten: Verlag Europa-Lehrmittel, 2013
- [BSW97] BEARD, R. W.; SARIDIS, G. N.; WEN, J. T.: Galerkin Approximations of the Generalized Hamilton-Jacobi-Bellman Equation. In: *Automatica* 33 (1997), Nr. 12, S. 2159–2177
- [BT96] BERTSEKAS, D. P.; TSITSIKLIS, J. M.: *Neuro-Dynamic Programming*. Belmont, Mass.: Athena Scientific, 1996
- [BU16] BOHN, C.; UNBEHAUEN, H.: *Identifikation dynamischer Systeme*. Wiesbaden: Springer Fachmedien, 2016
- [BYB94] BRADTKE, S. J.; YDSTIE, B. E.; BARTO, A. G.: Adaptive Linear Quadratic Control Using Policy Iteration. In: *American Control Conference (ACC)*, 1994, S. 3475–3479
- [CDG93] CICCARELLA, G.; DALLA MORA, M.; GERMANI, A.: A Luenberger-like Observer for Nonlinear Systems. In: *International Journal of Control* 57 (1993), Nr. 3, S. 537–556
- [CF14] CALANCA, A.; FIORINI, P.: Human-adaptive Control of Series Elastic Actuators. In: *Robotica* 32 (2014), Nr. 8, S. 1301–1316
- [Che78] CHENEY, E. W.: *Introduction to Approximation Theory*. 5. Auflage. New York: McGraw-Hill, 1978 (International Series in Pure and Applied Mathematics)
- [CHL19] COLLINS, J.; HOWARD, D.; LEITNER, J.: Quantifying the Reality Gap in Robotic Manipulation Tasks. In: *International Conference on Robotics and Automation (ICRA)*, 2019, S. 6706–6712
- [CJ10] CHOWDHARY, G.; JOHNSON, E.: Concurrent Learning for Convergence in Adaptive Control without Persistency of Excitation. In: *49th IEEE Conference on Decision and Control (CDC)*, 2010, S. 3674–3679
- [Coh13] COHN, D. L.: *Measure Theory*. 2. Auflage. Berlin, New York: Springer, 2013
- [DAMH19] DULAC-ARNOLD, G.; MANKOWITZ, D.; HESTER, T.: Challenges of Real-World Reinforcement Learning. In: *36th International Conference on Machine Learning*, 2019

- [DCd11] DESJARDINS, C.; CHAIB-DRAA, B.: Cooperative Adaptive Cruise Control: A Reinforcement Learning Approach. In: *IEEE Transactions on Intelligent Transportation Systems* 12 (2011), Nr. 4, S. 1248–1260
- [DCP96] DEVASIA, S.; CHEN, D.; PADEN, B.: Nonlinear Inversion-Based Output Tracking. In: *IEEE Transactions on Automatic Control* 41 (1996), Nr. 7, S. 930–942
- [Dev02] DEVASIA, S.: Should Model-Based Inverse Inputs be Used as Feedforward Under Plant Uncertainty? In: *IEEE Transactions on Automatic Control* 47 (2002), Nr. 11, S. 1865–1871
- [DHS17] DUŠEK, F.; HONC, D.; SHARMA, K. R.: Modelling of Ball and Plate System Based on First Principle Model and Optimal Control. In: *21st International Conference on Process Control*, 2017, S. 216–221
- [DJ09] DIERKS, T.; JAGANNATHAN, S.: Optimal Tracking Control of Affine Nonlinear Discrete-Time Systems with Unknown Internal Dynamics. In: *Joint 48th IEEE Conference on Decision and Control and 28th Chinese Control Conference*, 2009, S. 6750–6755
- [DJ10] DIERKS, T.; JAGANNATHAN, S.: Optimal Control of Affine Nonlinear Continuous-Time Systems. In: *American Control Conference (ACC)*, 2010, S. 1568–1573
- [DJ11] DIERKS, T.; JAGANNATHAN, S.: Online Optimal Control of Nonlinear Discrete-Time Systems Using Approximate Dynamic Programming. In: *Journal of Control Theory and Applications* 9 (2011), Nr. 3, S. 361–369
- [dKTB18] DE BRUIN, T.; KOBER, J.; TUYLS, K.; BABUŠKA, R.: Experience Selection in Deep Reinforcement Learning for Control. In: *Journal of Machine Learning Research* 19 (2018), Nr. 9, S. 1–56
- [DLL18] DU, Y.; LIU, C.; LI, Y.: Velocity Control Strategies to Improve Automated Vehicle Driving Comfort. In: *IEEE Intelligent Transportation Systems Magazine* 10 (2018), Nr. 1, S. 8–18
- [DLL⁺19] DUAN, J.; LI, S. E.; LIU, Z.; BUJARBARUAH, M.; CHENG, B.: Generalized Policy Iteration for Optimal Control in Continuous Time. In: *arXiv preprint 1909.05402* (2019), S. 1–10
- [Doy00] DOYA, K.: Reinforcement Learning in Continuous Time and Space. In: *Neural Computation* 12 (2000), Nr. 1, S. 219–245
- [DS10] DZJADYK, V. K.; SHEVCHUK, I. A.: *Theory of Uniform Approximation of Functions by Polynomials*. Berlin, New York: De Gruyter, 2010
- [EBHS07] ENNS, D.; BUGAJSKI, D.; HENDRICK, R.; STEIN, G.: Dynamic Inversion: An Evolving Methodology for Flight Control Design. In: *International Journal of Control* 59 (2007), Nr. 1, S. 71–91

- [EP10] ELДАР, Y. C.; POHL, V.: Recovering Signals from Lowpass Data. In: *IEEE Transactions on Signal Processing* 58 (2010), Nr. 5, S. 2636–2646
- [FC16] FU, Y.; CHAI, T.: Online Solution of Two-Player Zero-Sum Games for Continuous-Time Nonlinear Systems with Completely Unknown Dynamics. In: *IEEE Transactions on Neural Networks and Learning Systems* 27 (2016), Nr. 12, S. 2577–2587
- [FFH17] FLAD, M.; FRÖHLICH, L.; HOHMANN, S.: Cooperative Shared Control Driver Assistance Systems Based on Motion Primitives and Differential Games. In: *IEEE Transactions on Human-Machine Systems* 47 (2017), Nr. 5, S. 711–722
- [Fla16] FLAD, M.: *Kooperative Regelungskonzepte auf Basis der Spieltheorie und deren Anwendung auf Fahrerassistenzsysteme*. Karlsruhe, Karlsruher Institut für Technologie, Dissertation, 2016
- [FLMR95] FLIESS, M.; LÉVINE, J.; MARTIN, P.; ROUCHON, P.: Flatness and Defect of Nonlinear Systems: Introductory Theory and Examples. In: *International Journal of Control* 61 (1995), Nr. 6, S. 1327–1361
- [Föll16] FÖLLINGER, O.: *Regelungstechnik: Einführung in die Methoden und ihre Anwendung*. 12. Auflage. Berlin, Offenbach: VDE-Verlag, 2016
- [FOSH14] FLAD, M.; OTTEN, J.; SCHWAB, S.; HOHMANN, S.: Steering Driver Assistance System: A Systematic Cooperative Shared Control Design Approach. In: *IEEE International Conference on Systems, Man and Cybernetics*, 2014, S. 3585–3592
- [FWS⁺18] FEINBERG, V.; WAN, A.; STOICA, I.; JORDAN, M. I.; GONZALEZ, J. E.; LEVINE, S.: Model-Based Value Expansion for Efficient Model-Free Reinforcement Learning. In: *arXiv preprint 1803.00101* (2018)
- [FY16] FAN, Q.-Y.; YANG, G.-H.: Adaptive Actor-Critic Design-Based Integral Sliding-Mode Control for Partially Unknown Nonlinear Systems With Input Disturbances. In: *IEEE Transactions on Neural Networks and Learning Systems* 27 (2016), Nr. 1, S. 165–177
- [GBLB12] GRONDMAN, I.; BUSONIU, L.; LOPES, G. A. D.; BABUSKA, R.: A Survey of Actor-Critic Reinforcement Learning: Standard and Natural Policy Gradients. In: *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 42 (2012), Nr. 6, S. 1291–1307
- [Gee07] GEERING, H. P.: *Optimal Control with Engineering Applications*. Berlin, Heidelberg: Springer, 2007
- [GGH⁺18] GAITSGORY, V.; GRÜNE, L.; HÖGER, M.; KELLETT, C. M.; WELLER, S. R.: Stabilization of Strictly Dissipative Discrete Time Systems with Discounted Optimal Control. In: *Automatica* 93 (2018), S. 311–320

