

**PREDICTION OF TELECOMMUTING ENGAGEMENT THROUGH MACHINE  
LEARNING TO ENHANCE TRAVEL SURVEY DATA**

**Anna Sophie Reiffer, Corresponding Author**

Email: [anna.reiffer@kit.edu](mailto:anna.reiffer@kit.edu)

ORCID: 0000-0003-1764-0154

**Martin Kagerbauer**

[martin.kagerbauer@kit.edu](mailto:martin.kagerbauer@kit.edu)

ORCID: 0000-0002-4252-7874

**Peter Vortisch**

[peter.vortisch@kit.edu](mailto:peter.vortisch@kit.edu)

ORCID: 0000-0003-1647-2435

Word Count: 6873 words + 2 table(s)  $\times$  250 = 7373 words

Submission Date: August 2, 2023

## 1 ABSTRACT

2 This paper presents a novel approach to enhance household travel survey (HTS) data by predict-  
3 ing telecommuting engagement using machine learning (ML) classification techniques. The study  
4 aims to address the debate surrounding the impact of telecommuting on overall travel behavior,  
5 considering rebound effects and latent demand. While previous research has primarily relied on  
6 questionnaires or HTS data for analysis, few studies have successfully integrated telecommut-  
7 ing behavior into travel demand models. The intricate relationship between telecommuting and  
8 travel behavior has been a challenge, limiting the incorporation of telecommuting data into exist-  
9 ing models. This study fills this research gap by utilizing ML algorithms to predict telecommuting  
10 engagement based on one-day HTS data, employing features such as daily distances traveled and  
11 time spent at home.

12 Three feature selection algorithms, Boruta, VSURF, and Recursive Feature Elimination  
13 (RFE), were applied to identify the most relevant features for the ML models. Among the five clas-  
14 sification methods tested, the Random Forest (RF) model utilizing features selected by the Boruta  
15 algorithm demonstrated superior performance, achieving high accuracy, specificity, F1-Score, and  
16 Matthew's Correlation Coefficient (MCC). The Bayesian Network (BN) model, while performing  
17 well in sensitivity, underperformed in other metrics due to the unsuitability of continuous data.

18 To evaluate the proposed approach's generalization, the RF model was applied to a sepa-  
19 rate HTS dataset from the German Mobility Panel (MOP). The out-of-sample prediction achieved  
20 promising results, with a 76% accuracy in predicting telecommuting days. The approach presented  
21 in this study has potential applications in enhancing HTS data and can be extended to other data  
22 sources to improve activity-based models.

23  
24 *Keywords:* household travel survey data, machine learning, data fusion, data enhancement

## 1 INTRODUCTION

2 Telecommuting has been analyzed as a measure to reduce travel and subsequent emissions for al-  
3 most five decades (1) (1). Although studies consistently show that working from home decreases  
4 the number of commuting trips, it is still debated whether telecommuting decreases overall travel  
5 or if rebound effects and latent demand are high enough to offset any benefits. To this day, the  
6 majority of studies analyze travel behavior impacts of telecommuting based on questionnaires or  
7 household travel surveys (HTS) through descriptive or statistical model analysis (see, e.g., (2–6)).  
8 Despite the growing interest in incorporating telecommuting behavior into travel demand mod-  
9 els, few studies have successfully done so. This is mainly due to the intricate relationship between  
10 telecommuting and travel behavior, which even most activity-based models fail to account for. Fur-  
11 thermore, data on telecommuting engagement is limited, particularly when using HTS. Although  
12 some surveys ask respondents about the frequency at which they work from home, information on  
13 whether they worked from home on the day of the survey is only available in a few surveys, usually  
14 when the respective study is focused on in-home activities in addition to travel behavior (7). As  
15 emphasized by Asgari et al. (5), this information is crucial for accurately modeling telecommut-  
16 ing in travel demand models, and in light of the increase in work-from-home activities since the  
17 pandemic, it is essential that we examine this issue more closely. However, as appropriate data  
18 sources are scarce, new methods have to be explored on how existing data can be leveraged to  
19 inform timely research.

20 This paper aims to predict telecommuting engagement among respondents of a Household  
21 Travel Survey (HTS) by employing machine learning (ML) classification techniques. The study  
22 will utilize activity-travel behavior and sociodemographic data to train and assess various models  
23 using a one-day HTS. The efficacy of the trained models will be assessed by implementing them  
24 on a separate seven-day HTS dataset. This study's findings demonstrate that machine learning  
25 methods can be used to enhance HTS data using a secondary data source.

26 The need to account for non-travel activities has been addressed in previous studies on the  
27 design of advanced survey methods. Aschauer et al. (8) present a "Mobility-Activity-Expenditure-  
28 Diary" in which both travel and non-travel activities were reported. Schmid et al. (9) conducted a  
29 multi-stage survey to account for travel behavior and time allocation. A smartphone-based survey  
30 of travel, activities, and time use was presented by Alho et al. (10). Similarly, Winkler et al. (11)  
31 adopted a smartphone-based survey method in which travel is tracked passively, and respondents  
32 supplement the data by providing detailed information on activities and expenditures. Although  
33 the data collected in these surveys are undeniably valuable, all surveys offered monetary incentives  
34 to ensure high-quality responses. However, necessary funds are not always available. In an effort  
35 to reduce costs and make the most of existing data, researchers have explored the integration of  
36 data from various sources, including mobile phone surveys and web surveys (12). The utilization  
37 of smart card data can be augmented with HTS data (13), while the purpose of trips can be inferred  
38 through TNC data via HTS data (14). Moreover, HTS data can be supplemented with data from  
39 mobile phones (15) or social media (16). Studies have indicated the potential benefits of combining  
40 older, yet information-rich, HTS data with newer HTS data (17), as well as how the combination  
41 of data from two HTS can minimize biases and underrepresentation (18). All of these studies  
42 incorporate either weighting techniques (12, 13, 17) or econometric approaches (14–16, 18) to  
43 enhance the respective datasets. Machine learning techniques have primarily been used to improve  
44 data that has been collected passively, such as GPS-tracked trips. These techniques have been  
45 employed to identify the purpose of trips, the types of destinations, and the mode of transportation

1 used. (19).

