

Universal Distributional Decision-based Black-box Adversarial Attack with Reinforcement Learning

Yiran Huang¹, Yexu Zhou¹, Michael Hefenbrock¹, Till Riedel¹, Likun Fang¹, and Michael Beigl¹

Karlsruhe Institute of Technology, Karlsruhe, Germany
{yhuang, zhou, hefenbrock, riedel, fang, beigl}@teco.edu

Abstract. The vulnerability of the high-performance machine learning models implies a security risk in applications with real-world consequences. Research on adversarial attacks is beneficial in guiding the development of machine learning models on the one hand and finding targeted defenses on the other. However, most of the adversarial attacks today leverage the gradient or logit information from the models to generate adversarial perturbation. Works in the more realistic domain: decision-based attacks, which generate adversarial perturbation solely based on observing the output label of the targeted model, are still relatively rare and mostly use gradient-estimation strategies. In this work, we propose a pixel-wise decision-based attack algorithm that finds a distribution of adversarial perturbation through a reinforcement learning algorithm. We call this method Decision-based Black-box Attack with Reinforcement learning (DBAR). Experiments show that the proposed approach outperforms state-of-the-art decision-based attacks with a higher attack success rate and greater transferability.

Keywords: Adversarial attack · Decision attack · Reinforcement Learning.

1 Introduction

Many high-performing machine learning algorithms used in computer vision, speech recognition and other areas are susceptible to minimal changes of their inputs [1]. Despite their good performance, the vulnerability of machine learning models has raised widespread concerns. Small perturbation on road signs can have a serious impact on automated driving. These actions, which modify the benign input by imperceptible perturbations and thus manipulate the machine learning model to suit the attacker's interests, are called adversarial attacks.

Most adversarial attacks used to construct adversarial perturbation rely either on gradient information (white-box attack) or logit output (score-based attack) of the model. While these approaches help to study the limitations of current machine learning algorithms [15], they do not reflect the level of information a real attacker would have access to in most scenarios. In contrast, decision-based attacks consider limited access to the targeted model, allowing only the label information output by the model to be used. Such limited access is far more common in the real-world scenarios making its study more practical.

Most of decision-based attacks start with an adversarial example with large perturbation. Then, adversarial examples with smaller perturbations are gradually found by sample-based gradient estimation. Different attacks exploit the samples in different ways, and therefore the efficiency of the algorithms varies. Such gradient-estimation-based approaches, however, require a large number of queries to the targeted model, which affects the efficiency of the algorithm and makes it impossible to perform real-time attack. In addition, the perturbations generated by gradient-based approach are too specific to the particular targeted model and benign example, and therefore lack transferability. To address these shortcomings, we propose a novel pixel-wise decision-based attack approach, called Decision-based Black-box Attack with Reinforcement learning (DBAR), which is guided by rewards instead of gradients. We therefore phrase the search for adversarial perturbations as an reinforcement learning task. Depending on whether the learned agent is targeting a single or multiple benign examples, two different attacks are designed, i.e., context-free attack and context-aware attack.

Our contributions can be summarized as: (1) Context-free DBAR achieves state-of-the-art performance, and perturbations sampled from the discovered distribution are more transferable than those generated by the other decision based attacks. In addition, the context-free DBAR is an universal attack which can also attack time-series data and super pixel of image data. (2) The context-aware attack achieves an effective attack without any queries on the targeted model after training, which is not possible for most existing decision-based attacks. (3) The algorithm generates a distribution which can be used to sample multiple different attacks.

