

1 **BUILT ENVIRONMENT FACTORS AFFECTING BIKE SHARING RIDERSHIP:**  
2 **A DATA-DRIVEN APPROACH FOR MULTIPLE CITIES**

3

4

5

6 **David Duran-Rodas**

7 Technical University of Munich

8 Arcisstrasse 21, 80333 Munich

9 E.mail: david.duran@tum.de

10

11 **Emmanouil Chaniotakis**

12 Technical University of Munich

13 Arcisstrasse 21, 80333 Munich

14 E.mail: m.chaniotakis@tum.de

15

16 **Constantinos Antoniou**

17 Technical University of Munich

18 Arcisstrasse 21, 80333 Munich

19 E.mail: c.antoniou@tum.de

20

21

22 Word Count: 7054 words + 1 table  $\times$  250 = 7304 words

23

24

25

26

27

28

29 Submission Date: August 1, 2018

**1 ABSTRACT**

2 Bike sharing has been found to present environmental, economic and social benefits. Identifica-  
3 tion of factors influencing ridership is necessary for policy-making, as well as when examining  
4 transferability and aspects of performance and reliability. In this work, a data-driven method is  
5 formulated to correlate arrivals and departures of station-based bike sharing systems with built  
6 environment factors in multiple cities. Ridership data from stations of multiple cities are pooled  
7 in one data-set regardless of their geographic boundaries. The method bundles the collection,  
8 analysis, and processing of data, as well as, the models' estimation using statistical and machine  
9 learning techniques. The method was applied on a national level in six cities in Germany, and also,  
10 on an international level in three cities in Europe and North America. The results suggest that the  
11 models' performance did not depend on clustering cities by size but by the relative daily distri-  
12 bution of the rentals. Selected statistically significant factors were identified to vary temporally  
13 (e.g. nightclubs were significant during the night). The most influencing variables were related  
14 to the city population, distance to city center, leisure-related establishments and transport related  
15 infrastructure. This data-driven method can help as a support decision-making tool to implement  
16 or expand bike sharing systems.

17

18 *Keywords:* bike sharing, built environment, open-source, multiple cities comparison

## 1 INTRODUCTION

2 Bike sharing is defined as the shared use of a bicycle, where a user accesses a fleet of bicycles  
3 offered on public space (1). It is part of the shared economy social–economic phenomenon, where  
4 individuals or organizations prioritize use over ownership of items (2). Bike sharing systems have  
5 a long history, with the very first system launched in 1965. Its deployment was in Amsterdam  
6 with fifty free and unlocked bicycles. Theft and vandalism led to a coin–deposit system, also  
7 not successful, mainly due to the user’s anonymity. Nowadays, information and communications  
8 technology (ICT) enables wireless pick–up, drop–off, and a real–time GPS tracking of bicycles  
9 (3), which lead to the widespread of bike sharing to more than 1,600 cities around the world (4).  
10 (4). Europe and Asia are the continents with the majority of bike sharing systems worldwide. In  
11 2015, China presented the biggest fleet in the world with 753,508 bicycles, followed by France  
12 with 42,930, and Spain with 25,084 (4). Categorization of bike sharing systems can be defined by  
13 the use of stations or not: (a) station-based (SBBS), b) free-floating (FFBS) and c) a mix of the  
14 two (5).

15 The wide deployment and observed growing trends of bike sharing can be attributed, among  
16 others, to its associated social, economic and environmental benefits. These are related to creat-  
17 ing a larger cycling population, cost savings, increasing transit use, reducing greenhouse gases,  
18 decreasing congestion, creating environmental awareness, improving public health, among others  
19 (a comprehensive review of benefits attributed to bike sharing can be found in (3), (6) and (7)).  
20 However, not all systems were deployed successfully. Some were perceived as a public nuisance or  
21 were misused and vandalized (8). Possible reasons for a system failure were bicycles’ poor qual-  
22 ity, lack of funding, oversaturated market, delayed expansion, inconvenient system design, unfair  
23 fares, low political support (8, 9).

24 The identified benefits strongly suggest the necessity to further increase the use of bike  
25 sharing systems and to enable their deployment in more cities. At the same time, the unsuccessful  
26 deployment of some projects makes the examination of the factors that affect ridership and system  
27 reliability rather imperative. These two needs have been the driving force for a high number of  
28 studies on the influencing factors that affect the bike sharing usage (e.g., built environment, socio-  
29 demographic characteristics, system settings). Most studies analyze the influencing factors in a)  
30 multiple cities, with each city considered as one observation (10, 11) or b) single city at a local  
31 (station) level, where one city is analyzed and observations are based on an area of influence, e.g.  
32 near stations (12–15).

33 The multiple–cities approach suffers an exclusion of varying characteristics within a city,  
34 which provides an indication of how the system should be structured to enable a successful de-  
35 ployment. Conversely, the station level approach is performed in a single city and is bounded by  
36 the urban settings examined. The main issue with this approach is that it does not examine the  
37 system’s transferability but the ridership within a city.

38 Aiming at overcoming the above–discussed drawbacks of existing approaches, this paper  
39 contributes to the related literature, by focusing on the investigation of bike sharing systems as  
40 one entity regardless of the city they belong to. We present a multiple cities data-driven approach  
41 focusing on the comparison of general built environment characteristics on a station–level influ-  
42 encing bike sharing systems beyond geographic boundaries. As such, the data used for each city  
43 is pooled in one complete data-set, where each observation refers to one station’s area of influence  
44 (as defined in the Method Section). The influencing factors chosen to be investigated describe  
45 the characteristics of the built environment (guided by the high influence found in the majority of

1 previous studies).

2       The second main contribution of this study lies upon the modeling techniques used for the  
3 most influencing factors selection. As discussed in the Related Work Section, in most cases a  
4 predetermined set of factors is used. This set is hypothesized to contribute to a successful deploy-  
5 ment of bike sharing, and thus, could omit possible patterns revealed by an alternative selection  
6 approach. In this study, a data-driven approach is followed to allow the discovery of factors that  
7 might not be commonly addressed. This is done by using different linear and non-linear model-  
8 ing techniques which are evaluated upon modeling performance criteria such as goodness-of-fit,  
9 information criteria, and (cross-)validation.

10       Two applications of the methods discussed are included: a) a national application in six  
11 cities in Germany, Europe and b) an international application in three cities from Europe and  
12 North America with similar urban characteristics. The first application intends to illustrate the per-  
13 formance of different modeling techniques, while the second intends to illustrate the applicability  
14 of the methods in an international setting, while taking into account seasonality. In both cases,  
15 different validation techniques are exercised with very positive results. All the factors are defined  
16 using open-source data and by the derivation and deployment of an automated feature creation  
17 methodology.

## 18 **RELATED WORK**

19 The spatial-temporal factors influencing historical rentals of SBBS systems have been studied  
20 in cities all over the world with different sizes. The resulting factors have been compared be-  
21 tween cities with factors within a city scale (10, 11) or in a single city on a station scale (12–19).  
22 The modeling approach followed in most of the cases above can be summarized as a) the model  
23 estimation method used is mainly a linear regression using ordinary least squares (13, 15, 17–  
24 19); b) the dependent variable is the logarithm of the number or rates of arrivals and departures  
25 (12, 13, 15, 17, 19); c) the model assessment is usually performed using the indexes: log-likelihood  
26 (LL),  $R^2$  and AIC-BIC.

27       Regarding the multiple cities approach, De Chardon et al. (10) studied the trips per day per  
28 bicycle (TDB) in 75 SBBS systems in Europe, Israel, United States, Canada, Brazil and Australia.  
29 They used a robust regression to build the model with the logarithm of the TBD as the depen-  
30 dent variable. The resulting influencing variables were the operator's attributes, the compactness,  
31 the weather, the transportation infrastructure, as well as system-related characteristics, such as  
32 helmet requirement and the number of docks at stations. Faghih-Imani et al. (15) aggregated to  
33 an hourly value arrivals and departures in Barcelona and Seville into a Sub-City District level.  
34 They correlated both cities separately the of the logarithm of the dependent variable linearly to so-  
35 ciodemographic and socioeconomic variables and Points Of Interest (POI). Barcelona and Seville  
36 presented a similar pattern where the common influencing POIs were related to business, leisure,  
37 and restaurants.