2 This study contributes to the current body of research on leveraging existing data sources to  
3 improve the value of the data by inferring non-travel activities. By utilizing ML-based data fusion  
4 of multiple sources, we can determine if in-home activities are linked to telecommuting. This will  
5 improve the data required for integrating telecommuting into travel demand models.

6 The rest of this paper is structured as follows. We first present the two HTS used in this  
7 study. We further provide an overview of the machine learning classifiers analyzed in this research,  
8 including feature selection algorithms. Subsequently, we describe the performance metrics used to  
9 assess the suitability of each classifier for the proposed task. We go on to describe and discuss our  
10 findings. We conclude this paper by providing a summary and general implications of our study.

## 11 MATERIALS AND METHODS

12 In this section, we first provide an overview of the research design applied in this study. Further,  
13 we describe the travel diary data used in our models, including a descriptive analysis of the key  
14 variables. Subsequently, we elaborate on the classification methods we applied to predict telecom-  
15 muting participation on the survey day, including the utilized feature selection algorithms.

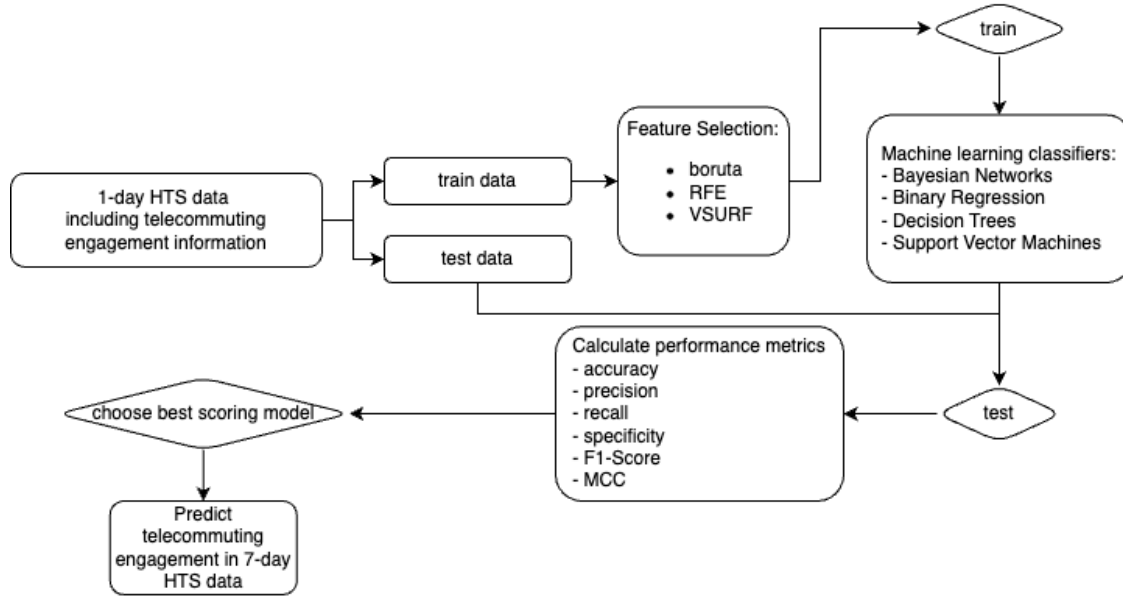
16 Considering that ML methods have not yet been utilized for data fusion of HTS, there is  
17 no apriori way of knowing which ML classifier will perform best. Similarly, it is unclear which  
18 features should be included or left out. We have, therefore, opted to train and test multiple machine-  
19 learning classifiers based on different feature configurations. The overall research framework em-  
20 ployed in this study is illustrated in figure 1. The data in this study stems from an HTS, which  
21 was split into a train and a test data set with a ratio of 70% (train) to 30% (test). Based on the test  
22 data, we applied three feature selection methods. For each of the three selected feature configura-  
23 tions, we trained five machine-learning classifiers. Subsequently, we predicted the telecommuting  
24 engagement in the test data. We determined the confusion matrix for each model and feature con-  
25 figuration pair and calculated six performance metrics. We chose the overall best scoring model,  
26 which was then used to predict telecommuting engagement in another HTS.

27 In the next subsections, we will provide more detail on the different parts of the research  
28 framework.

### 29 Household Travel Survey Data

30 This paper is based on data from a household travel survey conducted in the metropolitan area  
31 of Stuttgart, Germany. The survey was conducted in the fall of 2021 among 9,543 respondents  
32 in 4,567 (households). The scope of the survey was twofold: the data is used to update the travel  
33 demand model of the region, which is supplemented to account for telecommuting. Thus, the  
34 second scope of the survey was to gather detailed information on working from home behavior of  
35 respondents. The data used in this study stems from the travel diaries, which respondents kept for  
36 one day, and household- and person-level information. Based on the travel diaries, we determined  
37 variables of respondents' activity-travel patterns, namely duration of the activities *home*, *work*,  
38 *work-related*, *shopping/errands*, *education*, *escorting someone*, *roundtrips*, and *other*. The number  
39 of trips was determined by each category, as well as the overall travel time and distances traveled.  
40 Telecommuting engagement on the survey day was measured by one question, which asked where  
41 respondents worked on the given day. Additionally, we tested if the sociodemographic information  
42 of respondents influenced the models.