- [GHLL17] GU, S.; HOLLY, E.; LILICRAP, T.; LEVINE, S.: Deep Reinforcement Learning for Robotic Manipulation with Asynchronous Off-Policy Updates. In: *IEEE International Conference on Robotics and Automation*, 2017, S. 3389–3396
- [GJ15] GAO, W.; JIANG, Z.-P.: Linear Optimal Tracking Control: An Adaptive Dynamic Programming Approach. In: *American Control Conference (ACC)*, 2015, S. 4929–4934
- [GJ16] GAO, W.; JIANG, Z.-P.: Adaptive Dynamic Programming and Adaptive Optimal Output Regulation of Linear Systems. In: *IEEE Transactions on Automatic Control* 61 (2016), Nr. 12, S. 4164–4169
- [Gli11] GLIMCHER, P. W.: Understanding Dopamine and Reinforcement Learning: The Dopamine Reward Prediction Error Hypothesis. In: *Proceedings of the National Academy of Sciences of the United States of America* (2011), S. 15647–15654
- [GM86] GREEN, M.; MOORE, J. B.: Persistence of Excitation in Linear Systems. In: *Systems & Control Letters* 7 (1986), Nr. 5, S. 351–360
- [Gör17] GÖRGES, D.: Relations between Model Predictive Control and Reinforcement Learning. In: *IFAC-PapersOnLine* 50 (2017), Nr. 1, S. 4920–4928
- [Gör19] GÖRGES, D.: Distributed Adaptive Linear Quadratic Control using Distributed Reinforcement Learning. In: *IFAC-PapersOnLine* 52 (2019), Nr. 11, S. 218–223
- [Gos09] GOSAVI, A.: Reinforcement Learning: A Tutorial Survey and Recent Advances. In: *INFORMS Journal on Computing* 21 (2009), Nr. 2, S. 178–192
- [GSK19] GAVRISHCHAKA, V.; SENYUKOVA, O.; KOEPKE, M.: Synergy of Physics-Based Reasoning and Machine Learning in Biomedical Applications: Towards Unlimited Deep Learning with Limited Data. In: *Advances in Physics: X* 4 (2019), Nr. 1, S. 204–256
- [HDB94] HOUK, J. C.; DAVIS, J. L.; BEISER, D. G.: A Model of How the Basal Ganglia Generate and Use Neural Signals that Predict Reinforcement. In: *Models of Information Processing in the Basal Ganglia*. 1994, S. 249–270
- [Hey16] HEYDARI, A.: Analyzing Policy Iteration in Optimal Control. In: *American Control Conference (ACC)*, 2016, S. 5728–5733
- [Hir79] HIRSCHORN, R. M.: Invertibility of Nonlinear Control Systems. In: *SIAM Journal on Control and Optimization* 17 (1979), Nr. 2, S. 289–297
- [HJK20] HYATT, P.; JOHNSON, C. C.; KILLPACK, M. D.: Model Reference Predictive Adaptive Control for Large-Scale Soft Robots. In: *Frontiers in Robotics and AI* 7 (2020), 132
- [HL14] HUANG, Y.; LIU, D.: Neural-Network-Based Optimal Tracking Control Scheme for a Class of Unknown Discrete-Time Nonlinear Systems Using Iterative ADP Algorithm. In: *Neurocomputing* 125 (2014), S. 46–56

- [HLM03] HESPANHA, J. P.; LIBERZON, D.; MORSE, A.: Overcoming the Limitations of Adaptive Control by Means of Logic-Based Switching. In: *Systems & Control Letters* 49 (2003), Nr. 1, S. 49–65
- [HMT⁺21] HAMAYA, M.; MATSUBARA, T.; TERAMAE, T.; NODA, T.; MORIMOTO, J.: Design of Physical User-Robot Interactions for Model Identification of Soft Actuators on Exoskeleton Robots. In: *The International Journal of Robotics Research* 40 (2021), Nr. 1, S. 397–410
- [HS15] HAUSKNECHT, M.; STONE, P.: Deep Recurrent Q-Learning for Partially Observable MDPs. In: *AAAI Fall Symposium Series*, 2015, S. 1–7
- [HSSH17] HWANGBO, J.; SA, I.; SIEGWART, R.; HUTTER, M.: Control of a Quadrotor with Reinforcement Learning. In: *IEEE Robotics and Automation Letters* 2 (2017), Nr. 4, S. 2096–2103
- [HSW90] HORNIK, K.; STINCHCOMBE, M.; WHITE, H.: Universal Approximation of an Unknown Mapping and its Derivatives Using Multilayer Feedforward Networks. In: *Neural Networks* 3 (1990), Nr. 5, S. 551–560
- [HWL21] HA, M.; WANG, D.; LIU, D.: Generalized Value Iteration for Discounted Optimal Control with Stability Analysis. In: *Systems & Control Letters* 147 (2021), S. 1–7
- [HXH⁺19] HUANG, Z.; XU, X.; HE, H.; TAN, J.; SUN, Z.: Parameterized Batch Reinforcement Learning for Longitudinal Control of Autonomous Land Vehicles. In: *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 49 (2019), Nr. 4, S. 730–741
- [IS96] IOANNOU, P. A.; SUN, J.: *Robust Adaptive Control*. Upper Saddle River, New Jersey: Prentice-Hall Inc., 1996
- [Isi89] ISIDORI, A.: *Nonlinear Control Systems*. 2. Auflage. Berlin, Heidelberg: Springer, 1989
- [JALG18] JENKINS, B. M.; ANNASWAMY, A. M.; LAVRETSKY, E.; GIBSON, T. E.: Convergence Properties of Adaptive Systems and the Definition of Exponential Stability. In: *SIAM Journal on Control and Optimization* 56 (2018), Nr. 4, S. 2463–2484
- [JJ12] JIANG, Y.; JIANG, Z.-P.: Computational Adaptive Optimal Control for Continuous-Time Linear Systems with Completely Unknown Dynamics. In: *Automatica* 48 (2012), Nr. 10, S. 2699–2704
- [JJ14a] JIANG, Y.; JIANG, Z.-P.: Adaptive Dynamic Programming as a Theory of Sensorimotor Control. In: *Biological Cybernetics* 108 (2014), Nr. 4, S. 459–473
- [JJ14b] JIANG, Y.; JIANG, Z.-P.: Global Adaptive Dynamic Programming for Continuous-Time Nonlinear Polynomial Systems. In: *IFAC Proceedings Volumes* 47 (2014), Nr. 3, S. 9756–9761

- [JJ14c] JIANG, Y.; JIANG, Z.-P.: Robust Adaptive Dynamic Programming and Feedback Stabilization of Nonlinear Systems. In: *IEEE Transactions on Neural Networks and Learning Systems* 25 (2014), Nr. 5, S. 882–893
- [JKBD15] JOHNSON, M.; KAMALAPURKAR, R.; BHASIN, S.; DIXON, W. E.: Approximate N-Player Nonzero-Sum Game Solution for an Uncertain Continuous Nonlinear System. In: *IEEE Transactions on Neural Networks and Learning Systems* 26 (2015), Nr. 8, S. 1645–1658
- [JZLH19] JIANG, H.; ZHANG, H.; LUO, Y.; HAN, J.: Neural-Network-Based Robust Control Schemes for Nonlinear Multiplayer Systems with Uncertainties via Adaptive Dynamic Programming. In: *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 49 (2019), Nr. 3, S. 579–588
- [KAV⁺14] KATZOURAKIS, D. I.; ABBINK, D. A.; VELENIS, E.; HOLWEG, E.; HAPPEE, R.: Drivers Arms Time-Variant Neuromuscular Admittance During Real Car Test-Track Driving. In: *IEEE Transactions on Instrumentation and Measurement* 63 (2014), Nr. 1, S. 221–230
- [KB18] KNABNER, P.; BARTH, W.: *Lineare Algebra: Grundlagen und Anwendungen*. 2. Auflage. Berlin, Heidelberg: Springer, 2018
- [KBJ⁺21] KUUTTI, S.; BOWDEN, R.; JIN, Y.; BARBER, P.; FALLAH, S.: A Survey of Deep Learning Applications to Autonomous Vehicle Control. In: *IEEE Transactions on Intelligent Transportation Systems* 22 (2021), Nr. 2, S. 712–733
- [KBP13] KOBER, J.; BAGNELL, J. A.; PETERS, J.: Reinforcement Learning in Robotics: A Survey. In: *The International Journal of Robotics Research* 32 (2013), Nr. 11, S. 1238–1274
- [KCS03] KNUPLEZ, A.; CHOWDHURY, A.; SVECKO, R.: Modeling and Control Design for the Ball and Plate System. In: *IEEE International Conference on Industrial Technology*, 2003, S. 1064–1067
- [KDBD15] KAMALAPURKAR, R.; DINH, H.; BHASIN, S.; DIXON, W. E.: Approximate Optimal Trajectory Tracking for Continuous-Time Nonlinear Systems. In: *Automatica* 51 (2015), S. 40–48
- [Kha02] KHALIL, H. K.: *Nonlinear Systems*. 3. Auflage. Upper Saddle River: Prentice Hall, 2002
- [KHL⁺12] KHAN, S. G.; HERRMANN, G.; LEWIS, F. L.; PIPE, T.; MELHUSH, C.: Reinforcement Learning and Optimal Adaptive Control: An Overview and Implementation Examples. In: *Annual Reviews in Control* 36 (2012), Nr. 1, S. 42–59
- [KIP⁺18] KALASHNIKOV, D.; IRPAN, A.; PASTOR, P.; IBARZ, J.; HERZOG, A.; JANG, E.; QUILLEN, D.; HOLLY, E.; KALAKRISHNAN, M.; VANHOUCKE, V.; LEVINE, S.: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation. In: *Proceedings of Machine Learning Research* Bd. 87. 2018, S. 651–673