2 Related work

The definition of decision-based attack was first proposed in [1]. It starts with an example in the target category and optimizes the attack with random selection and validation. This method is simple and effective; however, it is inefficient because the information from the sampled examples is not fully utilized e.g., information from the 'worse' samples. Several methods attempt to bridge this gap. For example, [2] biases the sampling process by combining low-frequency noise with gradients from surrogate models. However, its performance depends on the transferability between the surrogate model and the target model. Similarly, transfer-based attacks [15] also rely on carefully chosen surrogate models. However they obtain an attack on the original model by attacking the surrogate model. Opt attack [5] transforms the adversarial attack problem into a continuous real-valued optimization problem, i.e., the direction and distance to the decision boundary. This optimization problem can be solved by any zeroth-order optimization algorithm. However, distance calculation and gradient estimation in large dimensions will consume a large number of queries, which reduces the efficiency of the algorithm. Evo attack [8] applies evolutionary algorithms to generate adversarial perturbations and employs some techniques to reduce the dimensionality of the search space. It uses a custom variant in normal distribution and update the variant with (1+1)-CMA-ES. However, the variance is sign-independent and the sampling is therefore unstable. Rays [4] uses the dichotomous method to find perturbations. Although it has achieved good results on many datasets, the effectiveness of the algorithm is difficult to prove as results depend strongly on the

test set. HSJ [3] estimate gradient in different way and achieve a decision-based attack. However, gradient estimation is time-consuming and, at the same time, reduces the transferability of the generated adversarial perturbations. In this paper, we try to solve the problem without estimating the gradient.

3 Methodology

We model the decision-based black-box adversarial attack problem as finding the adversarial distribution p_Θ with parameters Θ of an m -class deep classification model $\mathcal{M} : \mathbb{R}^d \rightarrow [m]$ that accepts an input $x \in [0, 1]^d$ and outputs $y \in [m] = \{1, \dots, m\}$. The objective function can be described as

$$\min_{\Theta} \left(\lambda \mathbb{E}_{\eta \sim p_\Theta} \|\eta\|_\infty - \mathcal{P}_{\eta \sim p_\Theta} (\mathcal{M}(x + \eta) \neq \mathcal{M}(x)) \right), \quad (1)$$

where η is the adversarial perturbation sampled from the distribution p_Θ , $\|\cdot\|_\infty$ denotes the l_∞ norm and \mathcal{P} evaluates a probability. The objective function consists of two components, the expected l_∞ norm of the perturbations sampled from the distribution and the attack success rate of the perturbations sampled from the distribution. λ is a parameter that trade-off between the expectation and the success rate. The goal of the problem definition is to find a distribution such that the center of the distribution found is as close as possible to the benign example x while the adversarial perturbations sampled from the distribution maintain a high attack success rate.

To solve (1) through a reinforcement learning algorithm, depending on whether single or multiple benign examples are considered, we design two different environments: a context-free environment and a context-aware environment, which correspond to context-free attack and context-aware attack. Both environments share the same setting except the transition model. The state and action space are both set to \mathbb{R}^d . The perturbation distribution p_Θ to optimized is regarded as the agent. Each adversarial perturbation η sampled from the distribution is an action and each benign example x is a state. Since the action is continuous, we model the agent with normal distribution $p_\Theta(\eta | x) = \mathcal{N}(\eta | \mu_\Theta(x), \text{diag}(\Sigma_\Theta(x)))$. A trajectory τ consists of fixed number of decision step. In each decision step, a perturbation (action) is sample from the distribution (agent) and send to the environment to get the reward and next state. To achieve the optimization goal, we define the reward function as

$$r(x, \eta) = \frac{2 \cdot \mathbb{1}_{\{\mathcal{M}(x+\eta) \neq \mathcal{M}(x)\}} - 1}{\|\eta\|_\infty}, \quad (2)$$

where $\mathbb{1}$ is the indicator function to identify whether adversarial perturbations mislead the classifier \mathcal{M} . When the attack is successful, $2 \cdot (\mathbb{1}_{\{\mathcal{M}(x+\eta) \neq \mathcal{M}(x)\}} - 1) = 1$, the algorithm can try to increases the reward by shrinking the perturbation. On the other hand, if the attack fails, the algorithm may try to find a successful attack by increasing the perturbation¹. The expectation and probability in (1) are approximated by Monte Carlo estimation using sampled trajectories.

¹ Eq. 2 can be modified to a target attack by setting the condition of indicator function to $\mathcal{M}(x + \eta) = \text{target}$.