43 We further used data from the German Mobility Panel (MOP). The MOP is a longitudi-



**FIGURE 1:** Overview of the research framework employed in this study

1 nal household travel survey that has been conducted annually since 1994. Because of its panel  
 2 design, the questionnaires are relatively rigid to ensure the comparability and compatibility of sur-  
 3 vey waves. Furthermore, respondents are asked to keep a travel diary for seven consecutive days,  
 4 which puts a high response burden on them, and person- and household-level questionnaires have  
 5 to be kept as concise as possible. Thus, only a few questions concerning working from home are  
 6 included. In the 2022 wave, one additional question was included asking how many days respon-  
 7 dents worked from home in the survey week; however, not on which specific days they did. We use  
 8 this data to evaluate how the models perform on out-of-sample data. The data was prepared similar  
 9 to the 1-day HTS data. However, in this data, respondents were not asked where they worked on  
 10 a given survey day. Instead, they provided the number of days worked from home in the survey  
 11 week. Thus, the models' ability to predict telecommuting engagement can only be evaluated at an  
 12 aggregate level.

### 13 *Data preparation*

14 Because work from home is limited to - axiomatically - to employed respondents and furthermore,  
 15 by the jobs they hold or their employer. Thus the number of respondents dropped considerably  
 16 in both surveys after applying the respective filters. However, we were still able to leverage a  
 17 large enough sample for our study. The prevalence of working from home is relatively low in both  
 18 surveys. In the training data set, 473 out of the total used sample of 1380 attested they had worked  
 19 from home on the survey date. This amounts to about 37%. A similar proportion is provided in  
 20 the MOP data. Overall, 3896 telecommuting days were reported out of the 10515 reported survey  
 21 days, which equals a proportion of 34%.

22 In machine learning classification tasks, dealing with unbalanced classes can present sig-  
 23 nificant challenges. Unbalanced classes refer to situations where the distribution of class labels is  
 24 highly imbalanced, with one class having a much larger number of instances than the others. When  
 25 this occurs, standard classifiers tend to be biased toward the majority class, leading to poor per-

formance on the minority class and reduced overall predictive accuracy. Consequently, the model might struggle to generalize well to real-world scenarios where the minority class is of primary interest. One common issue is that the classifier might achieve high accuracy simply by always predicting the majority class while completely overlooking the minority class. One strategy to mitigate this issue is to apply the synthetic minority over-sampling technique (SMOTE) (20). SMOTE works by creating synthetic samples of the minority class by interpolating between existing instances. It selects a random instance from the minority class, identifies its  $k$ -nearest neighbors, and then generates new samples by combining the features of the selected instance with those of its neighbors.

## Classification Algorithms

In this section, we provide a concise overview of each classification algorithm investigated in this study.

### Bayesian Networks

Bayesian networks are probabilistic graphical models that represent the dependencies between random variables using a directed acyclic graph (DAG) (21, 22). Each node in the graph corresponds to a random variable, and the edges between nodes encode conditional dependencies. The conditional probability distribution for each variable given its parents is modeled using Bayes' rule. Let  $X_i$  denote the  $i$ -th random variable, and  $\text{Pa}(X_i)$  be the set of parent nodes of  $X_i$  in the graph. The joint probability distribution of all variables can be written as:

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{Pa}(X_i)) \quad (1)$$

Inference in Bayesian networks involves computing probabilities or making predictions based on observed evidence. Popular algorithms for inference include variable elimination and Markov chain Monte Carlo (MCMC) methods.

### Binary Regression

Binary regression, also known as logistic regression, is a popular supervised learning algorithm for binary classification tasks. Given a dataset with input-output pairs  $(\mathbf{x}_i, y_i)$ , where  $\mathbf{x}_i$  is a feature vector and  $y_i \in \{0, 1\}$  is the binary class label, the goal is to learn a model that estimates the probability of the positive class, i.e.,  $P(y_i = 1 | \mathbf{x}_i)$ . The logistic regression model assumes a linear relationship between the features and the log-odds of the positive class:

$$\log \left( \frac{P(y_i = 1 | \mathbf{x}_i)}{1 - P(y_i = 1 | \mathbf{x}_i)} \right) = \mathbf{w}^T \mathbf{x}_i + b \quad (2)$$

where  $\mathbf{w}$  is the weight vector and  $b$  is the bias term. To obtain probabilistic predictions, the logistic function is applied to the output of the linear model:

$$P(y_i = 1 | \mathbf{x}_i) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x}_i + b)}} \quad (3)$$

### Decision Trees

Decision trees are non-linear, hierarchical models used for both classification and regression tasks (23). They recursively split the data into subsets based on the values of individual features, aiming

1 to maximize the information gain or Gini impurity at each split. Each internal node in the tree  
 2 represents a decision based on a feature, and each leaf node corresponds to a predicted class or  
 3 regression value. The decision tree can be represented as a set of rules, and the final prediction for  
 4 a given input is determined by following the path from the root to the appropriate leaf node.

#### 5 *Random Forest*

6 Random forests are ensemble learning methods that combine multiple decision trees to improve  
 7 predictive performance and reduce overfitting (24, 25). The key idea is to build a collection of  
 8 decision trees by training each tree on a random subset of the training data (bootstrap sampling)  
 9 and selecting a random subset of features at each split. The final prediction is made by aggregating  
 10 the predictions of all individual trees, often using majority voting for classification problems or  
 11 averaging for regression problems. Random forests tend to be more robust and accurate than indi-  
 12 vidual decision trees, and they can handle high-dimensional data and capture complex relationships  
 13 between variables.