- [KKD14] KAMALAPURKAR, R.; KLOTZ, J. R.; DIXON, W. E.: Concurrent Learning-Based Approximate Feedback-Nash Equilibrium Solution of N-Player Nonzero-Sum Differential Games. In: *IEEE/CAA Journal of Automatica Sinica* 1 (2014), Nr. 3, S. 239–247
- [KL15] KIUMARSI, B.; LEWIS, F. L.: Actor-Critic-Based Optimal Tracking for Partially Unknown Nonlinear Discrete-Time Systems. In: *IEEE Transactions on Neural Networks and Learning Systems* 26 (2015), Nr. 1, S. 140–151
- [KLJ17] KIUMARSI, B.; LEWIS, F. L.; JIANG, Z.-P.: H_∞ Control of Linear Discrete-Time Systems: Off-Policy Reinforcement Learning. In: *Automatica* 78 (2017), S. 144–152
- [KLL15] KIUMARSI, B.; LEWIS, F. L.; LEVINE, D. S.: Optimal Control of Nonlinear Discrete Time-Varying Systems Using a New Neural Network Approximation Structure. In: *Neurocomputing* 156 (2015), S. 157–165
- [KLM96] KAELBLING, L. P.; LITTMAN, M. L.; MOORE, A. W.: Reinforcement Learning: A Survey. In: *Journal of Artificial Intelligence Research* 4 (1996), Nr. 1, S. 237–285
- [KLM⁺14] KIUMARSI, B.; LEWIS, F. L.; MODARES, H.; KARIMPOUR, A.; NAGHIBI-SISTANI, M.-B.: Reinforcement Q-Learning for Optimal Tracking Control of Linear Discrete-Time Systems with Unknown Dynamics. In: *Automatica* 50 (2014), Nr. 4, S. 1167–1175
- [KLNSK15] KIUMARSI, B.; LEWIS, F. L.; NAGHIBI-SISTANI, M.; KARIMPOUR, A.: Optimal Tracking Control of Unknown Discrete-Time Linear Systems Using Input-Output Measured Data. In: *IEEE Transactions on Cybernetics* 45 (2015), Nr. 12, S. 2770–2779
- [KT03] KONDA, V. R.; TSITSIKLIS, J. N.: On Actor-Critic Algorithms. In: *SIAM Journal on Control and Optimization* 42 (2003), Nr. 4, S. 1143–1166
- [KT04] KULIKOV, G. G. (Hrsg.); THOMPSON, H. A. (Hrsg.): *Dynamic Modelling of Gas Turbines: Identification, Simulation, Condition Monitoring and Optimal Control*. London: Springer, 2004 (Advances in Industrial Control)
- [Kuč72] KUČERA, V.: The Discrete Riccati Equation of Optimal Control. In: *Kybernetika* 8 (1972), Nr. 5, S. 430–447
- [Kuč73] KUČERA, V.: A Review of the Matrix Riccati Equation. In: *Kybernetika* 9 (1973), Nr. 5, S. 42–61
- [KWD13] KAMALAPURKAR, R.; WALTERS, P.; DIXON, W.: Concurrent Learning-Based Approximate Optimal Regulation. In: *52nd IEEE Conference on Decision and Control (CDC)*, 2013, S. 6256–6261

- [KWD16] KAMALAPURKAR, R.; WALTERS, P.; DIXON, W. E.: Model-Based Reinforcement Learning for Approximate Optimal Regulation. In: *Automatica* 64 (2016), S. 94–104
- [Lan97] LANDELIUS, T.: *Reinforcement Learning and Distributed Local Model Synthesis*. Linköping, Linköping University, Dissertation, 1997
- [LCL⁺19] LI, J.; CHAI, T.; LEWIS, F. L.; DING, Z.; JIANG, Y.: Off-Policy Interleaved Q-Learning: Optimal Control for Affine Nonlinear Discrete-Time Systems. In: *IEEE Transactions on Neural Networks and Learning Systems* 30 (2019), Nr. 5, S. 1308–1320
- [Lév09] LÉVINE, J.: *Analysis and Control of Nonlinear Systems. A Flatness-Based Approach*. Berlin, Heidelberg: Springer, 2009
- [LG19] LI, G.; GÖRGES, D.: Ecological Adaptive Cruise Control and Energy Management Strategy for Hybrid Electric Vehicles Based on Heuristic Dynamic Programming. In: *IEEE Transactions on Intelligent Transportation Systems* 20 (2019), Nr. 9, S. 3526–3535
- [LG20] LI, G.; GÖRGES, D.: Ecological Adaptive Cruise Control for Vehicles With Step-Gear Transmission Based on Reinforcement Learning. In: *IEEE Transactions on Intelligent Transportation Systems* 21 (2020), Nr. 11, S. 4895–4905
- [LGM20] LI, G.; GÖRGES, D.; MU, C.: Integrated Adaptive Dynamic Programming for Data-Driven Optimal Controller Design. In: *Neurocomputing* 403 (2020), S. 143–152
- [LHP⁺16] LILICRAP, T. P.; HUNT, J. J.; PRITZEL, A.; HEES, N.; EREZ, T.; TASSA, Y.; SILVER, D.; WIERSTRA, D.: Continuous Control with Deep Reinforcement Learning. In: *4th International Conference on Learning Representations*, 2016, S. 1–14
- [Lib12] LIBERZON, D.: *Calculus of Variations and Optimal Control Theory*. Princeton: Princeton University Press, 2012
- [LK98] LIN, J.-S.; KANELAKOPOULOS, I.: Nonlinearities Enhance Parameter Convergence in Output-Feedback Systems. In: *IEEE Transactions on Automatic Control* 43 (1998), Nr. 2, S. 204–222
- [LK99] LIN, J.-S.; KANELAKOPOULOS, I.: Nonlinearities Enhance Parameter Convergence in Strict Feedback Systems. In: *IEEE Transactions on Automatic Control* 44 (1999), Nr. 1, S. 89–94
- [LL04] LEE, J. M.; LEE, J. H.: Approximate Dynamic Programming Strategies and Their Applicability for Process Control: A Review and Future Directions. In: *International Journal of Control Automation and Systems* 2 (2004), S. 263–278