In both environment settings, the next state is selected independent of the current state and action. In the context-free environment, only one state, the benign example, exists, while in the context-aware environment, the next state is selected by random sampling.

The objective function can be written as:

$$J(\Theta) = \int p_{\Theta}(\tau) \left(\sum_{t=0}^T r(x_t, \eta_t) \right) d\tau, \quad (3)$$

$$p_{\Theta}(\tau) \approx \prod_{t=0}^{T-1} p_{\Theta}(\eta_t | x_t).$$

To learn Θ , we need to calculate the gradient of the objective function (3). In the black-box adversarial attack, each reward in one trajectory is treated equally and does not depend on the actions in the other time step of the same trajectory. Therefore, when calculating the gradient, at time step t , terms that do not depend on the action η_t can be omitted. Using the log derivation trick and Monte Carlo sampling [19], the gradient of the objective function can be expressed as

$$\begin{aligned} \nabla J(\Theta) &\approx \sum_{i=0}^I \left[\left(\nabla_{\Theta} \log \left(\prod_{t=0}^{T-1} p_{\Theta}(\eta_{t,i} | x_{t,i}) \right) \right) \left(\sum_{t=0}^T r(x_{t,i}, \eta_{t,i}) \right) \right] \\ &\approx \sum_{i=0}^I \left[\sum_{t=0}^{T-1} r(x_{t,i}, \eta_{t,i}) \nabla_{\Theta} \log p_{\Theta}(\eta_{t,i} | x_{t,i}) \right] \\ &= \sum_{i=0}^{I \cdot T} r(x_i, \eta_i) \nabla_{\Theta} \log p_{\Theta}(\eta_i | x_i) \end{aligned}$$

So far, this gradient is valid only for the samples generated by p_{Θ} . We apply the importance sampling technique so that old trajectories can be reused. In addition, we limit the update step size as suggested in [18], since importance sampling only works when the update size is small. Together with the stable training trick mentioned in [14], the gradient can be expressed as

$$\nabla J(\Theta) = \sum_{i=0}^M (\nabla_{\Theta} \min(w_i(\Theta), \text{clip}(w_i(\Theta), 1 - \epsilon, 1 + \epsilon))) \cdot (r(x_i, \eta_i) - V(x_i))$$

$$\text{with } w_i(\Theta) = \frac{p_{\Theta}(\eta_i | x_i)}{p_{\Theta_{\text{old}}}(\eta_i | x_i)},$$

where $p_{\Theta_{\text{old}}}$ is the distribution that generates the training samples and p_{Θ} is the distribution, that is frequently updated. The parameter ϵ limits the update size and $V(x_i)$ is the expected reward given to a benign example x_i .

Algorithm 1 summarizes the process of generating adversarial distribution, where Actor: $\mathbb{R}^d \rightarrow \mathbb{R}^d \times \mathbb{R}^d$ and Critic: $\mathbb{R}^d \rightarrow \mathbb{R}$ are two neural networks with the same ResNet architecture except for the output layer.

Algorithm 1 Generating adversarial distribution through DBAR

Input: x_0 (benign example), N (number of iterations), M (number of samples in one iteration), K (number of training), L (size of minibatch), $init_mean$, $init_std$, ϵ