#### 14 *Support Vector Machines*

15 Support Vector Machines (SVMs) are powerful supervised learning algorithms used for both clas-  
 16 sification and regression tasks. SVMs aim to find the optimal hyperplane that best separates the  
 17 data points of different classes while maximizing the margin between the classes (26). In the  
 18 case of binary classification, given a training dataset  $(\mathbf{x}_i, y_i)$ , where  $\mathbf{x}_i$  is the feature vector and  
 19  $y_i \in \{-1, 1\}$  is the class label, SVMs find the weight vector  $\mathbf{w}$  and bias term  $b$  that define the  
 20 decision boundary:

$$\mathbf{w}^T \mathbf{x}_i + b = 0 \quad (4)$$

21 The margin is computed as the distance between the hyperplane and the closest data points  
 22 (support vectors) from each class. SVM aims to maximize this margin while penalizing misclas-  
 23 sifications. For non-linearly separable data, SVM can use kernel tricks to map the data into a  
 24 higher-dimensional space, where linear separation becomes possible. Common kernel functions  
 25 include polynomial, radial basis function (RBF), and sigmoid kernels.

#### 26 **Feature Selection**

27 We further present the feature selection algorithms applied. We have trained and tested all al-  
 28 gorithms on the dataset described above using the features determined by the respective feature  
 29 algorithm.

#### 30 *Boruta*

31 Boruta is a method for determining the importance of variables in a system using random forests.  
 32 The system involves replicating each descriptive variable and randomly permuting the values of  
 33 replicated variables across objects (27). The randomization is different for each run of the random  
 34 forest algorithm. The importance of each variable is computed for each run, and a statistical test  
 35 is performed to determine if the variable is significant or not. An attribute is considered important  
 36 for a single run if its level of importance is greater than the highest level of importance among all  
 37 randomized attributes. If a variable is deemed unimportant, it is removed from the system along  
 38 with its replicated mirror pair. The procedure is repeated for a predefined number of iterations or  
 39 until all attributes are either rejected or deemed important. The algorithm was applied using the R

1 package *boruta* (28).

## 2 *Variable Selection Using Random Forests - VSURF*

3 Another method based on RF is the VSURF algorithm, which is short for "Variable Selection  
4 Using Random Forests" (29). The procedure consists of two steps. In the first step, the variables  
5 are ranked based on their importance, estimating a threshold value for variable importance (VI)  
6 using the standard deviation of VI for less important variables and retaining only the variables  
7 with an averaged VI value above the threshold. In the second step, a sequence of ascending RF  
8 models is used to make predictions. Variables are added to each model only if they significantly  
9 decrease the error rate, using a threshold based on the out-of-bag (OOB) error decrease. The final  
10 set of variables comes from the last model. In this study, we applied the VSURF method using the  
11 R package with the same name (29).

## 12 *Recursive Feature Elimination - RFE*

13 Recursive Feature Elimination (RFE) was first introduced by Guyon et al. (30). It is a method  
14 for feature selection similar to backward feature elimination (31) but allows for the elimination  
15 of multiple variables simultaneously instead of having to eliminate one feature at a time through  
16 exhaustive enumeration. In the RFE procedure, a model is first built on all features. In the second  
17 step, a ranked feature list is created by ranking the combination of each feature. Lastly, features  
18 are eliminated if they do not meaningfully contribute to the model. We applied the RFE method  
19 using the R package *caret* (32).

## 20 **Performance Metrics**

21 We utilize several quantitative metrics to assess the performance of the ML classifiers. They all  
22 rely on the true positive (TP), true negative (TN), false positive (FP), or false negative (FN) values  
23 in one way or another. In the context of this study, these values are defined as:

- 24 • true positive (TP): model correctly classifies a telecommuting day
- 25 as a telecommuting day
- 26 • true negative (TN): model correctly classifies a non-telecommuting day
- 27 as a non-telecommuting day
- 28 • false positive (FP): model incorrectly classifies a non-telecommuting day
- 29 as a telecommuting day
- 30 • false negative (FN): model incorrectly classifies a telecommuting day
- 31 as a non-telecommuting day

## 32 *Accuracy*

33 The accuracy of a classification model is the percentage of sample objects that are correctly classi-  
34 fied and labeled. This is done by calculating the ratio of the total number of true predictions to the  
35 sum of all observations.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (5)$$

## 36 *Precision*

37 Precision, also known as the positive predictive value (PPV), is defined as the ratio of correctly  
38 classified positive cases over all classified positive cases (33).



$$Precision = PPV = \frac{TP}{TP + FP} \quad (6)$$

#### 1 Recall

2 Recall, also referred to as sensitivity or true positive rate (TPR), is defined as the ratio of correctly  
3 classified positive cases over all actually positive cases (33). Recall and precision are often trade-  
4 offs of each other.

$$Recall = Sensitivity = TPR = \frac{TP}{TP + FN} \quad (7)$$

#### 5 Specificity

6 Specificity, also known as the true negative rate (TNR), is determined like the TPR, except that  
7 negative cases are now relevant. The TNR is defined as the ratio of correctly classified negative  
8 cases over all actually negative cases (33).

$$Specificity = TNR = \frac{TN}{TN + FP} \quad (8)$$

#### 9 F1 Score

10 The F1-score is determined by calculating the harmonic mean of the precision (eq. 6) and recall  
11 (eq. 7) (33). The F1-score can take values between 0 and 1, where 1 constitutes a perfect clas-  
12 sification. As can be seen in the formula, this is achieved if the sum of false positives and false  
13 negatives is zero.

$$F1 = \frac{2 \cdot TP}{2 \cdot TP + FP + FN} = \frac{2 \cdot Precision \cdot TPR}{Precision + TPR} \quad (9)$$

#### 14 Matthew's Correlation Coefficient

15 Although many studies use accuracy as the gold standard of model evaluation, it is very sensitive to  
16 unbalanced data, which can lead to a false sense of model performance (33). A way to counteract  
17 the issue of class imbalance when evaluating model performance is to calculate the Matthews  
18 Correlation Coefficient (MCC) (34). The MCC is calculated similarly to the Pearson product-  
19 moment correlation coefficient based on the confusion matrix of the model. The MCC can take  
20 on values between -1 and 1, where -1 is the worst possible value ( $TP = TN = 0$ ) and 1 is the best  
21 possible value (i.e.,  $FP = FN = 0$ ).