- [LLHW16] LUO, B.; LIU, D.; HUANG, T.; WANG, D.: Model-Free Optimal Tracking Control via Critic-Only Q-Learning. In: *IEEE Transactions on Neural Networks and Learning Systems* 27 (2016), Nr. 10, S. 2134–2144
- [LLSX14] LU, R.; LI, Z.; SU, C.-Y.; XUE, A.: Development and Learning Control of a Human Limb with a Rehabilitation Exoskeleton. In: *IEEE Transactions on Industrial Electronics* 61 (2014), Nr. 7, S. 3776–3785
- [LLW14] LIU, D.; LI, H.; WANG, D.: Online Synchronous Approximate Optimal Learning Algorithm for Multi-Player Non-Zero-Sum Games with Unknown Dynamics. In: *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 44 (2014), Nr. 8, S. 1015–1027
- [LLW⁺17] LUO, B.; LIU, D.; WU, H.-N.; WANG, D.; LEWIS, F. L.: Policy Gradient Adaptive Dynamic Programming for Data-Based Optimal Control. In: *IEEE Transactions on Cybernetics* 47 (2017), Nr. 10, S. 3341–3354
- [LNY⁺15] LV, Y.; NA, J.; YANG, Q.; WU, X.; GUO, Y.: Online Adaptive Optimal Control for Continuous-Time Nonlinear Systems with Completely Unknown Dynamics. In: *International Journal of Control* 89 (2015), Nr. 1, S. 99–112
- [LP03] LAGOUDAKIS, M. G.; PARR, R.: Least-Squares Policy Iteration. In: *Journal of Machine Learning Research* 4 (2003), S. 1107–1149
- [LPC12] LEE, J. Y.; PARK, J. B.; CHOI, Y. H.: Integral Q-Learning and Explorized Policy Iteration for Adaptive Optimal Control of Continuous-Time Linear Systems. In: *Automatica* 48 (2012), Nr. 11, S. 2850–2859
- [LPC15] LEE, J. Y.; PARK, J. B.; CHOI, Y. H.: Integral Reinforcement Learning for Continuous-Time Input-Affine Nonlinear Systems with Simultaneous Invariant Explorations. In: *IEEE Transactions on Neural Networks and Learning Systems* 26 (2015), Nr. 5, S. 916–932
- [LS17] LEE, J. Y.; SUTTON, R. S.: Integral Policy Iterations for Reinforcement Learning Problems in Continuous Time and Space. In: *arXiv preprint 1705.03520* (2017)
- [Lun20a] LUNZE, J.: *Regelungstechnik 1*. 12., überarbeitete Auflage. Berlin, Heidelberg: Springer Vieweg, 2020
- [Lun20b] LUNZE, J.: *Regelungstechnik 2*. 10., überarbeitete und aktualisierte Auflage. Berlin, Heidelberg: Springer Vieweg, 2020
- [LV09] LEWIS, F.; VRABIE, D.: Reinforcement Learning and Adaptive Dynamic Programming for Feedback Control. In: *IEEE Circuits and Systems Magazine* 9 (2009), Nr. 3, S. 32–50
- [LVS12] LEWIS, F. L.; VRABIE, D. L.; SYRMOUS, V. L.: *Optimal Control*. 3. Auflage. Hoboken: Wiley, 2012

- [LVV12] LEWIS, F. L.; VRABIE, D.; VAMVOUDAKIS, K. G.: Reinforcement Learning and Feedback Control: Using Natural Decision Methods to Design Optimal Adaptive Controllers. In: *IEEE Control Systems Magazine* 32 (2012), Nr. 6, S. 76–105
- [LW14] LIU, D.; WEI, Q.: Policy Iteration Adaptive Dynamic Programming Algorithm for Discrete-Time Nonlinear Systems. In: *IEEE Transactions on Neural Networks and Learning Systems* 25 (2014), Nr. 3, S. 621–634
- [LWW⁺17] LIU, D.; WEI, Q.; WANG, D.; YANG, X.; LI, H.: *Adaptive Dynamic Programming with Applications in Optimal Control*. Cham: Springer International Publishing, 2017 (Advances in Industrial Control)
- [LYD17] LI, J.; YUAN, D.; DING, Z.: Optimal Tracking Control for Discrete-Time Systems by Model-Free Off-Policy Q-Learning Approach. In: *11th Asian Control Conference, 2017*, S. 7–12
- [LYWW15] LIU, D.; YANG, X.; WANG, D.; WEI, Q.: Reinforcement-Learning-Based Robust Controller Design for Continuous-Time Uncertain Nonlinear Systems Subject to Input Constraints. In: *IEEE Transactions on Cybernetics* 45 (2015), Nr. 7, S. 1372–1385
- [Mar84] MAREELS, I.: Sufficiency of Excitation. In: *Systems & Control Letters* 5 (1984), Nr. 3, S. 159–163
- [MB90] MATHÉLIN, M. de; BODSON, M.: Frequency Domain Conditions for Parameter Convergence in Multivariable Recursive Identification. In: *Automatica* 26 (1990), Nr. 4, S. 757–767
- [MBK19] MACCLUER, B. D.; BOURDON, P.; KRIETE, T.: *Differential Equations: Techniques, Theory, and Applications*. Providence: American Mathematical Society, 2019
- [MBTL12] MANNAVA, A.; BALAKRISHNAN, S. N.; TANG, L.; LANDERS, R. G.: Optimal Tracking Control of Motion Systems. In: *IEEE Transactions on Control Systems Technology* 20 (2012), Nr. 6, S. 1548–1558
- [MCLS02] MURRAY, J. J.; COX, C. J.; LENDARIS, G. G.; SAEKS, R.: Adaptive Dynamic Programming. In: *IEEE Transactions on Systems, Man and Cybernetics, Part C: Applications and Reviews* 32 (2002), Nr. 2, S. 140–153
- [Mey00] MEYER, C. D.: *Matrix Analysis and Applied Linear Algebra*. Philadelphia: SIAM, 2000
- [MKS⁺13] MNIH, V.; KAVUKCUOGLU, K.; SILVER, D.; GRAVES, A.; ANTONOGLU, I.; WIERSTRA, D.; RIEDMILLER, M.: Playing Atari with Deep Reinforcement Learning. In: *arXiv preprint 1312.5602* (2013), S. 1–9

- [MKS⁺15] MNIH, V.; KAVUKCUOGLU, K.; SILVER, D.; RUSU, A. A.; VENESS, J.; BELLEMARE, M. G.; GRAVES, A.; RIEDMILLER, M.; FIDJELAND, A. K.; OSTROVSKI, G.; PETERSEN, S.; BEATTIE, C.; SADIK, A.; ANTONOGLU, I.; KING, H.; KUMARAN, D.; WIERSTRA, D.; LEGG, S.; HASSABIS, D.: Human-Level Control Through Deep Reinforcement Learning. In: *Nature* 518 (2015), S. 529–533
- [ML13] MODARES, H.; LEWIS, F. L.: Online Solution to the Linear Quadratic Tracking Problem of Continuous-Time Systems Using Reinforcement Learning. In: *52nd IEEE Conference on Decision and Control*, 2013, S. 3851–3856
- [ML14a] MODARES, H.; LEWIS, F. L.: Linear Quadratic Tracking Control of Partially-Unknown Continuous-Time Systems Using Reinforcement Learning. In: *IEEE Transactions on Automatic Control* 59 (2014), Nr. 11, S. 3051–3056
- [ML14b] MODARES, H.; LEWIS, F. L.: Optimal Tracking Control of Nonlinear Partially-Unknown Constrained-Input Systems Using Integral Reinforcement Learning. In: *Automatica* 50 (2014), Nr. 7, S. 1780–1792
- [MM70] MENDEL, J. M.; MCLAREN, R. W.: Reinforcement-Learning Control and Pattern Recognition Systems. In: MENDEL, J. M. (Hrsg.); FU, K. S. (Hrsg.): *Adaptive, Learning and Pattern Recognition Systems*. New York: Academic Press, 1970, S. 287–318
- [MMP⁺10] MICHALOS, G.; MAKRIS, S.; PAPAKOSTAS, N.; MOURTZIS, D.; CHRYSOLOURIS, G.: Automotive Assembly Technologies Review: Challenges and Outlook for a Flexible and Adaptive Approach. In: *CIRP Journal of Manufacturing Science and Technology* 2 (2010), Nr. 2, S. 81–91
- [MN77] MORGAN, A. P.; NARENDRA, K. S.: On the Uniform Asymptotic Stability of Certain Linear Nonautonomous Differential Equations. In: *SIAM Journal on Control and Optimization* 15 (1977), Nr. 1, S. 5–24
- [Mor78] MORÉ, J. J.: The Levenberg-Marquardt Algorithm: Implementation and Theory. In: WATSON, G. A. (Hrsg.): *Numerical Analysis*. Berlin, Heidelberg: Springer, 1978, S. 105–116
- [MRLP16] MODARES, H.; RANATUNGA, I.; LEWIS, F. L.; POPA, D. O.: Optimized Assistive Human-Robot Interaction Using Reinforcement Learning. In: *IEEE Transactions on Cybernetics* 46 (2016), Nr. 3, S. 655–667
- [MSV08] MOARREF, M.; SAADAT, M.; VOSSOUGH, G.: Mechatronic Design and Position Control of a Novel Ball and Plate System. In: *16th Mediterranean Conference on Control and Automation*, 2008, S. 1071–1076
- [MSWS17] MU, C.; SUN, C.; WANG, D.; SONG, A.: Adaptive Tracking Control for a Class of Continuous-Time Uncertain Nonlinear Systems Using the Approximate Solution of HJB Equation. In: *Neurocomputing* 260 (2017), S. 432–442