38       On the other hand, considering a single city approach, Dung Tran et al. (14) developed  
39 models for bike sharing in Lyon using rain stations, restaurant, cinema and embankment roads,  
40 altitude, among others. They also found that the population density showed a positive effect in  
41 the morning and the number of jobs with a positive impact in the afternoon. Faghih-Imani and  
42 Eluru (12) correlated the hourly arrivals and departures for one month in the SBBS "CitiBike" in  
43 New York with temporal, spatial and weather variables. They concluded that the fit of the model  
44 improved significantly by adding temporally and spatially lagged dependent variables. The length

1 of bicycle routes, the presence of subway stations, the area of parks on weekends, and the number  
 2 of restaurants increased the usage of the system, while the length of railways decreased it.

3 Activities that are associated with bike sharing are commuting and leisure (20). These  
 4 activities are vastly related to the built environment. Thus, many studies have examined its rela-  
 5 tionship to bike sharing upon the transport infrastructure, Points Of Interest (POI) and the land use  
 6 categories (Table 1).

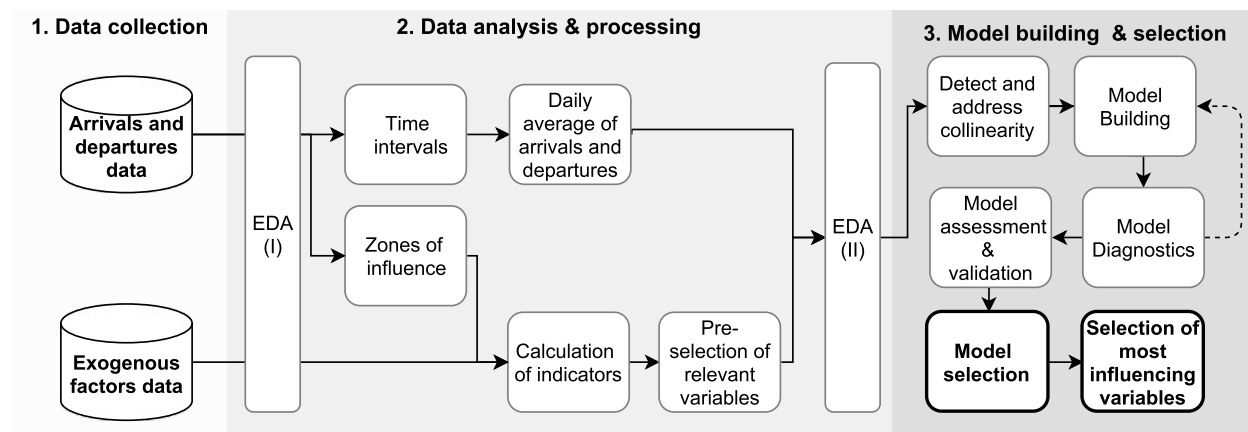
**TABLE 1:** Built environment factors most influencing bike sharing historical data

CATEGORY	VARIABLE	REFERENCES								
		(10)	(12)	(16)	(17)	(13)	(18)	(14)	(15)	(19)
<b>Transport</b>	Cycling infrastructure	✓	✓							
	Railways length		✓							
	Subway stations		✓	✓		✓		✓	✓	
	Rail stations							✓		
<b>POIs</b>	Universities					✓	✓			✓
	Student residence							✓		
	Restaurants		✓		✓		✓	✓		
	Cinema							✓		
	CBD				✓		✓			
	Number of business				✓		✓		✓	
<b>Land use</b>	Parks		✓							
	Residential land use			✓						
	Parking land use			✓						
	Bodies of water				✓					

7 To the best of the authors' knowledge, there is no data-driven method to measure built en-  
 8 vironment variables by assigning them automatically different types of indicators as quantity in the  
 9 area of influence of the stations or proximity to the stations. Contrary to De Chardon et al. (10), we  
 10 analyzed multiple cities, but on a station scale. Also, instead of comparing the influencing factors  
 11 of multiple cities (15), we model the cities together. Finally, in the literature review, there was not  
 12 a comparison of different linear and non-linear modeling techniques to define which technique fits  
 13 better the bicycle usage and the influencing built environment.

## 14 **METHOD**

15 The proposed method aims at building models automatically in different temporal scales to identify  
 16 the built environment variables that influence the historical rentals of SBBS systems in multiple  
 17 cities. The method goes through three main components: 1) automated data collection, 2) auto-  
 18 mated data analysis and processing, and 3) automated model building and selection of the modeling  
 19 technique with the better fitting results, and automated selection of the most influencing variables  
 20 (Figure 1).



**FIGURE 1:** Methodological framework

### 1 Data collection, analysis, and processing

2 Data collection is performed on historical arrivals and departures from bike sharing systems (de-  
 3 pendent variables) and the built environment (independent variables) in multiple cities. The inde-  
 4 pendent variables are points, lines, and polygons of the built environment: e.g., points of interest,  
 5 public transport stations, railways, roadways, waterways, land use, and natural features.

#### 6 *Ridership Data*

7 The historical ridership data (in terms of arrivals and departures to and from a station) are explored  
 8 to define time intervals to build models independent of time. This is performed to allow homogene-  
 9 ity in terms of dependent variables and to correct the effects of the time of the day. A clustering  
 10 analysis is carried out to determine which days of the week illustrate significantly different rider-  
 11 ship patterns. In each cluster, different periods are identified based on the hourly distribution of  
 12 the rentals based on peak and off-peak times. The cumulative ridership variable (dependent) and  
 13 built environment variables are aggregated on a spatial scale, based on zones of influence. These  
 14 zones are defined as the maximum area of influence that an individual is willing to walk to reach a  
 15 bike-sharing station. Their boundaries are defined as the intersection of the Thiessen polygons of  
 16 the stations, human-made and natural barriers and a buffer circumference from the stations repre-  
 17 senting the maximum walking distance [200 to 400 meters (10, 13, 14, 16–18, 21)] that a station  
 18 can attract or produce.

#### 19 *Built environment Data*

20 Built environment data is downloaded from an open-source database on a city level. Each built  
 21 environment variable is assigned two indicators in each zone of influence: 1) proximity-based  
 22 indicators (minimum distance from a station to the examined spatial feature inside the zone of  
 23 influence, and 2) quantity or presence of the variable in a zone of influence. The selection of the  
 24 appropriate type of indicator is decided based on some basic hypotheses. Let  $v$  be a random (inde-  
 25 pendent) variable used to describe a particular built environment distribution across observations.  
 26 Also, let  $\sigma_v$  represent its standard deviation. A variable is defined as static if the standard deviation  
 27 is smaller than a threshold  $t$  ( $\sigma_v < t$ ). Under the above hypotheses, indicators will be introduced  
 28 in the model as dummy variables indicating "presence" rather than quantity. A sensitivity analysis  
 29 is carried out to determine the value of the threshold of the standard deviation. Only the variables

1 that are present in all the cities of the study are considered. These variables are explored to exclude  
 2 those presenting inconsistencies or irrelevant to the influence of bike sharing ridership.

3 Finally, Pearson and Spearman correlation tests are carried out to determine the variables  
 4 that are collinear. If two variables are collinear, the variable that influences more the rentals is  
 5 considered (multiple regression models are estimated).

## 6 **Model building and selection**

7 Model building and selection is based on a sequential model definition and model validation pro-  
 8 cess. The aim pursued is to build linear and non-linear models to identify the models that better  
 9 fit the dataset, while a) being parsimonious without a substantial loss of their fitting performance,  
 10 b) avoiding over-fitting and c) including variable selection or variable categorization for compu-  
 11 tational efficiency, given a large number of independent variables. Mathematical transformations  
 12 of the variables are considered to handle heteroscedasticity or non-linearity issues and to improve  
 13 the models. Most common transformations are the square root, logarithmic, inverse, exponential,  
 14 arcsine (22), or Box-Cox transformation (23).

### 15 *Model Structures*

16 The model building techniques to be examined are stepwise Ordinary Least Square regression  
 17 (stepwise OLS) (24), Generalized Linear Models (GLM) with a lasso selection technique (25), and  
 18 Gradient Boosting Machine (GBM) (26).