Output: Actor

```

1: Initialization: Initialize Actor( $\cdot$ ) with  $init\_mean$  and  $init\_std$ , Critic( $\cdot$ ),  $x \leftarrow x_0$ 
2: for  $i \leftarrow 1$  to  $N$  do
3:    $B \leftarrow []$ 
4:   for  $j \leftarrow 1$  to  $M$  do
5:      $\mu, \mathbf{I} \cdot \sigma^2 \leftarrow \text{Actor}(x)$ 
6:      $r, p, x' \leftarrow$  sample action from  $\mathcal{N}(\mu, \mathbf{I} \cdot \sigma^2)$ , calculate its log-probability and applied it
       to the environment to get reward and the next benign example
7:      $B \leftarrow B \cup \{r : r, lp : lp, x : x\}, x \leftarrow x'$ 
8:   end for
9:   for  $k \leftarrow 1$  to  $K$  do
10:     $\{B_1, \dots, B_{\lfloor M/L \rfloor}\} \leftarrow$  generate mini-batch from  $B$ 
11:    for  $b \leftarrow \{B_1, \dots, B_{\lfloor M/L \rfloor}\}$  do
12:       $\mu', \mathbf{I} \cdot (\sigma')^2 \leftarrow \text{Actor}(b[x])$ 
13:       $lp' \leftarrow$  compute log-probability of  $b[x]$  in  $\mathcal{N}(\mu', \mathbf{I} \cdot (\sigma')^2)$ 
14:       $v \leftarrow \text{Critic}(b[x])$ 
15:       $a \leftarrow b[r] - v$ 
16:       $w \leftarrow \exp(lp' - b[lp])$ 
17:       $loss\_actor \leftarrow \min(w, \text{clip}(w, 1 - \epsilon, 1 + \epsilon)) \cdot a$ 
18:       $loss\_critic \leftarrow \text{MSE}(b[r], v)$ 
19:      update Actor with the gradient of  $loss\_actor$ 
20:      update Critic with the gradient of  $loss\_critic$ 
21:    end for
22:  end for
23: end for
24: return Actor

```

4 Experiments

In this section, we perform experiments to investigate the following questions: (i) How does the context-free DBAR algorithm perform on image datasets compared to state-of-the-art decision-based attack methods? (ii) Can context-free DBAR be applied to time-series datasets? (iii) Are the perturbations discovered by the context-free DBAR algorithm more transferable than those discovered by state-of-the-art decision-based attack methods? (iv) Can context-aware DBAR perform real-time attacks after training? (v) How do the hyper-parameters affect the performance of context-free DBAR?

4.1 Experiment Setting

Baselines and hyper-parameters: To evaluate DBAR, we compare it with the following decision-based attacks: (i) the state-of-the-art Decision-based black-box Boundary attack [1] and the HopSkipJump Attack [3] for image datasets. (ii) the Universal White-box attack FGSM [21] and BIM [13] for time series datasets.

All attacks are implemented by the python package Foolbox [17]. We use the default hyper-parameter settings for all attacks with a fixed random seed. We limit the maximum number of queries for all the attacks to 20000. For the DBAR algorithm, the number of training epoch K is set to $K = 10$, with a mini-batch size of $L = 10$ and $\epsilon = 0.02$. Additionally $init_mean = 0$ and $init_std = 0.5$.

Datasets and models: We carried out attacks over the following image datasets with varied dimensions and dataset sizes: CIFAR10 [12], CIFAR100 [12], STL10 [6], Caltech101 [10]. The pixel values of all images are normalized to $[0, 1]$. We also attack models with different structures such as ResNet20 [11] with 272474 parameters and VGG11 [20] with 9756426 parameters. Both models were obtained from Pytorch [16]. In addition, we carried out attacks over the publicly available time series UCR archiv dataset [7] and attack the time series ResNet-ts model as defined by [9].

4.2 Adversarial examples for image and time series data

We perform non-target attacks on all the image datasets mentioned above and summarise the attack success rate (ASR) of each methods in Table 1. Concretely, all attacks are applied to 1000 correctly classified test examples from each dataset. If an adversarial example can mislead the classification model and is in a 0.04 ($10/255$) l_∞ neighborhood of the benign example, we denote this attack as a success. From the result, we see that, context-free DBAR achieves better results on most of the datasets, except for the STL10 dataset on the VGG11 model. Struggling with finding the first success attack is probably the main reason for the failure of the attack. This happens when the $init_std$ is set too small for the given data set. The influence of the hyper-parameters on the attack performance an run-time is analyzed in experiment 4.4.