- [Mun06] MUNOS, R.: Policy Gradient in Continuous Time. In: *Journal of Machine Learning Research* 7 (2006), Nr. 27, S. 771–791
- [NA87] NARENDRA, K. S.; ANNASWAMY, A. M.: Persistent Excitation in Adaptive Systems. In: *International Journal of Control* 45 (1987), Nr. 1, S. 127–160
- [NA05] NARENDRA, K. S.; ANNASWAMY, A. M.: *Stable Adaptive Systems*. 2. Auflage. Mineola, New York: Dover Publications, 2005
- [NC15] NA, X.; COLE, D. J.: Game-Theoretic Modeling of the Steering Interaction Between a Human Driver and a Vehicle Collision Avoidance Controller. In: *IEEE Transactions on Human-Machine Systems* 45 (2015), Nr. 1, S. 25–38
- [NCH08] NG, L.; CLARK, C. M.; HUISSOON, J. P.: Reinforcement Learning of Adaptive Longitudinal Vehicle Control for Dynamic Collaborative Driving. In: *IEEE Intelligent Vehicles Symposium*, 2008, S. 907–912
- [Ngu18] NGUYEN, N. T.: *Model-Reference Adaptive Control*. Cham: Springer International Publishing, 2018
- [Niv09] NIV, Y.: Reinforcement Learning in the Brain. In: *Journal of Mathematical Psychology* 53 (2009), Nr. 3, S. 139–154
- [NKJS04] NG, A. Y.; KIM, H. J.; JORDAN, M. I.; SASTRY, S.: Autonomous Helicopter Flight via Reinforcement Learning. In: THRUN, S. (Hrsg.); SAUL, L. K. (Hrsg.); SCHÖLKOPF, B. (Hrsg.): *Advances in Neural Information Processing Systems*. Cambridge: MIT Press, 2004, S. 799–806
- [NLW⁺19] NA, J.; LI, G.; WANG, B.; HERRMANN, G.; ZHAN, S.: Robust Optimal Control of Wave Energy Converters Based on Adaptive Dynamic Programming. In: *IEEE Transactions on Sustainable Energy* 10 (2019), Nr. 2, S. 961–970
- [NMH⁺15] NA, J.; MAHYUDDIN, M. N.; HERRMANN, G.; REN, X.; BARBER, P.: Robust Adaptive Finite-Time Parameter Estimation and Control for Robotic Systems. In: *International Journal of Robust and Nonlinear Control* 25 (2015), Nr. 16, S. 3045–3071
- [NP96] NEVISTIC, V.; PRIMBS, J. A.: *Constrained Nonlinear Optimal Control: A Converse HJB Approach*. Pasadena, California Institute of Technology, Technischer Bericht, 1996
- [NW06] NOCEDAL, J.; WRIGHT, S. J.: *Numerical Optimization*. 2. Auflage. New York: Springer Science and Business Media LLC, 2006
- [Pav27] PAVLOV, I. P.: *Conditioned Reflexes: An Investigation of the Physiological Activity of the Cerebral Cortex*. Übersetzt durch G. V. Anrep. London: Oxford University Press, 1927

- [Pla19] PLATTFORM INDUSTRIE 4.0; BUNDESMINISTERIUM FÜR WIRTSCHAFT UND ENERGIE (Hrsg.): *Technologieszenario Künstliche Intelligenz in der Industrie 4.0*. 2019
- [PLB15] PAPAGEORGIOU, M.; LEIBOLD, M.; BUSS, M.: *Optimierung: Statische, dynamische, stochastische Verfahren für die Anwendung*. 4. Auflage. Berlin, Heidelberg: Springer, 2015
- [PMB13] PASCANU, R.; MIKOLOV, T.; BENGIO, Y.: On the Difficulty of Training Recurrent Neural Networks. In: *Proceedings of the 30th International Conference on Machine Learning*, 2013, S. 1310–1318
- [PMK99] PESHKIN, L.; MEULEAU, N.; KAEHLING, L. P.: Learning Policies with External Memory. In: *Proceedings of the 16th International Conference on Machine Learning*, 1999, S. 307–314
- [Pow11] POWELL, P. D.: Calculating Determinants of Block Matrices. In: *arXiv preprint 1112.4379* (2011), S. 1–11
- [Pra17] PRALY, L.: *Convergence of the Gradient Algorithm for Linear Regression Models in the Continuous and Discrete Time Cases*. Paris, Mines ParisTech, Technischer Bericht, 2017
- [PRH19] PUCETTI, L.; RATHGEBER, C.; HOHMANN, S.: Actor-Critic Reinforcement Learning for Linear Longitudinal Output Control of a Road Vehicle. In: *IEEE Intelligent Transportation Systems Conference*, 2019, S. 2907–2913
- [PSA17] PADOAN, A.; SCARCIOTTI, G.; ASTOLFI, A.: A Geometric Characterization of the Persistence of Excitation Condition for the Solutions of Autonomous Systems. In: *IEEE Transactions on Automatic Control* 62 (2017), Nr. 11, S. 5666–5677
- [PT12] PIETQUIN, O.; TANGO, F.: A Reinforcement Learning Approach to Optimize the Longitudinal Behavior of a Partial Autonomous Driving Assistance System. In: *Proceedings of the 20th European Conference on Artificial Intelligence*, 2012, S. 987–992
- [QZL13] QIN, C.; ZHANG, H.; LUO, Y.: Online Optimal Tracking Control of Continuous-Time Linear Systems with Unknown Dynamics by Using Adaptive Dynamic Programming. In: *International Journal of Control* 87 (2013), Nr. 5, S. 1000–1009
- [QZLY19] QU, Q.; ZHANG, H.; LUO, C.; YU, R.: Robust Control Design for Multi-Player Nonlinear Systems with Input Disturbances via Adaptive Dynamic Programming. In: *Neurocomputing* 334 (2019), S. 1–10
- [Rec19] RECHT, B.: A Tour of Reinforcement Learning: The View from Continuous Control. In: *Annual Review of Control, Robotics, and Autonomous Systems* 2 (2019), S. 253–279

- [RPS07] RIEDMILLER, M.; PETERS, J.; SCHAAL, S.: Evaluation of Policy Gradient Methods and Variants on the Cart-Pole Benchmark. In: *IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning*, 2007, S. 254–261
- [RW72] RESCORLA, R. A.; WAGNER, A. R.: A Theory of Pavlovian Conditioning: Variations on the Effectiveness of Reinforcement and Non-Reinforcement. In: BLACK, A. H. (Hrsg.); PROKASY, W. F. (Hrsg.): *Classical Conditioning II: Current Research and Theory*. New York: Appleton-Century-Crofts, 1972, S. 64–99
- [SB89] SASTRY, S.; BODSON, M.: *Adaptive Control: Stability, Convergence, and Robustness*. Englewood Cliffs: Prentice-Hall, 1989
- [SB18] SUTTON, R. S.; BARTO, A. G.: *Reinforcement Learning: An Introduction*. 2. Auflage. Cambridge, Mass.: MIT Press, 2018
- [SBW92] SUTTON, R. S.; BARTO, A. G.; WILLIAMS, R. J.: Reinforcement Learning is Direct Adaptive Optimal Control. In: *IEEE Control Systems Magazine* 12 (1992), Nr. 2, S. 19–22
- [SCN⁺04] SMETS, I. Y.; CLAES, J. E.; NOVEMBER, E. J.; BASTIN, G. P.; VAN IMPE, J. F.: Optimal Adaptive Control of (Bio)chemical Reactors: Past, Present and Future. In: *Journal of Process Control* 14 (2004), Nr. 7, S. 795–805
- [Sco04] SCOTT, S. H.: Optimal Feedback Control and the Neural Basis of Volitional Motor Control. In: *Nature Reviews Neuroscience* 5 (2004), Nr. 7, S. 532–546
- [SDM97] SCHULTZ, W.; DAYAN, P.; MONTAGUE, R.: A Neural Substrate of Prediction and Reward. In: *Science* 275 (1997), Nr. 5306, 1593–1599
- [SHM⁺16] SILVER, D.; HUANG, A.; MADDISON, C. J.; GUEZ, A.; SIFRE, L.; VAN DEN DRIESCHE, G.; SCHRITTWIESER, J.; ANTONOGLU, I.; PANNEERSHELVA, V.; LANCTOT, M.; DIELEMAN, S.; GREWE, D.; NHAM, J.; KALCHBRENNER, N.; SUTSKEVER, I.; LILLICRAP, T.; LEACH, M.; KAVUKCUOGLU, K.; GRAEPEL, T.; HASSABIS, D.: Mastering the Game of Go with Deep Neural Networks and Tree Search. In: *Nature* 529-503 (2016), Nr. 7587, 484–503
- [SLH⁺14] SILVER, D.; LEVER, G.; HEES, N.; DEGRIS, T.; WIERSTRA, D.; RIEDMILLER, M.: Deterministic Policy Gradient Algorithms. In: *Proceedings of the 31st International Conference on Machine Learning*, 2014
- [SLW17] SONG, R.; LEWIS, F. L.; WEI, Q.: Off-Policy Integral Reinforcement Learning Method to Solve Nonlinear Continuous-Time Multiplayer Nonzero-Sum Games. In: *IEEE Transactions on Neural Networks and Learning Systems* 28 (2017), Nr. 3, S. 704–713