19 Stepwise OLS for variables selection is chosen based on its wide use in the pertinent liter-  
 20 ature for similar types of problems (11, 13, 15, 17–19). The core of stepwise OLS is the multiple  
 21 linear regression, which is iteratively used to build a model using an observations vector  $Y$  that is  
 22 linearly related to a matrix  $X$  (independent variables) and  $\varepsilon$  residuals [ $Y = X \cdot \beta + \varepsilon$  (24)]. Stepwise  
 23 regression addresses the subset selection of a large number of  $k$  parameters. There are three types  
 24 of stepwise selection procedures: 1) forward selection, 2) backward selection 3) both directions  
 25 (27). The forward selection initiates with only the constant term (i.e., no parameters) and adds  
 26 variables based on a comparison criterion. The backward elimination process, in contrast, starts  
 27 with a full equation and excludes the uncorrelated parameters. The stepwise method in both direc-  
 28 tions sequentially adds or deletes parameters. It starts with a forward selection, but at each step, it  
 29 can remove a parameter. Its advantage is if a non-significant parameter is included in the process,  
 30 it might be eliminated later.

The selection of the parameters is based on criteria to compare the regression in each step.  
 Commonly used criteria are the Akaike Information Criterion (AIC) and Bayes Information Crite-  
 rion (BIC) (24). AIC (28) is defined as:

$$AIC = n * \ln(MSE) + 2 \cdot k \quad (1)$$

where  $n$  denotes the number of observations,  $MSE$  the mean squared error and  $k$  the number  
 of parameters. A direct implication of using AIC is that for two models with the same error, AIC  
 would penalize the one with more parameters. However, the use of AIC tends to improve with  
 a larger number of  $k$  parameters, thus it is commonly accused of being prone to overfit models  
 selection. BIC (29) tends to control the overfitting of AIC. It is proportional to AIC but it uses a  
 logarithmic factor for the effect that the number of variables has:

$$BIC = n * \ln(MSE) + \ln(n) \cdot k \quad (2)$$

1 Generalized linear models (GLM) are an extension of OLS. Usually, this family of models  
 2 is based on the maximum likelihood estimation. GLM assume that the error  $\varepsilon$  presents a distri-  
 3 bution from the exponential family, such as binomial, Poisson, Gaussian. Also, they consider the  
 4 mean function  $\mu_i$  as a function of the linear observations [ $h(\mu_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}$   
 5 where,  $h(\mu_i)$  is a function that links  $\mu_i$  with the observation  $Y_i$ ] (24). The least absolute shrinkage  
 6 and selection operator (lasso) technique (25) shrinks the coefficients  $\beta$  increasing stability while  
 7 retaining the best variables. Lasso assumes that  $X_{i,j}$  are standardized with a mean of zero and a stan-  
 8 dard deviation of 1. Then, it minimizes the sum of the squared differences between the observation  
 9 and the linear regression [ $\hat{\beta}_{lasso} = \underset{\beta}{\operatorname{argmin}} (\sum_{i=1}^n (Y_i - \beta_0 - \sum_{j=1}^k X_{ij} \beta_j)^2 + \lambda \sum_{j=1}^k |\beta_j|)$ ]. Models  
 10 generated with lasso are easy to interpret. The selection of  $\lambda$  is calculated after a cross-validation  
 11 test to select the  $\lambda$  that presents the smallest error (27).

12 Gradient Boosting Machine (GBM) is a machine learning algorithm that performs regres-  
 13 sion, classification, and ranking (26). It is a mix of boosting and gradient steepest descent. Boost-  
 14 ing is a procedure to reduce the variance of a model. It involves the creation of multiple  $B$  training  
 15 sets. Then, it builds a prediction for each training set  $\hat{f}^1(x), \hat{f}^2(x), \dots, \hat{f}^B(x)$  and it fits different  
 16 decision trees to each copy. Each tree is a modified version of the original data set, and they grow  
 17 sequentially by using the information of the previously grown tree. The residuals are fit to the de-  
 18 cision tree, rather than a single decision tree to the data. We choose the sample data that modeled  
 19 poorly in the system before, i.e., in areas where the system is not performing well. Then, the resid-  
 20 uals are updated after adding the new decision tree into the fit function. Finally, it combines all the  
 21 trees to create a single model. A faster approximation to find the model is to consider a differen-  
 22 tiable loss criterion that can be derived by numerical optimization. Regarding the loss function, a  
 23 Gaussian function is used for numerical efficiency to minimize the squared error and the Laplace  
 24 for minimizing the absolute error.

25 GBM is considered in the study because the dataset might fit better in a nonlinear model.  
 26 It uses the input arguments: loss function, number of iterations, terminal nodes of each tree and  
 27 shrinkage factor (30). A sensitivity analysis has to be carried out to determine these values. In  
 28 addition to the resulting model, GBM provides a ranking list of the variables with their relative  
 29 influence normalized to sum one hundred. To carry out a variable selection from the ranking list,  
 30 mean square errors (MSEs) are calculated starting from highest ranked variable and then adding a  
 31 subsequent variable until a non-significant difference of the MSE is present. GLM and GBM have  
 32 shown high fitting performance in similar applications in the literature, e.g., (31).

33 Model building is carried out with a training set and the model validation with a testing set  
 34 for each time unit and for the linear and non-linear regression models. After the models are built,  
 35 two types of criteria are used to assess them: a) Indirect methods: lowest number of predictors,  
 36 lowest Mean Square Error (MSE), lowest BIC, and greatest goodness of fit measures ( $R^2$  and  
 37 adjusted  $R^2 - R_{adj}^2$ ); and b) Best validation results: selection of the model that adequately predicts  
 38 the arrivals and departures on a validation dataset.

### 39 APPLICATION

40 The data-driven method was applied in two cases: 1) national level in six German cities and 2)  
 41 international level with three cities in three countries in Europe and North America. The national  
 42 level application provides evidence on the applicability and performance of the different model  
 43 structures and estimation techniques, allowing for a more comprehensive evaluation of the impact  
 44 that different techniques have to the identification of the factors affecting ridership. Germany has



1 been used as the national case, due to the bike sharing fleet (fifth largest fleet in the world; approx.  
2 12,000 shared bicycles) (4) and data availability. The international level application builds on the  
3 first application and focuses on the extraction of conclusions for the application of the methods in  
4 an international comparison.