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP) \cdot (TP + FN) \cdot (TN + FP) \cdot (TN + FN)}} \quad (10)$$

## 22 RESULTS AND DISCUSSION

23 We first assess the results of the feature selection algorithms, which are presented in 1. Most  
24 strikingly is the similarity between boruta and rfe. The only difference in the two feature set is car  
25 access (included in rfe but not boruta) and the number of work-related trips (included in boruta but  
26 not in rfe). The vsurf feature set is the smallest, consisting of eight features. Features included  
27 in all three feature sets are *daily distance traveled*, *travel time*, *the hour of the last trip of the day*,

1 *the number of work trips the time spent at home, spent for leisure activities, spent shopping, and*  
 2 *on work-related activities..* While it is expected that time spent at home is an important feature,  
 3 we initially assumed that this would also be the case for time spent at work, which was deemed  
 4 unimportant by the vsurf algorithm.

**TABLE 1:** Selected features by selection method. Values indicate if feature was selected by the method (1) or not (0)

value	boruta	rfe	vsurf
age	1	1	0
car access	0	1	0
telecommuting Frequency	1	1	0
distance traveled	1	1	1
travel time	1	1	1
escorting someone	1	1	0
first trip of the day	1	1	0
last trip of the day	1	1	1
home	1	1	0
leisure	1	1	0
other	1	1	0
round trip	1	1	0
shopping	1	1	0
work	1	1	1
work-related	1	0	0
time use escorting someone	1	1	0
time use home	1	1	1
time use leisure	1	1	1
time use round trip	1	1	0
time use shopping	1	1	1
time use work	1	1	0
time use work-related	1	1	1

5 After selecting the features, we trained each model with the three feature sets on 70% of the  
 6 1-day HTS data. Subsequently, we tested the models on the remaining 30% by predicting whether  
 7 a respondent in the test data worked from home on the respective survey day. In order to measure  
 8 the performance of these models, we calculated the confusion matrix of predicted and real values,  
 9 which provided us with values for true positive, true negative, false positive, and false negative.  
 10 Figure 2 presents the confusion matrices for each classifier and each feature selection method.



**FIGURE 2:** Confusion matrix by classifier and feature selection algorithm

1 The highest true positive values are predicted by the random forest model with 238 cor-  
 2 rectly classified true values based on the boruta and rfe feature sets, and 231 based on the vsurf  
 3 feature set, respectively. All models have only few false negative vales, with bayesian networks  
 4 based on the rfe feature set. This seems to come at the price of a very high value for false positive  
 5 predictions. These are comparatively low in all other models, with random forest, again, perform-  
 6 ing best. Finally, true positive predictions are highest in the Bayesian network model based on  
 7 the rfe feature set and lowest also for the Bayesian network when considering the boruta feature  
 8 set. This is an interesting finding, as the boruta and rfe feature sets are almost identical (see Table  
 9 1), highlighting the need for pre-processing of data as Bayesian networks are often unsuited for  
 10 continuous data or outliers (35). All other models show almost identical confusion matrices for  
 11 these two feature sets.

12 To further assess the performance of the model, we put the values of each confusion matrix  
 13 into context with each other. The performance metrics for each model based on the values in the  
 14 confusion matrices are presented in table 2.

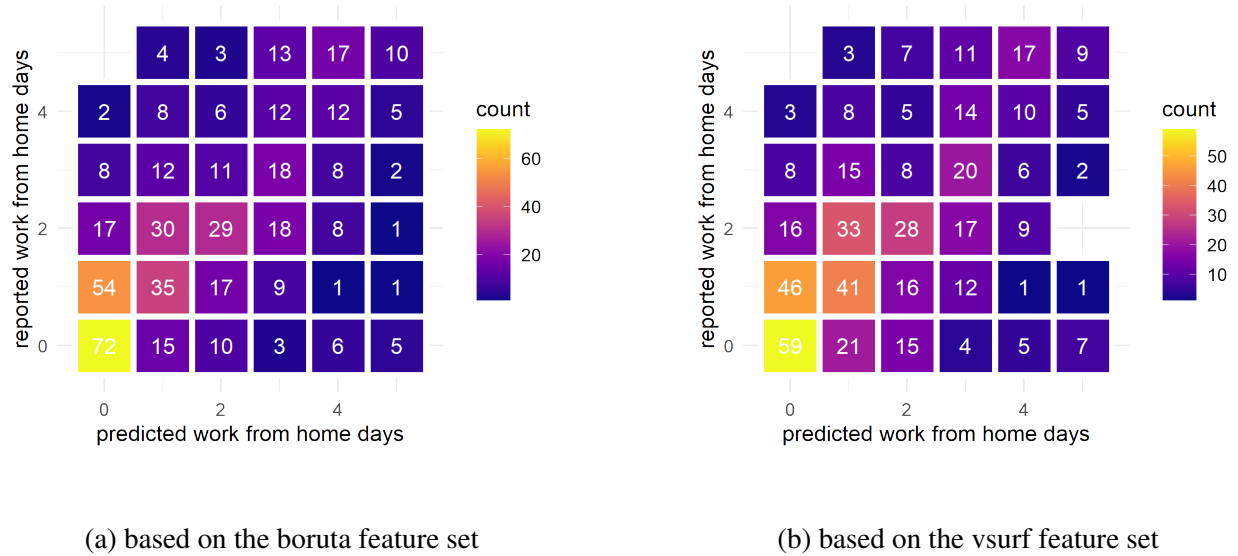
**TABLE 2:** Performance metrics by classification model and feature selection algorithm