- [SMC⁺11] STOCKAR, S.; MARANO, V.; CANOVA, M.; RIZZONI, G.; GUZZELLA, L.: Energy-Optimal Control of Plug-in Hybrid Electric Vehicles for Real-World Driving Cycles. In: *IEEE Transactions on Vehicular Technology* 60 (2011), Nr. 7, S. 2949–2962
- [Spi68] SPIEGEL, M. R.: *Mathematical Handbook of Formulas and Tables*. New York, St. Louis: McGraw-Hill, 1968
- [SSS⁺17] SILVER, D.; SCHRITTWIESER, J.; SIMONYAN, K.; ANTONOGLIOU, I.; HUANG, A.; GUEZ, A.; HUBERT, T.; BAKER, L.; LAI, M.; BOLTON, A.; CHEN, Y.; LILICRAP, T.; HUI, F.; SIFRE, L.; VAN DEN DRIESSCHE, G.; GRAEPEL, T.; HASSABIS, D.: Mastering the Game of Go Without Human Knowledge. In: *Nature* 550 (2017), Nr. 7676, 354–371
- [SSW18] SHI, W.; SONG, S.; WU, C.: High-Level Tracking of Autonomous Underwater Vehicles Based on Pseudo Averaged Q-Learning. In: *IEEE International Conference on Systems, Man, and Cybernetics*, 2018, S. 4138–4143
- [Sut88] SUTTON, R. S.: Learning to Predict by the Methods of Temporal Differences. In: *Machine Learning* 3 (1988), Nr. 1, S. 9–44
- [SWL19] SONG, R.; WEI, Q.; LI, Q. F.: *Studies in Systems, Decision and Control*. Bd. 166: *Adaptive Dynamic Programming: Single and Multiple Controllers*. New York, Berlin: Springer, 2019
- [Tao03] TAO, G.: *Adaptive Control Design and Analysis*. Hoboken: John Wiley and Sons Inc., 2003
- [TBBH10] TELEKE, S.; BARAN, M. E.; BHATTACHARYA, S.; HUANG, A. Q.: Optimal Control of Battery Energy Storage for Wind Farm Dispatching. In: *IEEE Transactions on Energy Conversion* 25 (2010), Nr. 3, S. 787–794
- [TCTH19] TANG, D.; CHEN, L.; TIAN, Z. F.; HU, E.: Modified Value-Function-Approximation for Synchronous Policy Iteration with Single-Critic Configuration for Nonlinear Optimal Control. In: *International Journal of Control* 94 (2019), Nr. 5, S. 1–13
- [Tho11] THORNDIKE, E. L.: *Animal Intelligence: Experimental Studies*. New York: The Macmillan Company, 1911
- [Tod04] TODOROV, E.: Optimality Principles in Sensorimotor Control. In: *Nature Neuroscience* 7 (2004), Nr. 9, S. 907–915
- [Vam15] VAMVOUDAKIS, K. G.: Non-Zero Sum Nash Q-Learning for Unknown Deterministic Continuous-Time Linear Systems. In: *Automatica* 61 (2015), S. 274–281
- [Vam16] VAMVOUDAKIS, K. G.: Optimal Trajectory Output Tracking Control with a Q-Learning Algorithm. In: *American Control Conference (ACC)*, 2016, S. 5752–5757

- [Vam17] VAMVOUDAKIS, K. G.: Q-Learning for Continuous-Time Linear Systems: A Model-Free Infinite Horizon Optimal Control Approach. In: *Systems & Control Letters* 100 (2017), S. 14–20
- [van97] VAN NIEUWSTADT, M. J.: *Trajectory Generation for Nonlinear Control Systems*. Pasadena, California Institute of Technology, Dissertation, 1997
- [van12] VAN HASSELT, H.: Reinforcement Learning in Continuous State and Action Spaces. In: WIERING, M. (Hrsg.); VAN OTTERLO, M. (Hrsg.): *Reinforcement Learning: State-of-the-Art*. Berlin, Heidelberg: Springer, 2012, S. 207–251
- [VL09] VRABIE, D.; LEWIS, F.: Neural Network Approach to Continuous-Time Direct Adaptive Optimal Control for Partially Unknown Nonlinear Systems. In: *Neural Networks* 22 (2009), Nr. 3, S. 237–246
- [VL10] VAMVOUDAKIS, K. G.; LEWIS, F. L.: Online Actor-Critic Algorithm to Solve the Continuous-Time Infinite Horizon Optimal Control Problem. In: *Automatica* 46 (2010), Nr. 5, S. 878–888
- [VL11] VAMVOUDAKIS, K. G.; LEWIS, F. L.: Multi-Player Non-Zero-Sum Games: On-line Adaptive Learning Solution of Coupled Hamilton-Jacobi Equations. In: *Automatica* 47 (2011), Nr. 8, S. 1556–1569
- [VLV13] VRABIE, D. L.; LEWIS, F. L.; VAMVOUDAKIS, K. G.: *IET Control Engineering*. Bd. 81: *Optimal Adaptive Control and Differential Games by Reinforcement Learning Principles*. London: Institution of Electrical Engineers, 2013
- [VMH16] VAMVOUDAKIS, K. G.; MIRANDA, M. F.; HESPANHA, J. P.: Asymptotically Stable Adaptive-Optimal Control Algorithm With Saturating Actuators and Related Persistence of Excitation. In: *IEEE Transactions on Neural Networks and Learning Systems* 27 (2016), Nr. 11, S. 2386–2398
- [VMKL17] VAMVOUDAKIS, K. G.; MODARES, H.; KIUMARSI, B.; LEWIS, F. L.: Game Theory-Based Control System Algorithms with Real-Time Reinforcement Learning: How to Solve Multiplayer Games Online. In: *IEEE Control Systems Magazine* 37 (2017), Nr. 1, S. 33–52
- [VPAKL09] VRABIE, D.; PASTRAVANU, O.; ABU-KHALAF, M.; LEWIS, F. L.: Adaptive Optimal Control for Continuous-Time Linear Systems Based on Policy Iteration. In: *Automatica* 45 (2009), Nr. 2, S. 477–484
- [VVL09b] VRABIE, D.; VAMVOUDAKIS, K.; LEWIS, F.: Adaptive Optimal Controllers Based on Generalized Policy Iteration in a Continuous-Time Framework. In: *17th Mediterranean Conference on Control and Automation*, 2009, S. 1402–1409
- [Wan20] WANG, D.: Robust Policy Learning Control of Nonlinear Plants With Case Studies for a Power System Application. In: *IEEE Transactions on Industrial Informatics* 16 (2020), Nr. 3, S. 1733–1741