### 5 **Multiple National Cities approach**

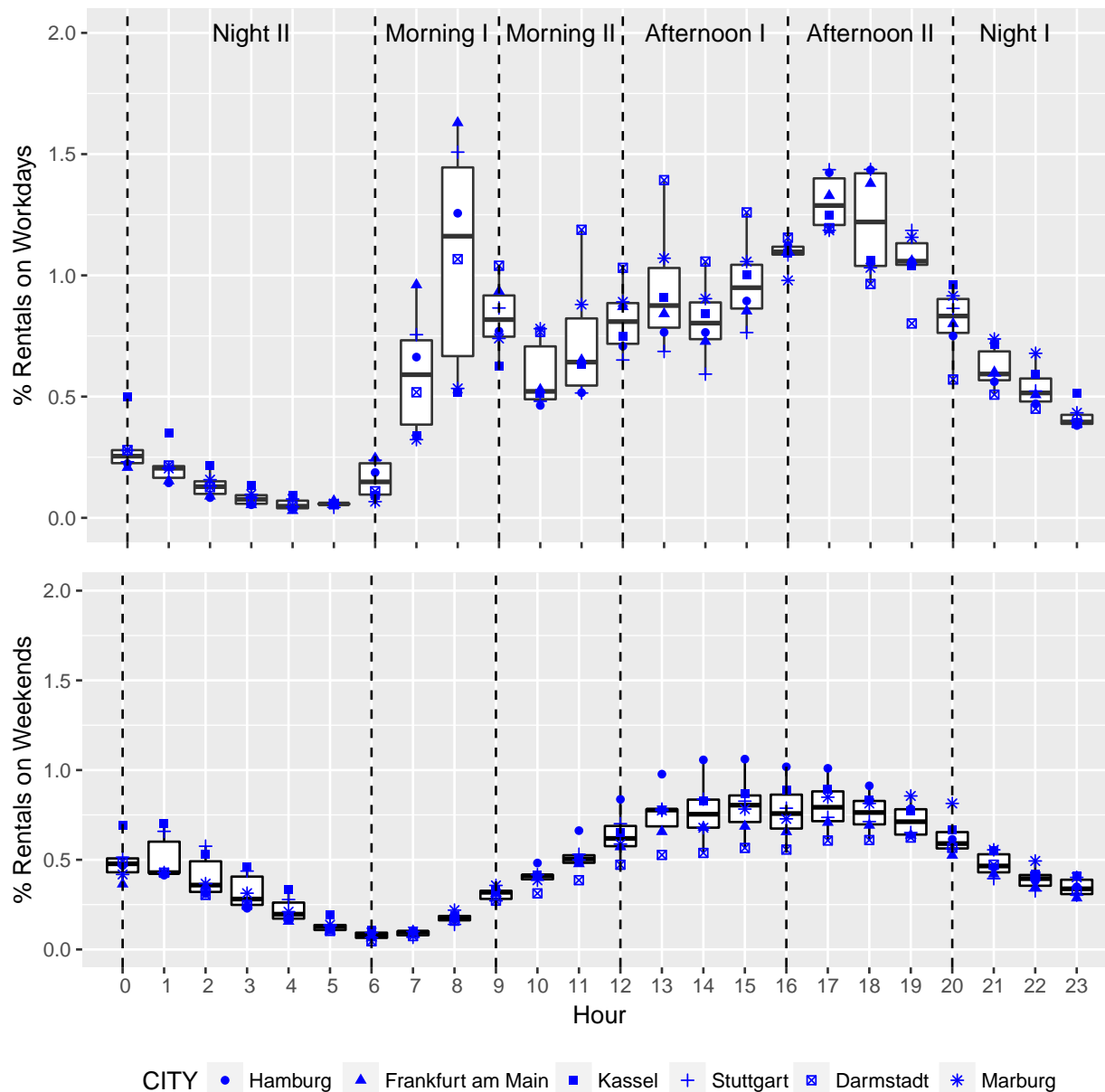
6 This approach includes six German cities for the SBBS system "Call a bike" (32): Hamburg, Frank-  
7 furt am Main, Stuttgart, Kassel, Darmstadt, and Marburg. Arrivals and departures of the bicycles  
8 were downloaded from the Open-Data-Portal offered by the German train company (Deutsche  
9 Bahn) under the link: <http://data.deutschebahn.com/dataset/data-call-a-bike> on June  
10 2017. The dataset included the rentals in fifty cities in Germany for approximately 3.5 years. The  
11 majority of the data, however, referred to the six selected cities, because of their high usage of bike  
12 sharing (>250,000 rentals in total, around 3.5 GB). In total, 10.5 million rentals were included  
13 in the dataset referring to the period between 01.01.2014 and 15.05.2017 (1232 days). Around  
14 73% of the rentals referred to the city of Hamburg, followed by Frankfurt with 12%. Peaks were  
15 identified in the summertime (May to July). It is worth mentioning that Wednesdays and Thurs-  
16 days showed the highest ridership, which was found to decrease during the weekends. Regarding  
17 the hourly distribution, there was a different trend between workdays and weekends (Figure 3). In  
18 workdays, there were two peak periods at 8:00 and at 17:00 (based on the median values). Figure 2  
19 shows the spatial distribution of the intensity of the rentals. Each area represents the frequency of  
20 a station with the help of Voronoi Diagrams for better visualization.

### 21 *Data analysis and processing*

22 The rentals were clustered into days of the week using Pearson correlation analysis. The three re-  
23 sulting clusters were workdays, Saturdays, and Sundays. Arrivals and departures were aggregated  
24 into time intervals representing peak and off-peak periods in the morning, afternoon and night  
25 (Figure 3). The built environment variables were downloaded from the collaborative open-source  
26 dataset OpenStreetMaps (33). Unclassified roads and a selection of variables, which were found  
27 to be inaccurately positioned or irrelevant, were excluded (e.g., vending machines, wastebaskets).  
28 The distance to the city center was also considered as a built environment variable since it was  
29 present in the literature review. For the zones of influence, a 300 meters buffer ratio was used  
30 as being the most common value used in the literature. Four indicators were assigned to around  
31 200 types of spatial features (around **800** spatial variables). A threshold value selected of  $SD = 5$   
32 was selected after a sensitivity analysis to determine if the indicator of a variable is related to the  
33 quantity or presence. A total of **194** variables were examined with **144** non-collinear variables to  
34 be selected after Pearson and Spearman correlation tests. A threshold value of 0.7 was considered  
35 as explained in Zhao et al (11).

### 36 *Model building, diagnosis and validation*

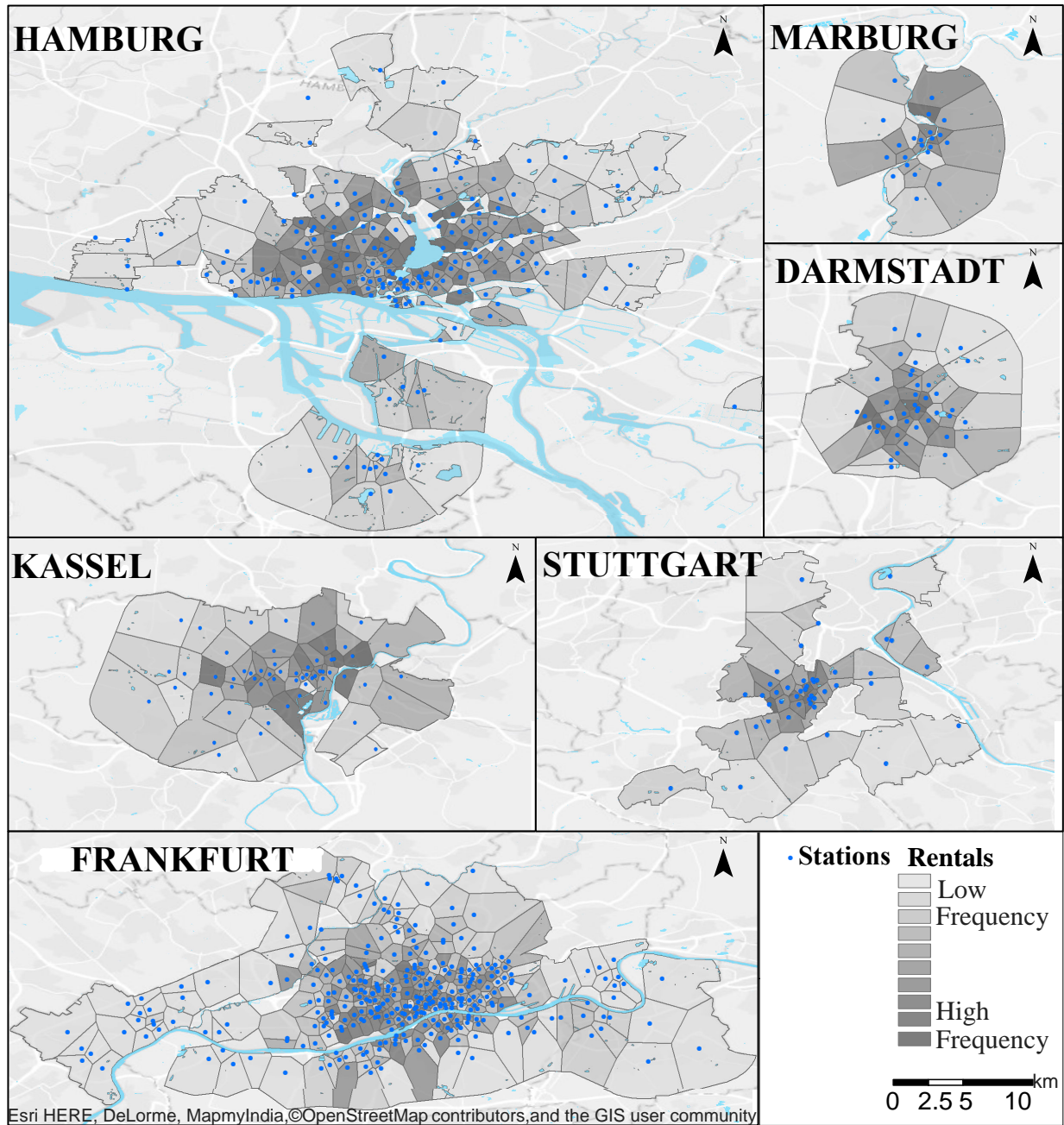
37 Aiming at examining applicability and performance of the different model structures and estima-  
38 tion techniques, all methods discussed in the Methods section were used (OLS, GLM with lasso  
39 and GBM). In all cases, the relationship between arrivals and departures to 144 non-collinear built  
40 environment variables was examined. The city population was used to weight ridership (for dif-  
41 ferent cities' sizes) (34). Apart from model fitting and model diagnostics, model validation was  
42 performed by dividing ridership data into a training set including the zone of influence of 5 cities



**FIGURE 2:** Hourly distribution and definition of times intervals

1 and a testing set of one city’s zone. Validation was performed on a city level (and not using a ran-  
 2 dom sample of zones) in order to examine how well the models would perform in a German city  
 3 without a bike-sharing system. The city of Kassel was chosen for validation due to its high rider-  
 4 ship and the fact that Hamburg and Frankfurt were not considered because they involved together  
 5 around 76% of the zones of influence.