classifier	feature selection	accuracy	precision	sensitivity	specificity	F1-Score	mcc
Bayesian Network	rfe	0.55	0.43	0.96	0.33	0.59	0.34
Random Forest	rfe	0.88	0.79	0.91	0.87	0.84	0.76
Support Vector Machnie	rfe	0.86	0.73	0.92	0.82	0.81	0.71
Binary Regression	rfe	0.86	0.75	0.91	0.84	0.82	0.72
Decision Trees	rfe	0.86	0.72	0.94	0.81	0.82	0.72
Bayesian Network	vsurf	0.86	0.73	0.92	0.83	0.82	0.71
Random Forest	vsurf	0.87	0.75	0.94	0.84	0.83	0.74
Support Vector Machnie	vsurf	0.87	0.74	0.95	0.83	0.83	0.74
Binary Regression	vsurf	0.85	0.72	0.94	0.81	0.81	0.71
Decision Trees	vsurf	0.86	0.72	0.94	0.81	0.82	0.72
Bayesian Network	boruta	0.77	0.65	0.70	0.80	0.68	0.50
Random Forest	boruta	0.89	0.79	0.92	0.88	0.85	0.77
Support Vector Machnie	boruta	0.87	0.74	0.93	0.83	0.82	0.73
Binary Regression	boruta	0.86	0.74	0.91	0.84	0.82	0.71
Decision Trees	boruta	0.86	0.72	0.94	0.81	0.82	0.72

Overall, we can see that almost all models achieved high rates of accuracy with only two models achieving an accuracy below 0.85. In both cases, bayesian network classification performed much worse compared to the other models. The recall/sensitivity metric shows even higher values, with, again, only the bayesian network model based on the features selected through the boruta algorithm reaching values below 0.91. Specificity/the true negative rate is not as high as the previous two metrics, but overall, almost consistently values of over .80 are reached. Regarding the F1-Score and the Matthew's Correlation Coefficient (MCC), a similar trend concerning Bayesian networks is detectable: overall, the metric values are relatively high and over 0.80 and 0.70, respectively. However, both metrics are much lower for the Bayesian network model based on the features selected through rfe and boruta. Our analysis indicates that the random forest model utilizing features from the boruta selection exhibited the most superior performance overall. This model achieved the highest accuracy, specificity, F1-Score, and mcc values in comparison to the other models. However, it did present a comparatively low sensitivity value. On the other hand, the bayesian network model based on rfe feature selection demonstrated the highest sensitivity value, but underperformed in all other metrics.

To further evaluate the performance of our proposed approach to data enhancement, we predicted the telecommuting engagement in a separate HTS. For this purpose, we leveraged data from the German Mobility Panel in which respondents keep a travel diary for seven days. As the random forest models performed best, we used those for prediction. To evaluate how the comparatively large number of features from the best model (boruta feature set) compares to the smaller feature set from the vsurf algorithm, we conducted two out-of-sample predictions. Because information on telecommuting engagement is provided at the week-level and not the day-l, we cannot calculate the aforementioned performance metrics. Instead, we predicted the telecommuting engagement for each day and added them over the week for each respondents to get a comparative measure. At the aggregate level over the entire dataset, the model based on the boruta feature set performs slightly better. This model detects 803 telecommuting days out of the 915 (87.7%) reported days. Whereas the model based on the smaller vsurf feature set predicts 792 work from home days (85.6%) The confusion matrices by the number of telecommuting days are presented

- 1 in Figure 3. On the left (Figure 3a), the results based on the boruta feature set are depicted, and
- 2 on the right (Figure 3b) those on the vsurf feature set. Values above the diagonal are likely false
- 3 negatives while values below the diagonal are most likely false positives.



**FIGURE 3:** Predicted and actual number of work-from-home days in German Mobility Panel dataset.

4 The two models perform very similar. Both have a high rate of predicted non-telecommuting  
 5 days over reported non-telecommuting days but also a relatively large rate of predicted non-  
 6 telecommuting days over the reported one day of telework per week. In the latter case, vsurf  
 7 performs slightly better than boruta. This performance difference is negated for a larger number of  
 8 telework days per week.

9 Overall, the application of the models on the MOP dataset shows promising results and  
 10 shows that the approach is viable to be applied to other HTS data. Random forest models are best  
 11 suited for this approach and perform well even on a relatively small feature set. To the best of our  
 12 knowledge, this is the first study testing different machine learning models to enhance in-home  
 13 activity information in HTS data. Thus, we cannot compare our results to other studies. However,  
 14 other studies comparing the performance of classification methods also find random forests to be  
 15 one of the best performing methods (36, 37).

## 16 CONCLUSION

17 In this study, we evaluated how machine learning methods can be leveraged to classify telecom-  
 18 muting engagement as an in-home activity to enhance data from a household travel survey. We  
 19 evaluated five classification methods, namely Bayesian Networks (BN), Binary Regression (BR),  
 20 Decision Trees (DT), Random Forest (RF), and Support Vector Machines (SVM). Overall, the NB  
 21 reached the highest specificity; however, it was outperformed by the other models on all other met-  
 22 rics. Overall, the RF model based on features selected using the *boruta* method performed best.  
 23 The out-of-sample prediction based on data from the German Mobility Panel shows promising

1 outcomes with an aggregate ratio of predicted telecommuting days over reported telecommuting  
2 days of 76%.

3       The research framework proposed is versatile and can be applied to various data sources  
4 and research questions. The models mainly utilized features that focused on the travel patterns  
5 of the respondents, making them suitable for GPS-based data. Variables that were particularly  
6 important in the models included daily distances traveled and time spent at home, which can easily  
7 be derived from GPS-based data. However, because additional features were still important in  
8 model prediction, we anticipate that the performance of ML models might be slightly lower than  
9 our application. Moreover, the method presented in this paper can be utilized to improve HTS data.  
10 The approach can be applied to time-use data, which will provide more detailed information on  
11 activities and enhance HTS data.