		Model Accuracy	BA	HSJ	DBAR			Model Accuracy	BA	HSJ	DBAR
Cifar10	ResNet20	0.91	1.00	1.00	1.00	STL10	ResNet20	0.83	1.00	1.00	1.00
	VGG11	0.90	0.71	0.78	0.84		VGG11	0.84	0.59	0.90	0.78
Cifar100	ResNet20	0.66	0.95	0.95	0.96	CalTech101	ResNet20	0.79	0.93	1.00	1.00
	VGG11	0.61	0.63	0.80	0.84		VGG11	0.68	0.51	0.80	0.82

Table 1. Attack success rate (ASR) of three different decision-based attack methods: Boudary Attack (BA), HopSkipJump Attack (HSJ) and the proposed context-free DBAR, against model ResNet20 and VGG11 on four different image datasets with varies sizes.

DBAR is universal in the sense that it is able to attack any example in form of \mathcal{R}^d . To prove this, we compare the performance of the proposed context-free DBAR against ResNet-ts model on the UCR open source time series datasets with the two popular white-box attacks FGSM and BIM. When the size of the perturbation found by an attack is smaller than 0.1, as proposed in [9], and the adversarial example can mislead the classification, we regard the attack as success. The parameter settings of this experiment are different from other experiments because the time series data are not normalized to $[0, 1]$. Furthermore, note that DBAR, as opposed to FGSM and BIM is a black-box attack. We remove the step limit and set the initial standard deviation to 80 to avoid struggling with finding the first successful attack.

	Model Accuracy	FGSM	BIM	DBAR		Model Accuracy	FGSM	BIM	DBAR
50words	0.73	0.77	0.88	0.91	DistalAge	0.80	0.78	0.79	0.64
Adiac	0.83	0.96	0.98	0.99	FaceAll	0.86	0.10	0.15	0.20
Beef	0.77	0.74	0.87	0.87	FaceUCR	0.95	0.17	0.20	0.19
Car	0.93	0.76	0.92	0.80	ElectricDevices	0.74	0.34	0.58	1.00
Diatom	0.30	0.00	0.00	0.41	ItalyPowerDemand	0.96	0.04	0.04	0.20

Table 2. Attack success rate (ASR) of three different adversarial attack methods: FGSM [21], BIM [13] and the proposed context-free DBAR, against ResNet-ts on ten different time series datasets with varies size and dimensions.

The results can be seen in Table 2. Although compared to the white-box algorithms, the context-free DBAR achieves better results on six datasets, ties on one dataset, and worse results on three datasets. As for the image datasets, the ASR score on the time series data is independent of the accuracy of the target model. And probably because of the different principles of generating attacks, DBAR performs well on some datasets where BIM performs very poorly, e.g., Diatom, ItalyPowerDemand. At the same time, the opposite situation also exists, see DistalAge. It is important to note that although FGSM has the worst results, it is attacking in real time, while both BIM and the proposed context-free DBAR require multiple iterations. However, a similar real-time attack can be achieved by context-aware DBAR. This is demonstrated in Section 4.3.

4.3 Transferability of the perturbation distribution and real-time attack

The high transferability of the perturbation allows attacks generated against one platform to be applied to the other. To demonstrate the transferability of the perturbation found by the proposed method, we run the experiment in Section 4.2 on the Cifar datasets again and apply the attacks generated against the ResNet20 model to the VGG11 model and vice versa. The ASR score is given in Table 3. The performance of the proposed method is significantly better than that of the other methods. In particular, the perturbations generated against the VGG11 model on Cifar100 dataset have a near-average ASR score against the ResNet20 model. Besides, we can see that perturbations generated against simple model (with fewer parameters, ResNet20) are difficult to perform success attack against the more complex VGG11 model.

	Targeted model	Boundary Attack	HSJ	DBAR		Targeted model	Boundary Attack	HSJ	DBAR
Cifar10	ResNet20	0.04	0.04	0.13	Cifar100	ResNet20	0.14	0.14	0.24
	VGG11	0.10	0.10	0.36		VGG11	0.14	0.09	0.45

Table 3. Attack success rate (ASR) of attacks generated for VGG11 and applied to ResNet20 and vice versa. Targeted model denotes the model used to generate the perturbation (attack).