6 Stepwise OLS was considered in both directions, while BIC was chosen as a selection cri-  
 7 terion. For GLM with lasso models, a Gaussian distribution was considered because it fit better  
 8 the training data. A k-folds cross-validation (35) was implemented to calculate the shrinkage fac-  
 9 tor that helped the models to fit better the data. Concerning GBM, K-fold cross-validation was

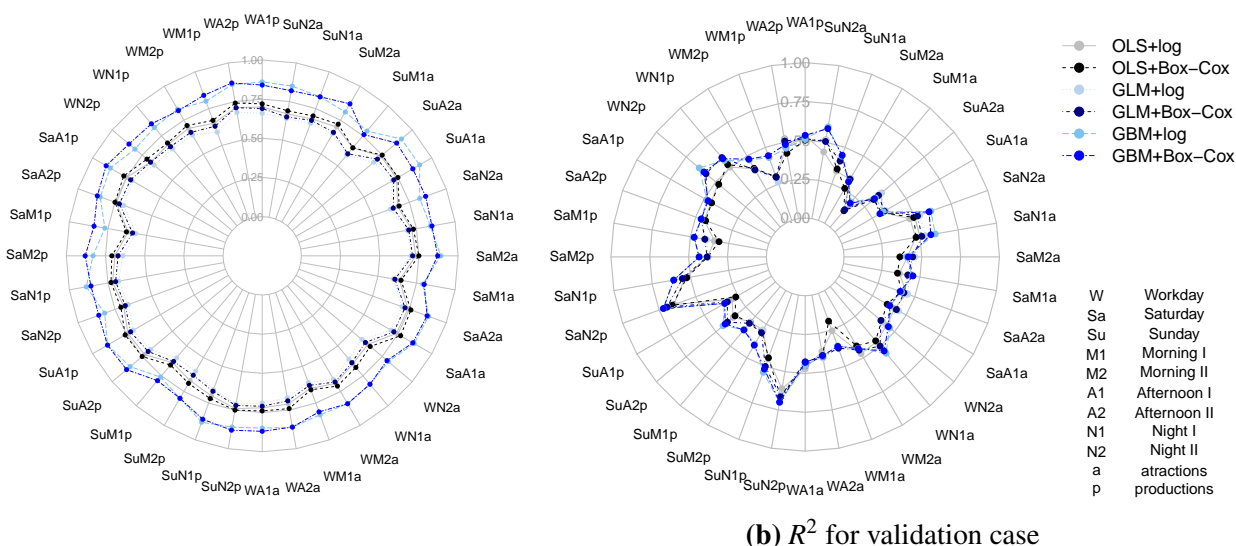


**FIGURE 3:** Spatial rentals distribution in cities of the study

1 realized to find the better number of trees or iterations with an input of 5 folds, a shrinkage factor  
 2 of 0.0001, and an interaction depth of 6. The presence of heteroscedasticity and nonnormality  
 3 in stepwise OLS and GLM led to the selection of logarithmic and Box-Cox transformations. Al-  
 4 though these properties were not identified using GBM, the transformations were also carried out  
 5 for matters of completeness. Outliers analysis was performed that indicated that zones with zero  
 6 arrivals and departures should be removed.

7 In total 324 models were built, considering arrivals and departures for three cases (workday,

1 Saturday, Sunday), 6 time intervals (morning, afternoon and night at peak and off-peak periods),  
 2 3 regression modeling techniques (stepwise OLS, GLM, GBM) and 3 transformation techniques  
 3 (no transformation, logarithmic and Box-Cox).  
 4 Regarding the parsimony the models, stepwise OLS selected the fewest number of variables  
 5 with an average in all temporal scales of 15.55 variables, followed by GLM with 26.13 and finally,  
 6 GBM with 39.90. According to the fitting results in the training cities (Figure 4a),  $R^2_{adj}$  in the  
 7 three regression methods trend together over different time periods. This indicates a rather indif-  
 8 ference to time goodness-of-fit. Between the regression techniques, GBM usually presented higher  
 9  $R^2_{adj}$  values. According to the validation performed with the city of Kassel, (Figure 4b) shows a  
 10 slight difference between the  $R^2$  values of different regression techniques, but there was a signif-  
 11 icant difference according to the time. Afternoons and nights showed the highest performances,  
 12 especially during the weekends. Finally, in all cases, a logarithmic and a Box-Cox transformation  
 13 illustrated a better goodness-of-fit with the logarithmic transformation to be slightly better.



(a)  $R^2$  adjusted from fit models

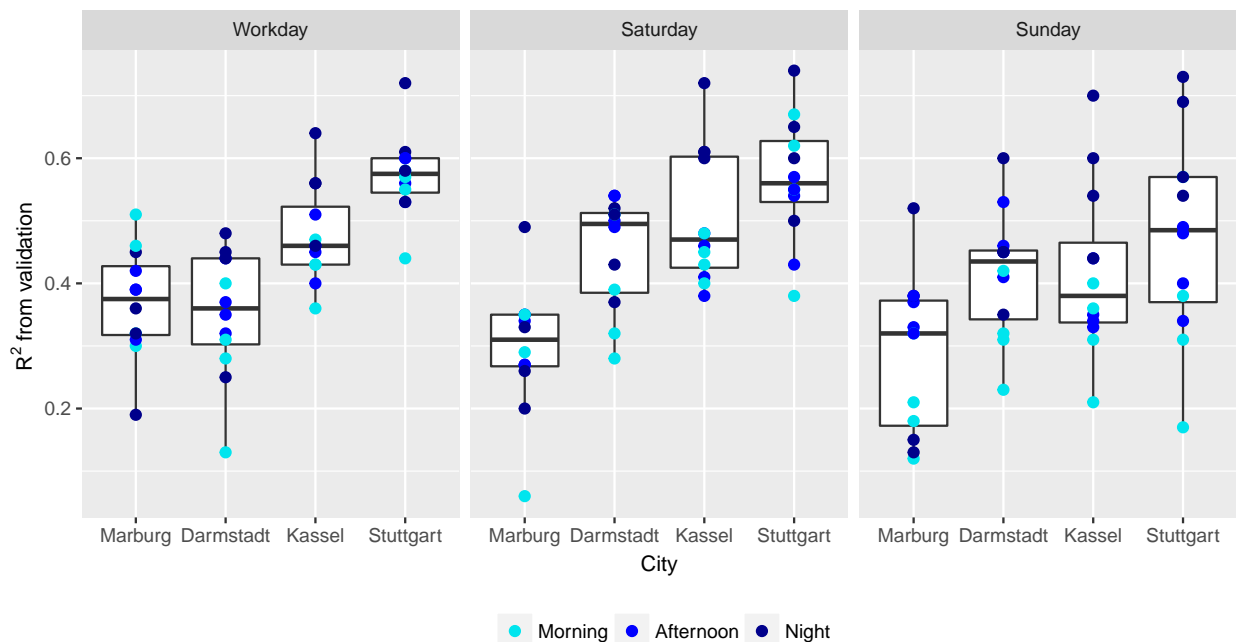
(b)  $R^2$  for validation case

**FIGURE 4:** Comparison of the  $R^2$  adjusted from the fit values of different models

14 Based on the above results, we also performed cross-validation, by dividing ridership data  
 15 in a training set of 5 cities' zones and a test set of one city's zone, for the case of GBM with a  
 16 logarithmic transformation. Hamburg and Frankfurt were excluded from this analysis since they  
 17 represent the majority of the zones of influence. The results presented in Figure 5 illustrate a rather  
 18 high performance in most cases, with workdays to have less variation of the validation scores  
 19 than on weekends. The city of Stuttgart, as a testing set, was the only case that showed a better  
 20 performance than the city of Kassel.