12       There are also some limitations worth noting. In our study, we applied the same feature  
13 selection algorithms for all machine learning methods without any feature engineering. This ap-  
14 proach was adopted to evaluate the performance of unprocessed data. However, this technique  
15 proved to be especially limiting for Bayesian network models. Therefore, future work will need to  
16 explore more advanced feature engineering techniques to optimize the performance of these mod-  
17 els. Further, the surveys used were not solely conducted for the purpose of this study, meaning this  
18 was not a controlled experiment. It would be interesting to repeat this effort with, e.g., GPS-based  
19 data that was enriched with information on telecommuting engagement as a controlled study to  
20 further evaluate the proposed approach.

21       Future work will, thus, include evaluating how feature engineering can improve the models.  
22 Further, we intend to apply the approach to other similar survey data to enhance existing data  
23 sources. Namely, we will evaluate how time-use data can be merged with household travel survey  
24 data to make the best of both worlds.

## 25 **ACKNOWLEDGEMENTS**

26 The data on which the models are trained was collected within the project "Traffic reduction  
27 through new forms of work and mobility technologies", funded by the German Federal Ministry  
28 of Education and Research under grant number 01UV2091C.

## 29 **AUTHOR CONTRIBUTIONS**

30 The authors confirm the contribution of the paper as follows: study conception and design: A.  
31 Reiffer; data collection: A. Reiffer, M. Kagerbauer; analysis and interpretation of results: A.  
32 Reiffer; draft manuscript preparation: A. Reiffer; supervision: M. Kagerbauer, P. Vortisch. All  
33 authors reviewed the results and approved the final version of the manuscript.

## 1 REFERENCES

- 2 1. Nilles, J., F. Carlson, P. Gray, and G. Hanneman, Telecommuting - An Alternative to  
3 Urban Transportation Congestion. *IEEE Transactions on systems, man, and cybernetics*,  
4 Vol. SMC-6, No. 2, 1976, pp. 77–84.
- 5 2. Elldér, E., Telework and Daily Travel: New Evidence from Sweden. *Journal of Transport*  
6 *Geography*, Vol. 86, 2020, p. 102777.
- 7 3. He, S. Y. and L. Hu, Telecommuting, Income, and out-of-Home Activities. *Travel Be-*  
8 *haviour and Society*, Vol. 2, No. 3, 2015, pp. 131–147.
- 9 4. Shi, Y., S. R. Sorrell, and T. J. Foxon, Do Teleworkers Have Lower Transport Emissions?  
10 What Are the Most Important Factors? *SSRN Electronic Journal*, 2022.
- 11 5. Asgari, H., X. Jin, and A. Mohseni, Choice, Frequency, and Engagement: Framework  
12 for Telecommuting Behavior Analysis and Modeling. *Transportation Research Record:*  
13 *Journal of the Transportation Research Board*, Vol. 2413, No. 1, 2014, pp. 101–109.
- 14 6. Reiffer, A., M. Magdolen, L. Ecke, and P. Vortisch, Effects of COVID-19 on Telework  
15 and Commuting Behavior: Evidence from 3 Years of Panel Data. *Transportation Research*  
16 *Record: Journal of the Transportation Research Board*, Vol. 2677, No. 4, 2023, pp. 478–  
17 493.
- 18 7. Asgari, H., X. Jin, and Y. Du, Examination of the Impacts of Telecommuting on the Time  
19 Use of Nonmandatory Activities. *Transportation Research Record: Journal of the Trans-*  
20 *portation Research Board*, Vol. 2566, No. 1, 2016, pp. 83–92.
- 21 8. Aschauer, F., R. Hössinger, K. W. Axhausen, B. Schmid, and R. Gerike, Implications  
22 of Survey Methods on Travel and Non-Travel Activities: A Comparison of the Austrian  
23 National Travel Survey and an Innovative Mobility-Activity-Expenditure Diary (MAED).  
24 *European Journal of Transport and Infrastructure Research*, Vol. 18, No. 1, 2018, pp.  
25 4–35.
- 26 9. Schmid, B., M. Balac, and K. W. Axhausen, Post-Car World: Data Collection Methods  
27 and Response Behavior in a Multi-Stage Travel Survey. *Transportation*, Vol. 46, No. 2,  
28 2019, pp. 425–492.
- 29 10. Alho, A., C. Cheng, D. T. Hieu, T. Sakai, F. Zhao, M. Ben-Akiva, and L. Cheah, Online and  
30 In-Person Activity Logging Using a Smartphone-Based Travel, Activity, and Time-Use  
31 Survey. *Transportation Research Interdisciplinary Perspectives*, Vol. 13, 2022, p. 100524.
- 32 11. Winkler, C., A. Meister, B. Schmid, and K. W. Axhausen, TimeUse+: Testing a Novel  
33 Survey for Understanding Travel, Time Use, and Expenditure Behavior, 2022, p. 17 p.
- 34 12. Verreault, H. and C. Morency, Integration of a Phone-Based Household Travel Survey and  
35 a Web-Based Student Travel Survey. *Transportation*, Vol. 45, No. 1, 2018, pp. 89–103.
- 36 13. Grapperon, A., B. Farooq, and M. Trépanier, *Information Fusion of Smart Card Data*  
37 *with Travel Survey*. CIRRELT - Interuniversity Research Centre on Enterprise Networks,  
38 Logistics and Transportation, 2016.
- 39 14. Hossain, S. and K. N. Habib, Inferring the Purposes of Using Ride-Hailing Services  
40 through Data Fusion of Trip Trajectories, Secondary Travel Surveys, and Land Use Data.  
41 *Transportation Research Record: Journal of the Transportation Research Board*, Vol.  
42 2675, No. 9, 2021, pp. 558–573.
- 43 15. Bwambale, A., C. F. Choudhury, S. Hess, and M. S. Iqbal, Getting the Best of Both  
44 Worlds: A Framework for Combining Disaggregate Travel Survey Data and Aggregate