- [Wat13] WATSON, J. B.: Psychology as the Behaviorist Views It. In: *Psychological Review* 20 (1913), Nr. 2, S. 158–177
- [Wat89] WATKINS, C. J. C. H.: *Learning from Delayed Rewards*, King's College, Cambridge, Dissertation, 1989
- [WD92] WATKINS, C. J.; DAYAN, P.: Q-Learning. In: *Machine Learning* 8 (1992), S. 279–292
- [Wei85] WEIERSTRASS, K.: Über die analytische Darstellbarkeit sogenannter willkürlicher Functionen einer reellen Veränderlichen. In: *Sitzungsberichte der Königlich Preussischen Akademie der Wissenschaften zu Berlin*. 1885, S. 633–639
- [Wer77] WERBOS, P. J.: Advanced Forecasting Methods for Global Crisis Warning and Models of Intelligence. In: *General Systems* Bd. XXII, 1977, S. 25–38
- [Wer92] WERBOS, P. J.: Approximate Dynamic Programming for Real-Time Control and Neural Modeling. In: *Handbook of Intelligent Control* (1992), S. 493–526
- [Wer99] WERBOS, P. J.: Stable Adaptive Control Using New Critic Designs. In: *9th Workshop on Virtual Intelligence/Dynamic Neural Networks*. 1999
- [Wer13] WERBOS, P. J.: Reinforcement Learning and Approximate Dynamic Programming (RLADP) – Foundations, Common Misconceptions and Challenges Ahead. In: LEWIS, F. L. (Hrsg.); LIU, D. (Hrsg.): *Reinforcement Learning and Approximate Dynamic Programming for Feedback Control*. Piscataway, Hoboken: John Wiley and Sons Inc., 2013, S. 3–30
- [WHL17] WANG, D.; HE, H.; LIU, D.: Adaptive Critic Nonlinear Robust Control: A Survey. In: *IEEE Transactions on Cybernetics* 47 (2017), Nr. 10, S. 3429–3451
- [WHQ20] WANG, D.; HA, M.; QIAO, J.: Self-Learning Optimal Regulation for Discrete-Time Nonlinear Systems Under Event-Driven Formulation. In: *IEEE Transactions on Automatic Control* 65 (2020), Nr. 3, S. 1272–1279
- [WL12] WU, H.-N.; LUO, B.: Neural Network Based Online Simultaneous Policy Update Algorithm for Solving the HJI Equation in Nonlinear H_∞ Control. In: *IEEE Transactions on Neural Networks and Learning Systems* 23 (2012), Nr. 12, S. 1884–1895
- [WLL14] WANG, D.; LIU, D.; LI, H.: Policy Iteration Algorithm for Online Design of Robust Control for a Class of Continuous-Time Nonlinear Systems. In: *IEEE Transactions on Automation Science and Engineering* 11 (2014), Nr. 2, S. 627–632
- [WLLS17] WEI, Q.; LIU, D.; LIN, Q.; SONG, R.: Discrete-Time Optimal Control via Local Policy Iteration Adaptive Dynamic Programming. In: *IEEE Transactions on Cybernetics* 47 (2017), Nr. 10, S. 3367–3379

- [WLMZ18] WANG, D.; LIU, D.; MU, C.; ZHANG, Y.: Neural Network Learning and Robust Stabilization of Nonlinear Systems With Dynamic Uncertainties. In: *IEEE Transactions on Neural Networks and Learning Systems* 29 (2018), Nr. 4, S. 1342–1351
- [WLZZ16] WANG, D.; LIU, D.; ZHANG, Q.; ZHAO, D.: Data-Based Adaptive Critic Designs for Nonlinear Robust Optimal Control with Uncertain Dynamics. In: *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 46 (2016), Nr. 11, S. 1544–1555
- [WSH⁺16] WANG, Z.; SCHAUL, T.; HESSEL, M.; VAN HASSELT, H.; LANCTOT, M.; DE FREITAS, N.: Dueling Network Architectures for Deep Reinforcement Learning. In: *Proceedings of the 33rd International Conference on Machine Learning* Bd. 48, 2016, 1995–2003
- [WXL⁺14] WANG, J.; XU, X.; LIU, D.; SUN, Z.; CHEN, Q.: Self-Learning Cruise Control Using Kernel-Based Least Squares Policy Iteration. In: *IEEE Transactions on Control Systems Technology* 22 (2014), Nr. 3, S. 1078–1087
- [WY18] WANG, J.-S.; YANG, G.-H.: Output-Feedback Control of Unknown Linear Discrete-Time Systems With Stochastic Measurement and Process Noise via Approximate Dynamic Programming. In: *IEEE Transactions on Cybernetics* 48 (2018), Nr. 7, S. 1977–1988
- [WZL09] WANG, F.-Y.; ZHANG, H.; LIU, D.: Adaptive Dynamic Programming: An Introduction. In: *IEEE Computational Intelligence Magazine* 4 (2009), Nr. 2, S. 39–47
- [WZL16] WANG, T.; ZHANG, H.; LUO, Y.: Infinite-Time Stochastic Linear Quadratic Optimal Control for Unknown Discrete-Time Systems Using Adaptive Dynamic Programming Approach. In: *Neurocomputing* 171 (2016), S. 379–386
- [WZLD15] WANG, B.; ZHAO, D.; LI, C.; DAI, Y.: Design and Implementation of an Adaptive Cruise Control System Based on Supervised Actor-Critic Learning. In: *5th International Conference on Information Science and Technology*, 2015, S. 243–248
- [XZLJ16] XIAO, G.; ZHANG, H.; LUO, Y.; JIANG, H.: Data-Driven Optimal Tracking Control for a Class of Affine Non-Linear Continuous-Time Systems with Completely Unknown Dynamics. In: *IET Control Theory & Applications* 10 (2016), Nr. 6, S. 700–710
- [YDZY20] YANG, Y.; DONG, J.; ZHANG, S.; YIN, Y.: Data-Driven Nonzero-Sum Game for Discrete-Time Systems Using Off-Policy Reinforcement Learning. In: *IEEE Access* 4 (2020), Nr. 8, S. 14074–14088

- [YHL13] YANG, X.; HUANG, Y.; LIU, D.: Neural-Network-Based Online Optimal Control for Uncertain Non-Linear Continuous-Time Systems with Control Constraints. In: *IET Control Theory & Applications* 7 (2013), Nr. 17, S. 2037–2047
- [YLLL16] YANG, X.; LIU, D.; LUO, B.; LI, C.: Data-Based Robust Adaptive Control for a Class of Unknown Nonlinear Constrained-Input Systems via Integral Reinforcement Learning. In: *Information Sciences* 369 (2016), S. 731–747
- [YSH⁺17] YU, R.; SHI, Z.; HUANG, C.; LI, T.; MA, Q.: Deep Reinforcement Learning Based Optimal Trajectory Tracking Control of Autonomous Underwater Vehicle. In: *36th Chinese Control Conference*, 2017, S. 4958–4965
- [YWM⁺19] YANG, Y.; WANG, L.; MODARES, H.; DING, D.; YIN, Y.; WUNSCH, D.: Data-Driven Integral Reinforcement Learning for Continuous-Time Non-Zero-Sum Games. In: *IEEE Access* 7 (2019), S. 82901–82912
- [ZCL13] ZHANG, H.; CUI, L.; LUO, Y.: Near-Optimal Control for Nonzero-Sum Differential Games of Continuous-Time Nonlinear Systems Using Single-Network ADP. In: *IEEE Transactions on Cybernetics* 43 (2013), Nr. 1, S. 206–216
- [ZCZL11] ZHANG, H.; CUI, L.; ZHANG, X.; LUO, Y.: Data-Driven Robust Approximate Optimal Tracking Control for Unknown General Nonlinear Systems Using Adaptive Dynamic Programming Method. In: *IEEE Transactions on Neural Networks* 22 (2011), Nr. 12, S. 2226–2236
- [ZD98] ZHOU, K.; DOYLE, J. C.: *Essentials of Robust Control*. Upper Saddle River: Prentice Hall, 1998
- [ZDJ14] ZARGARZADEH, H.; DIERKS, T.; JAGANNATHAN, S.: Adaptive Neural Network-Based Optimal Control of Nonlinear Continuous-Time Systems in Strict-Feedback Form. In: *International Journal of Adaptive Control and Signal Processing* 28 (2014), S. 305–324
- [ZDJ15] ZARGARZADEH, H.; DIERKS, T.; JAGANNATHAN, S.: Optimal Control of Nonlinear Continuous-Time Systems in Strict-Feedback Form. In: *IEEE Transactions on Neural Networks and Learning Systems* 26 (2015), Nr. 10, S. 2535–2549
- [Zim84] ZIMMER, K.: Determination of Proper Output Variables for the Decoupling of Nonlinear Systems. In: *Electronics Letters* 20 (1984), Nr. 25–26, S. 1052–1053
- [ZQJL14] ZHANG, H.; QIN, C.; JIANG, B.; LUO, Y.: Online Adaptive Policy Learning Algorithm for H-infinity State Feedback Control of Unknown Affine Nonlinear Discrete-Time Systems. In: *IEEE Transactions on Cybernetics* 44 (2014), Nr. 12, S. 2706–2718
- [ZWL08] ZHANG, H.; WEI, Q.; LUO, Y.: A Novel Infinite-Time Optimal Tracking Control Scheme for a Class of Discrete-Time Nonlinear Systems via the Greedy HDP Iteration Algorithm. In: *IEEE Transactions on Systems, Man, and Cybernetics. Part B: Cybernetics* 38 (2008), Nr. 4, S. 937–942