To evaluate the capabilities for real time attacks, we apply context-aware DBAR against ResNet20 on the Cifar10 and Cifar100. The results can be seen in Fig. 1. We can find that after 60 iterations, the adversarial perturbation obtained by sampling are

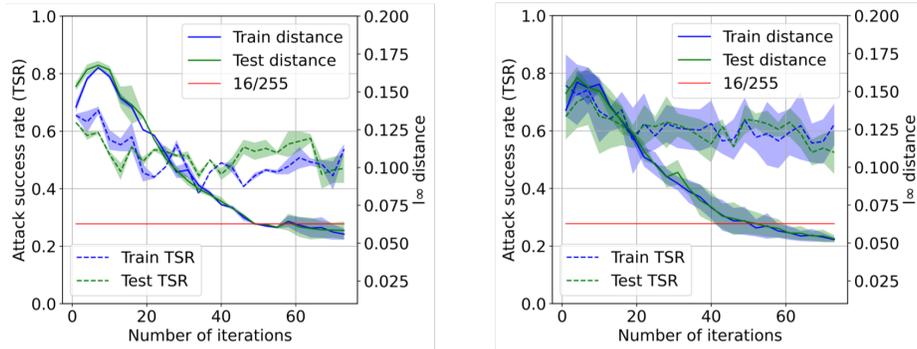


Fig. 1. Performance of context-aware DBAR against ResNet20 on Cifar10 (left) and Cifar100 (right) datasets. The picture on the left shows the performance of Cifar10 and the one on the right shows the performance of Cifar100.

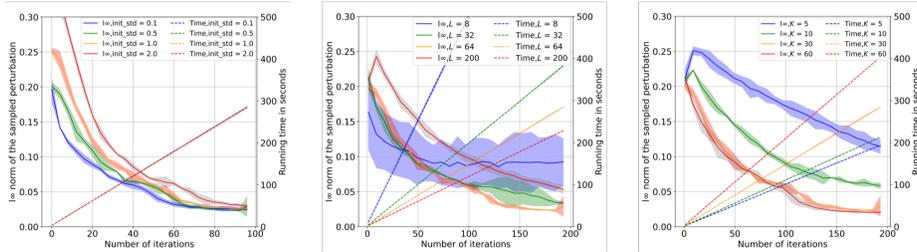


Fig. 2. Performance comparison of context-free DBAR with different hyper-parameter settings. The l_∞ norm is shown as a solid line in the plot, while the run time is shown as a dashed line.

able to implement effective attacks that have a success rate of about 50% on Cifar10 and 60% on Cifar100. We observe higher l_∞ norms for the perturbations at the beginning, which is most likely caused by the different update directions between the different benign examples. Ideally, the algorithm should reduce the size of the perturbation while increasing the success rate of the training. However, the attack success rate decreases as the training progresses, although the decrease is not significant. The algorithm improves the rewards obtained by reducing the perturbation size. This problem can be mitigated by increasing the contribution of the success attack in the reward function, see (2).

4.4 Impact of the parameter selection

In this experiment, we analyse the effect of following hyper-parameters on the performance of context-free DBAR regarding *init_std*, number of training epoch in each iteration K and the size of mini batch L . The baseline hyper-parameters used in the experiments are set as follows as $init_std = 0.5$, $L = 64$, $K = 30$. In each trial we modify one of the above parameters and summarise the result in Fig. 2. Evidently, the parameters

have a great impact on the convergence speed, stability and runtime of DBAR. The leftmost plot in Fig. 2 shows the effect of the standard deviation of the initial distribution $init_std$. The larger the parameter, the larger the perturbation sampled from the initial distribution. In general, the larger the initial perturbation, the higher the probability that it will successfully attack a benign example. When the $init_std$ is too small and the initial perturbation fails to attack the benign example, the algorithm will struggle in looking for the first successful attack. The size of the training batches L affects the stability, where lower L results in less stability.

Since the number of samples per iteration is constant, the smaller the batch size, the more times the agent is updated and the longer each iteration takes.