- 1 Mobile Phone Data for Trip Generation Modelling. *Transportation*, Vol. 48, No. 5, 2021,  
2 pp. 2287–2314.
- 3 16. Maghrebi, M., A. Abbasi, T. H. Rashidi, and S. T. Waller, Complementing Travel Diary  
4 Surveys with Twitter Data: Application of Text Mining Techniques on Activity Location,  
5 Type and Time. In *2015 IEEE 18th International Conference on Intelligent Transportation*  
6 *Systems*, IEEE, Gran Canaria, Spain, 2015, pp. 208–213.
- 7 17. Verreault, H. and C. Morency, Methodology to Add Value to Ageing Travel Survey Data.  
8 *Transportation Research Record: Journal of the Transportation Research Board*, 2023, p.  
9 036119812311591.
- 10 18. Wang, K., S. Hossain, and K. N. Habib, A Hybrid Data Fusion Methodology for House-  
11 hold Travel Surveys to Reduce Proxy Biases and Under-Representation of Specific Sub-  
12 Group of Population. *Transportation*, Vol. 49, No. 6, 2022, pp. 1801–1836.
- 13 19. Koushik, A. N., M. Manoj, and N. Nezamuddin, Machine Learning Applications in  
14 Activity-Travel Behaviour Research: A Review. *Transport Reviews*, Vol. 40, No. 3, 2020,  
15 pp. 288–311.
- 16 20. Chawla, N. V., K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, SMOTE: Synthetic  
17 Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, Vol. 16,  
18 2002, pp. 321–357.
- 19 21. Pearl, J., *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*.  
20 The Morgan Kaufmann Series in Representation and Reasoning, Morgan Kaufmann, San  
21 Francisco, Calif, rev. 2. ed., transferred to digital printing ed., 1997.
- 22 22. Neapolitan, R. E., *Probabilistic Reasoning in Expert Systems: Theory and Algorithms*.  
23 publisher not identified, Place of publication not identified, 2012.
- 24 23. Dattatreya, G. R. and L. Kanal, Decision Trees in Pattern Recognition. University of Mary-  
25 land, 1984.
- 26 24. Tin Kam Ho, Random Decision Forests. In *Proceedings of 3rd International Conference on*  
27 *Document Analysis and Recognition*, IEEE Comput. Soc. Press, Montreal, Que., Canada,  
28 1995, Vol. 1, pp. 278–282.
- 29 25. Tin Kam Ho, The Random Subspace Method for Constructing Decision Forests. *IEEE*  
30 *Transactions on Pattern Analysis and Machine Intelligence*, Vol. 20, No. 8, Aug./1998,  
31 pp. 832–844.
- 32 26. Cortes, C. and V. Vapnik, Support-Vector Networks. *Machine Learning*, Vol. 20, No. 3,  
33 1995, pp. 273–297.
- 34 27. Kursa, M. B., A. Jankowski, and W. R. Rudnicki, Boruta – A System for Feature Selection.  
35 *Fundamenta Informaticae*, Vol. 101, No. 4, 2010, pp. 271–285.
- 36 28. Kursa, M. B. and W. R. Rudnicki, Feature Selection with the Boruta Package. *Journal of*  
37 *Statistical Software*, Vol. 36, 2010, pp. 1–13.
- 38 29. Genuer, R., J.-M. Poggi, and C. Tuleau-Malot, VSURF: An R Package for Variable Selec-  
39 tion Using Random Forests. *The R Journal*, Vol. 7, No. 2, 2015, pp. 19–33.
- 40 30. Guyon, I., J. Weston, S. Barnhill, and V. Vapnik, Gene Selection for Cancer Classification  
41 Using Support Vector Machines. *Machine Learning*, Vol. 46, No. 1/3, 2002, pp. 389–422.
- 42 31. Kohavi, R. and G. H. John, Wrappers for Feature Subset Selection. *Artificial Intelligence*,  
43 Vol. 97, No. 1-2, 1997, pp. 273–324.
- 44 32. Kuhn, M., Building Predictive Models in R Using the **Caret** Package. *Journal of Statistical*  
45 *Software*, Vol. 28, No. 5, 2008.



- 1 33. Chicco, D. and G. Jurman, The Advantages of the Matthews Correlation Coefficient  
2 (MCC) over F1 Score and Accuracy in Binary Classification Evaluation. *BMC Genomics*,  
3 Vol. 21, No. 1, 2020, p. 6.
- 4 34. Baldi, P., S. Brunak, Y. Chauvin, C. A. F. Andersen, and H. Nielsen, Assessing the Ac-  
5 curacy of Prediction Algorithms for Classification: An Overview. *Bioinformatics*, Vol. 16,  
6 No. 5, 2000, pp. 412–424.
- 7 35. Cheng, J. and R. Greiner, Comparing Bayesian Network Classifiers, 2013.
- 8 36. Zhang, C., C. Liu, X. Zhang, and G. Almpanidis, An Up-to-Date Comparison of State-  
9 of-the-Art Classification Algorithms. *Expert Systems with Applications*, Vol. 82, 2017, pp.  
10 128–150.
- 11 37. Chen, R.-C., C. Dewi, S.-W. Huang, and R. E. Caraka, Selecting Critical Features for Data  
12 Classification Based on Machine Learning Methods. *Journal of Big Data*, Vol. 7, No. 1,  
13 2020, p. 52.