21 *Factors Affecting Bike Sharing*

22 Aiming at constructing an overview of the factors found to affect ridership, the occurrence of  
 23 parameters in all model structures were used. Figure 6 presents the most often selected variables  
 24 by the regression techniques per time interval. Darker blue indicates higher selection frequency.



**FIGURE 5:**  $R^2$  from validation by testing other cities (GBM with log transformation)

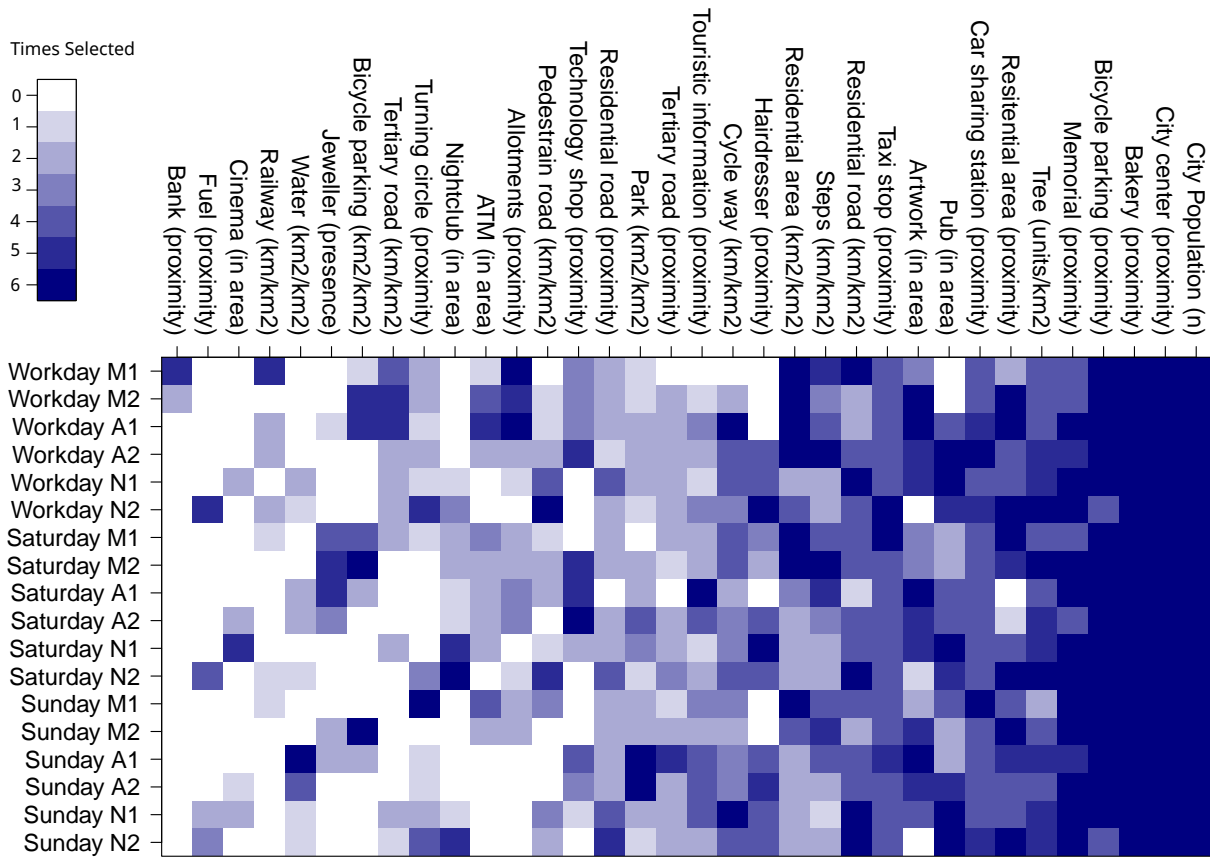
- 1 The most repetitive variables are the city population, the distance to the city center, bakeries,
- 2 bicycle parking, memorials, residential areas and car sharing stations.

### 3 Multiple international cities approach

4 The international application focused on the SBBS systems “Call a Bike” in Hamburg ([www.callabike-interaktiv.de/de/staedte/hamburg](http://www.callabike-interaktiv.de/de/staedte/hamburg)), “Divvy” in Chicago  
 5 ([www.divvybikes.com](http://www.divvybikes.com)) and “Bixi” in Montreal ([montreal.bixi.com](http://montreal.bixi.com)). The main objectives of  
 6 this application were the exploration of model transferability and the extraction of conclusions  
 7 for the application of the methods for different city structures on an international level. These  
 8 three cities were chosen since they share common characteristics as representative cities in their  
 9 countries with border limited by a body of water. However, in terms of population, Montreal and  
 10 Hamburg have relatively same inhabitants, but Chicago has around one million more inhabitants.  
 11 Thus, we are referring to mainly large cities, with a rather high population that could have different  
 12 travel characteristics and ridership patterns. As a consequence, the analysis was performed from  
 13 the beginning guiding a somewhat different modeling and validation approach. Bixi-Montreal  
 14 data (36) were collected from April 2014 until November 2017 with a data size of 734 MB and  
 15 545 stations. Divvy system works with 585 stations, where 1.75 GB data (37) was collected from  
 16 June 2013 until December 2017. Finally, 2.5 GB of Call a bike rentals were collected in Hamburg  
 17 from April 2014 until May 2017 in 207 stations (38).

### 19 Data analysis, and processing

20 The approach followed for the data analysis and processing was the same as the one for the national  
 21 case, with the exception that the rentals data were aggregated at an additional seasonal level. The  
 22 seasonality was added to analyze its effect on the resulting models. Chicago presented 9.93 rentals



**FIGURE 6:** Repetitive outcome variables influencing the arrivals using the modeling techniques after logarithmic and Box-Cox transformations

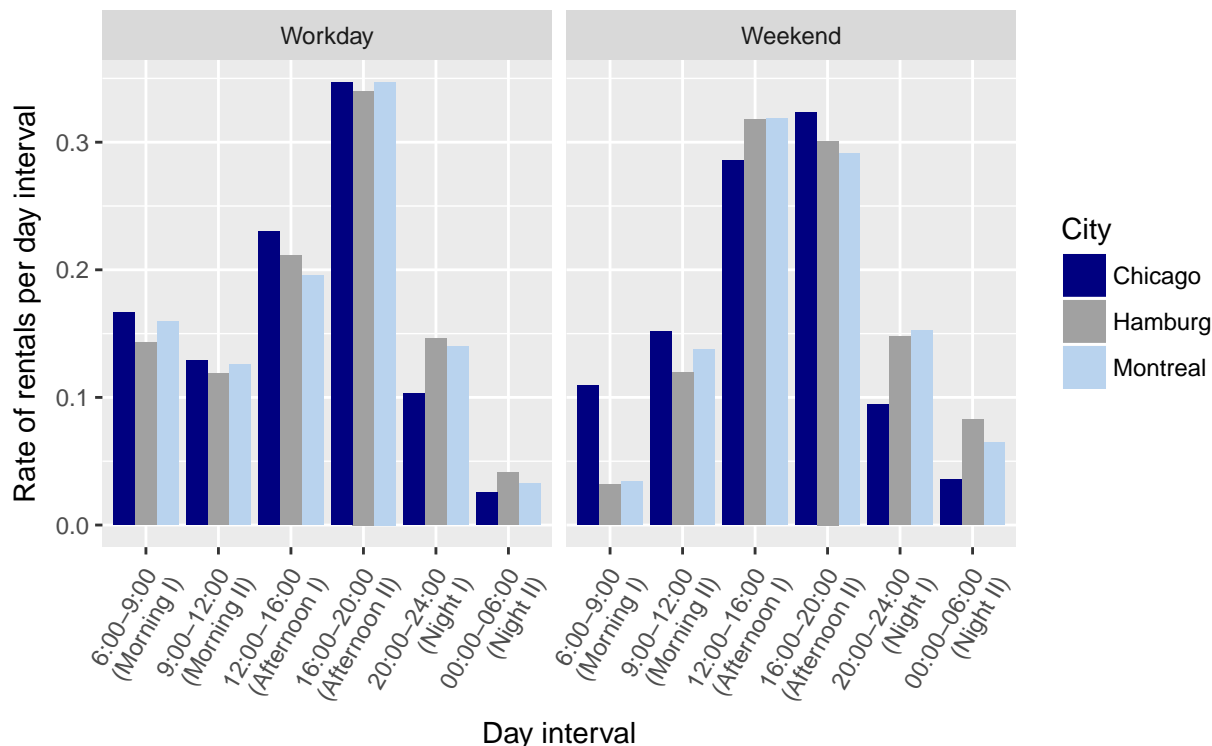
1 per day interval per station, while Hamburg 24.59 and Montreal 16.47. Figure 7 shows the distri-  
 2 bution of rentals per day interval in Chicago, Hamburg, and Montreal. These three cities present a  
 3 relatively similar distribution with an exception on the day interval "Morning I" on weekends. Fig-  
 4 ure 8 illustrates some examples of the spatial distribution of the rentals per time interval. It shows  
 5 higher ridership close to bodies of water. Concerning independent variables, **154** built environ-  
 6 ment variables were present in the three examined cities following the procedure as in the national  
 7 approach, where finally **113** non-collinear variables were considered for the modeling procedure.