In addition, due to the importance sampling, the log-probability of actions needs to be recalculated before updating agent, which further increases the runtime. These are reflected in the plot in the center of Fig. 2. The rightmost plot shows that larger number of training K in one iteration converges faster. However, since there are more updates per iteration, it also takes longer to run.

5 Conclusion

In this paper, we formulate the decision-based black-box adversarial attack as a reinforcement learning task and search for the adversarial attack based on reward criterion. We have experimentally demonstrated the feasibility of the proposed algorithm. In addition, we have shown some advantages of the algorithm such as the ability to generate attacks for different kinds of data such as images and time series, the transferability of attacks between different models and the ability for real-time attack. Besides, there is still room for further exploration in the proposed approach through using different reward functions or investigating the usage bounds by attacking very small perturbation or high resolution images.

6 Acknowledgements

This work was partially funded by the Ministry of The Ministry of Science, Research and the Arts Baden-Wuerttemberg as part of the SDSC-BW and by the German Ministry for Research as well as by Education as part of SDI-C (Grant 01IS19030A)

References

1. Brendel, W., Rauber, J., Bethge, M.: Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. arXiv preprint arXiv:1712.04248 (2017)
2. Brunner, T., Diehl, F., Le, M.T., Knoll, A.: Guessing smart: Biased sampling for efficient black-box adversarial attacks. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4958–4966 (2019)
3. Chen, J., Jordan, M.I., Wainwright, M.J.: Hopskipjumpattack: A query-efficient decision-based attack. In: 2020 IEEE Symposium on Security and Privacy (SP). pp. 1277–1294. IEEE (2020)

4. Chen, J., Gu, Q.: Rays: A ray searching method for hard-label adversarial attack. In: Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. pp. 1739–1747 (2020)
5. Cheng, M., Le, T., Chen, P.Y., Yi, J., Zhang, H., Hsieh, C.J.: Query-efficient hard-label black-box attack: An optimization-based approach. arXiv preprint arXiv:1807.04457 (2018)
6. Coates, A., Ng, A., Lee, H.: An analysis of single-layer networks in unsupervised feature learning. In: Proceedings of the fourteenth international conference on artificial intelligence and statistics. pp. 215–223. JMLR Workshop and Conference Proceedings (2011)
7. Dau, H.A., Bagnall, A., Kamgar, K., Yeh, C.C.M., Zhu, Y., Gharghabi, S., Ratanamahatana, C.A., Keogh, E.: The ucr time series archive. *IEEE/CAA Journal of Automatica Sinica* **6**(6), 1293–1305 (2019)
8. Dong, Y., Su, H., Wu, B., Li, Z., Liu, W., Zhang, T., Zhu, J.: Efficient decision-based black-box adversarial attacks on face recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7714–7722 (2019)
9. Fawaz, H.I., Forestier, G., Weber, J., Idoumghar, L., Muller, P.A.: Adversarial attacks on deep neural networks for time series classification. In: 2019 International Joint Conference on Neural Networks (IJCNN). pp. 1–8. IEEE (2019)
10. Griffin, G., Holub, A., Perona, P.: Caltech-256 object category dataset (2007)
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
12. Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images (2009)
13. Kurakin, A., Goodfellow, I., Bengio, S.: Adversarial machine learning at scale. arXiv preprint arXiv:1611.01236 (2016)
14. Mnih, V., Badia, A.P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., Kavukcuoglu, K.: Asynchronous methods for deep reinforcement learning. In: International conference on machine learning. pp. 1928–1937. PMLR (2016)
15. Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z.B., Swami, A.: Practical black-box attacks against machine learning. In: Proceedings of the 2017 ACM on Asia conference on computer and communications security. pp. 506–519 (2017)
16. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* **32**, 8026–8037 (2019)
17. Rauber, J., Brendel, W., Bethge, M.: Foolbox: A python toolbox to benchmark the robustness of machine learning models. arXiv preprint arXiv:1707.04131 (2017)
18. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347 (2017)
19. Shapiro, A.: Monte carlo sampling methods. *Handbooks in operations research and management science* **10**, 353–425 (2003)
20. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
21. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199 (2013)