-
- [ZZWZ16] ZHAO, D.; ZHANG, Q.; WANG, D.; ZHU, Y.: Experience Replay for Optimal Control of Nonzero-Sum Game Systems with Unknown Dynamics. In: *IEEE Transactions on Cybernetics* 46 (2016), Nr. 3, S. 854–865
- [ZZXS17] ZHANG, K.; ZHANG, H.; XIAO, G.; SU, H.: Tracking Control Optimization Scheme of Continuous-Time Nonlinear System via Online Single Network Adaptive Critic Design Method. In: *Neurocomputing* 251 (2017), S. 127–135

Eigene Veröffentlichungen und Tagungsbeiträge

- [BKI21] BÜHRLE, E.; KÖPF, F.; INGA, J.; HOHMANN, S.: Adaptive Optimal Trajectory Tracking Control of Continuous-Time Systems. In: *European Control Conference (ECC)*, 2021, S. 1987–1994
- [IBKH20] INGA, J.; BISCHOFF, E.; KÖPF, F.; HOHMANN, S.: Inverse Dynamic Games Based on Maximum Entropy Inverse Reinforcement Learning. In: *arXiv preprint 1911.07503v2* (2020), S. 1–8
- [IKFH17] INGA, J.; KÖPF, F.; FLAD, M.; HOHMANN, S.: Individual Human Behavior Identification Using an Inverse Reinforcement Learning Method. In: *IEEE International Conference on Systems, Man and Cybernetics (SMC)*, 2017, S. 99–104
- [IRKH17] INGA, J.; ROTHFUSS, S.; KÖPF, F.; HOHMANN, S.: *Inverse Optimierung in linear-quadratischen dynamischen Spielen*. Günzburg, 2017 (GMA-Fachausschuss 1.50 Workshop „Grundlagen vernetzter Systeme“)
- [KEFH18] KÖPF, F.; EBBERT, S.; FLAD, M.; HOHMANN, S.: Adaptive Dynamic Programming for Cooperative Control with Incomplete Information. In: *IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, 2018, S. 2632–2638
- [KFH19] KÖPF, F.; FLAD, M.; HOHMANN, S.: *Kooperative Optimalregelung unter eingeschränkter Kommunikation und Information*. Günzburg, 2019 (GMA-Fachausschuss 1.50 Workshop „Grundlagen vernetzter Systeme“)
- [KIB⁺19] KASTNER, A.; INGA, J.; BLAUTH, T.; KÖPF, F.; FLAD, M.; HOHMANN, S.: Model-Based Control of a Large-Scale Ball-on-Plate System With Experimental Validation. In: *IEEE International Conference on Mechatronics (ICM)*, 2019, S. 257–262
- [KIR⁺17] KÖPF, F.; INGA, J.; ROTHFUSS, S.; FLAD, M.; HOHMANN, S.: Inverse Reinforcement Learning for Identification in Linear-Quadratic Dynamic Games. In: *IFAC-PapersOnLine* 50 (2017), Nr. 1, S. 14902–14908
- [KKBH23] KARG, P.; KÖPF, F.; BRAUN, C. A.; HOHMANN, S.: Excitation for Adaptive Optimal Control of Nonlinear Systems in Differential Games. In: *IEEE Transactions on Automatic Control* (2023)
- [KKF⁺20] KUHN, E.; KUSMARTSEV, F.; FLAD, M.; HOHMANN, S.; KÖPF, F.: *Bewegungsabhängiges Stabilisierungsunterstützungssystem: Patent DE 10 2019 210 232*. 2020
- [KKIH21] KÖPF, F.; KILLE, S.; INGA, J.; HOHMANN, S.: Adaptive Optimal Trajectory Tracking Control Applied to a Large-Scale Ball-on-Plate System. In: *IEEE American Control Conference (ACC)*, 2021, S. 1352–1357
- [KNFH20] KÖPF, F.; NITSCH, A.; FLAD, M.; HOHMANN, S.: Partner Approximating Learners (PAL): Simulation-Accelerated Learning with Explicit Partner Modeling in Multi-Agent Domains. In: *IEEE International Conference on Control, Automation and Robotics (ICCAR)*, 2020, S. 746–752

- [KPRH20] KÖPF, F.; PUC CETTI, L.; RATHGEBER, C.; HOHMANN, S.: Reinforcement Learning for Speed Control with Feedforward to Track Velocity Profiles in a Real Vehicle. In: *IEEE International Conference on Intelligent Transportation Systems (ITSC)*, 2020, S. 1–8
- [KRP⁺20] KÖPF, F.; RAMSTEINER, S.; PUC CETTI, L.; FLAD, M.; HOHMANN, S.: Adaptive Dynamic Programming for Model-Free Tracking of Trajectories with Time-Varying Parameters. In: *International Journal of Adaptive Control and Signal Processing* 34 (2020), Nr. 7, S. 839–856
- [KTFH20] KÖPF, F.; TESFAZGI, S.; FLAD, M.; HOHMANN, S.: Deep Decentralized Reinforcement Learning for Cooperative Control. In: *IFAC-PapersOnLine* 53 (2020), Nr. 2, S. 1555–1562
- [KWFH20] KÖPF, F.; WESTERMANN, J.; FLAD, M.; HOHMANN, S.: Adaptive Optimal Control for Reference Tracking Independent of Exo-System Dynamics. In: *Neurocomputing* 405 (2020), S. 173–185
- [LKS⁺18] LEMMER, M.; KÖPF, F.; SCHWAB, S.; FLAD, M.; HOHMANN, S.: Modeling of Human-Centered Cooperative Control by Means of Tracking in Discrete-Time Linear-Quadratic Differential Games. In: *IEEE International Conference on Artificial Intelligence and Knowledge Engineering (AIKE)*, 2018, S. 156–161
- [PKRH20] PUC CETTI, L.; KÖPF, F.; RATHGEBER, C.; HOHMANN, S.: Speed Tracking Control Using Online Reinforcement Learning in a Real Car. In: *IEEE International Conference on Control, Automation and Robotics (ICCAR)*, 2020, S. 392–399
- [RIK⁺17] ROTHFUSS, S.; INGA, J.; KÖPF, F.; FLAD, M.; HOHMANN, S.: Inverse Optimal Control for Identification in Non-Cooperative Differential Games. In: *IFAC-PapersOnLine* 50 (2017), Nr. 1, S. 14909–14915

Betreute studentische Arbeiten

- [Bla18] BLAUTH, T.: *Kamerabasierte Messung der Ballposition eines Ball-auf-Platte-Systems*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Bachelorarbeit, 2018
- [Bra18] BRAUN, C. A.: *Entwicklung und Analyse Reinforcement-Learning-basierter Regler für Differentialspiele mit unbekanntem Kooperationspartner*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2018
- [Büh20] BÜHRLE, E.: *Reinforcement Learning zur adaptiven Zustandsfolgeregelung zeitkontinuierlicher Systeme*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2020
- [Den17] DENGLER, K.: *Implementierung und Untersuchung eines Expectation-Maximization-Algorithmus zur Trajektorienplanung in endlichen Zustandsräumen*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Bachelorarbeit, 2017
- [Ebb17] EBBERT, S.: *Bewertung von Methoden zur Regelung kooperativer Systeme mit unbekanntem Handlungspartner*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2017
- [Gaj18] GAJEK, S.: *Development of a Solution Concept for Nonlinear Cooperative Control Problems by Means of Dynamic Game Theory and Methods of Reinforcement Learning*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2018
- [Kar19] KARG, P.: *Strukturierte Anregung Reinforcement-Learning-basierter kooperativer Regler*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2019
- [Kil20] KILLE, S.: *Anwendung Reinforcement-Learning-basierter Trajektorienfolgeregler auf ein Ball-auf-Platte-System*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2020
- [Lem18] LEMMER, M.: *Entwurf und Analyse zeitdiskreter Linear-Quadratischer Zwei-Spieler-Folgeregelungen*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2018
- [Nit18] NITSCH, A.: *Entwicklung kooperativer Regler basierend auf Identifikationsmethoden und Deep Reinforcement Learning*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2018
- [Ram19] RAMSTEINER, S.: *Parameterbasierte Folgeregelung unter Verwendung von Reinforcement-Learning-Methoden*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2019

-
- [Tes19] TESFAZGI, S.: *Deep Decentralised Reinforcement Learning for Cooperative Control*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2019
- [Wes18] WESTERMANN, J.: *Optimale Folgeregelung für unbekannte lineare Systeme basierend auf Q-Learning mit gleitendem Horizont*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2018
- [Wie19] WIECZOREK, D. H.: *Untersuchung und Modellierung menschlicher Lernprozesse in haptisch gekoppelten kooperativen Szenarien*, Fakultät für Elektrotechnik und Informationstechnik, Karlsruher Institut für Technologie, Masterarbeit, 2019