8 *Model building, validation and variables selection*

9 Stepwise OLS with a logarithmic transformation was used for the model building. The choice  
 10 of using only one method was based on the computational time required to estimate all models  
 11 and because in the national application it was the most parsimonious method while preserving  
 12 relatively good fitting results. 72 models were built (one for each of four seasons, six day intervals,  
 13 and workdays, Saturdays, and Sundays). On validation, we considered an alternative approach of  
 14 the national case study by training 70% of the stations, and 30% for validation, without taking into  
 15 account the cities boundaries.

16 Model fitting and validation scores for Hamburg, Chicago an Montreal are shown in Fig-  
 17 ure 9.  $R^2_{adj}$  resulted of 0.68 as an average in the 72 models, and 0.63 as a  $R^2$  score in the testing





**FIGURE 7:** Distribution of rentals per day interval in Chicago, Hamburg and Montreal

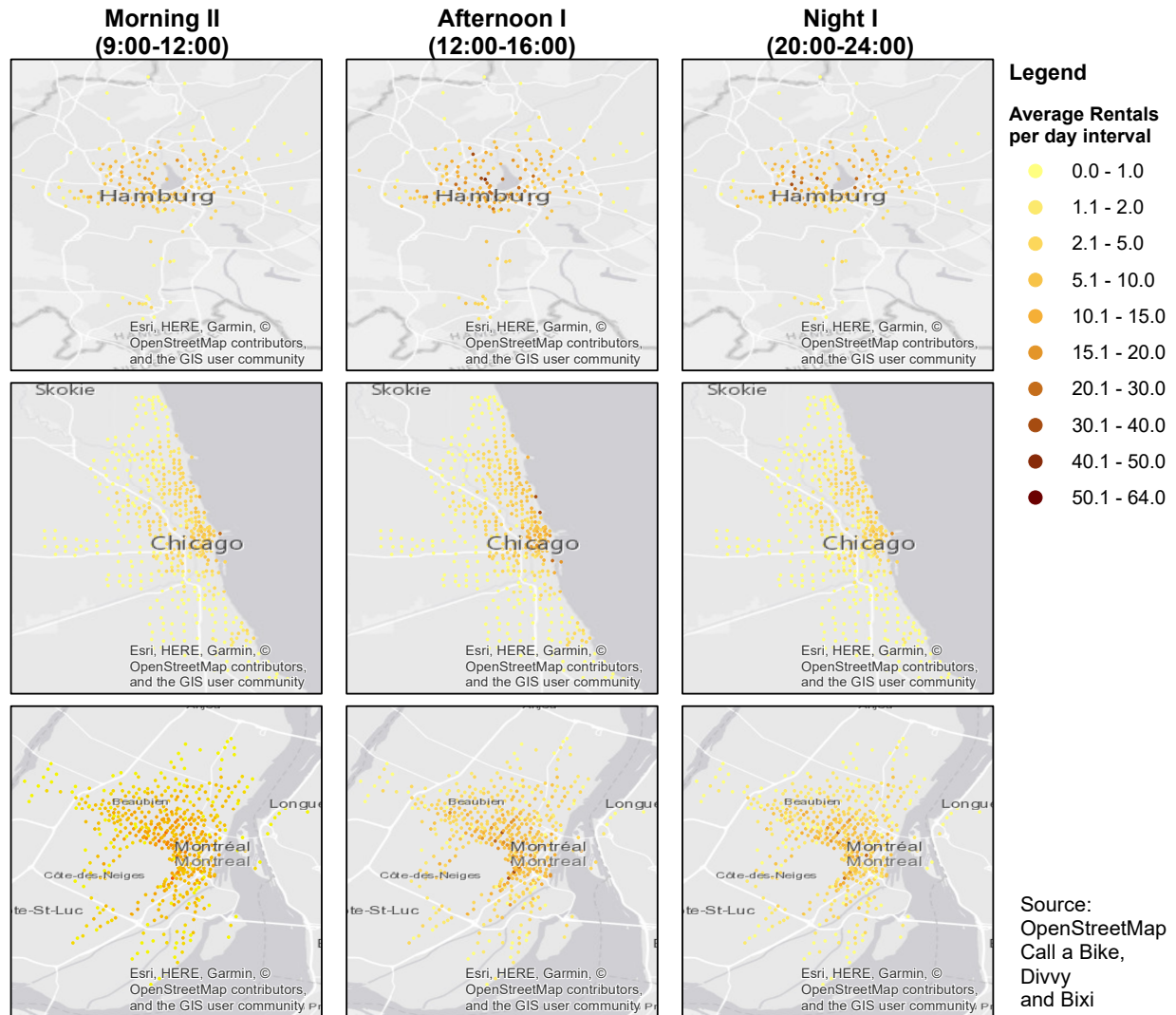
1 process. We run five times the model building as a cross-validation process, where the  $R^2$  from the  
 2 validation varied on average around the third decimal. The lower validation results were during  
 3 mornings on the weekends during all seasons, while higher values were associated with summer  
 4 and winter.

5 As an example of the resulting influencing factors in summer and winter on workdays,  
 6 Figure 10 presents the t-scores of the resulting models. On average 18.5 variables were selected  
 7 per model. The most common selected variables were the population and the proximity to colleges  
 8 and marina (area serving of leisure boats and yachts), bus stations, and restaurants and cafes.  
 9 Mainly in summer, several variables of high significance which have a negative influence during  
 10 morning and afternoon have a positive influence during the night. Land use influencing ridership  
 11 was mainly residential and parks.

## 12 DISCUSSION

13 A data-driven method using exclusively open-source data was applied in two case studies con-  
 14 sidering multiple cities in a national and an international level. From around 800 possible built  
 15 environment variables, the 144 most relevant and non-collinear variables were selected for the  
 16 model building for the national approach, while 113 were selected for the international approach.

17 Concerning model applicability, linear and non-linear modeling techniques were tested  
 18 in the national approach. GBM was the regression method that best fit the data, followed by  
 19 GLM. GLM and GBM required cross-validation tests to select the input arguments that helped  
 20 to build models. However, stepwise OLS was parsimonious with fewer input arguments, and  
 21 its results were easier to interpret. The three modeling techniques presented similar validation

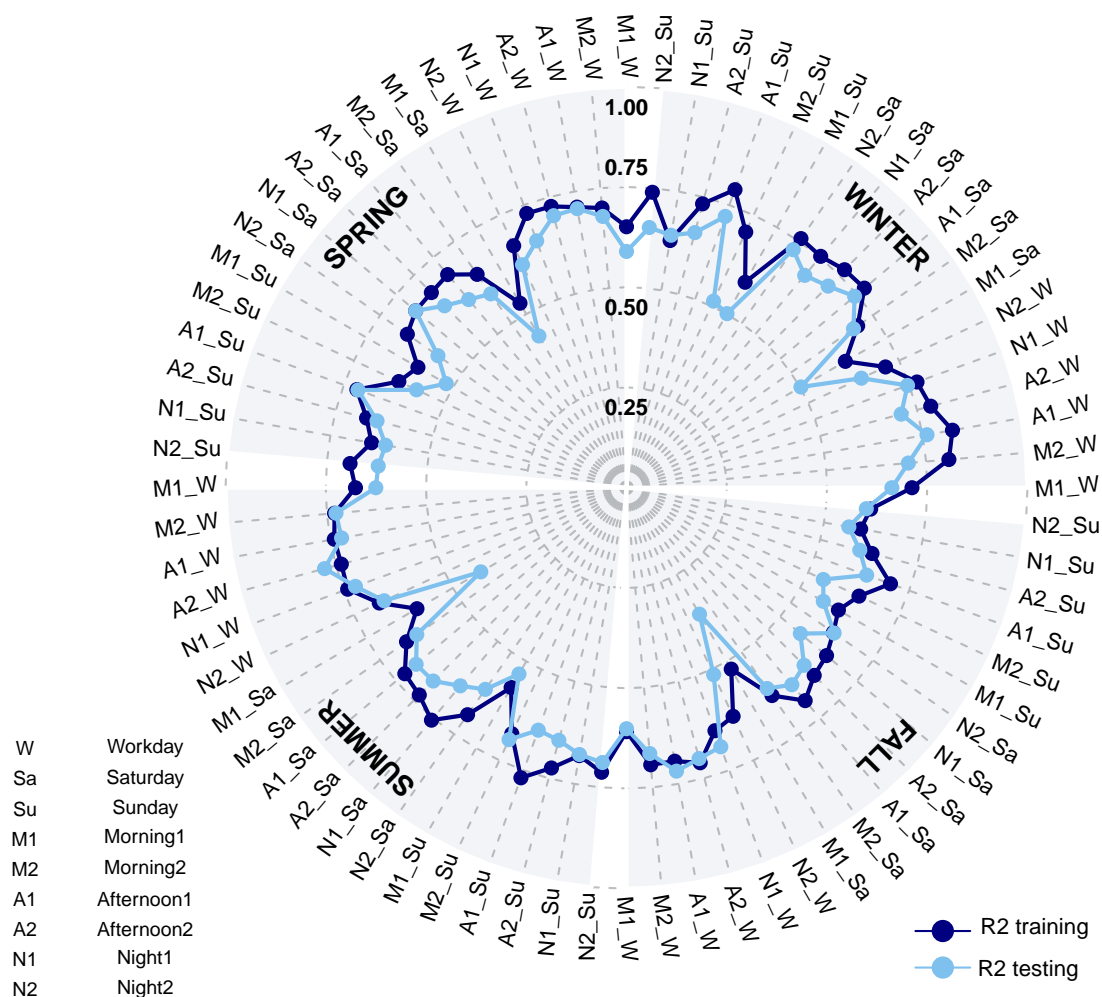


**FIGURE 8:** Spatial distribution of rentals per day interval in Chicago, Hamburg and Montreal

1 results. Logarithmic and Box-Cox transformations helped the models to predict better the arrivals  
 2 and departures and to select logical variables that would influence the shared bicycles ridership.  
 3 Generally, for the three regression methods, the logarithmic transformation performed a higher  $R^2$   
 4 in the validation phase.

5 The advantage of stepwise OLS and GLM was that a variable selection process was implicit  
 6 in the methods, but for GBM a variable selection process had to be developed to select those with  
 7 more influence from the ranking list. The most influencing variables in all of the built models and  
 8 through all time intervals were the population of the city and the distance from the city center (old  
 9 town) to the stations. The population of the city helped to weight the models to have a common  
 10 scale that was not biased if the city was large like Hamburg or small like Marburg. The distance to  
 11 the city center played a significant role for ridership as seen in Figure 2. The third most influencing  
 12 variable is the distance to bakeries. If a station is close to a bakery, this increases the probability of  
 13 higher ridership at that station.





**FIGURE 9:** Model fitting and validation scores for Hamburg, Chicago and Montreal (70% of the total stations for training)

1 In all developed models, several selected variables were logically correlated to bike shar-  
 2 ing ridership, were similar to literature review findings (Table 1), and they were coherent on the  
 3 authors’ expectations in influencing the arrivals and departures of bike sharing. For instance, the  
 4 most influencing variables are related to leisure activities, parks, green areas, and bodies of water  
 5 on the weekends, banks in the morning, gas stations, pubs, cinemas, clubs at night, shops on Satur-  
 6 days, and memorials outside of working hours. Just a few transport-related variables significantly  
 7 influenced the models. Distance to a car sharing station was significant for all time intervals as  
 8 well as bicycle parking. The only public transport variable displayed was railways during workday  
 9 mornings. It is worth mentioning that the tram and metro variables were not considered, because  
 10 they were not present in most of the studied cities.

11 According to the international approach, stepwise OLS with a logarithmic transformation  
 12 was chosen after the benefits identified in the national approach. On average, 18.5 variables were  
 13 selected per model. Urban structure was found to play an important role based on the distance  
 14 from all stations to the marina and colleges and also land use represented by residential area and

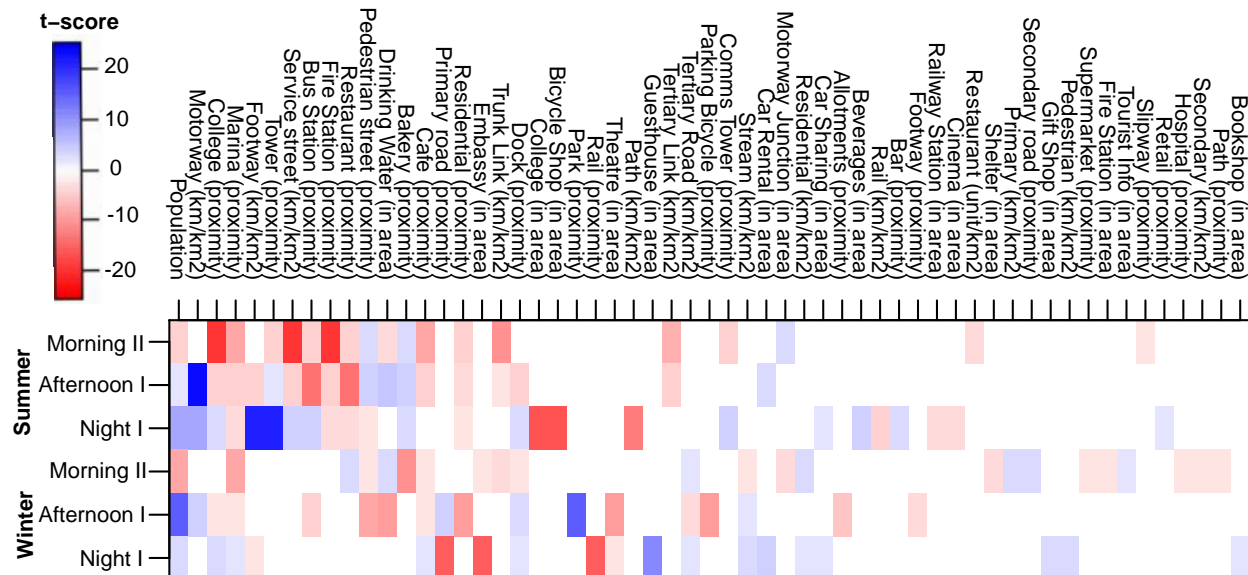


FIGURE 10: Example of t-scores for the international approach (Workdays)

1 parks. Summer and winter presented different factors, for example, bakeries, bus stations, and  
 2 restaurants are more significant in summer and not in winter. Logical variables were present as the  
 3 influence of bar and railway station in summer night or colleges influencing negatively at night.  
 4 An important observation is a correlation with car sharing stations during the night representing a  
 5 possible correlation between car and bike sharing.

6 Both approaches showed that the modeling validation results were correlated to the hourly  
 7 ridership distribution. Similar relative ridership hourly distribution was associated with higher  
 8 scores. For example, in the international approach in morning on weekends showed the most  
 9 different distribution between the cities (Figure 7), presenting the worst modeling performance.  
 10 On the other hand, in the national approach, the models that fit better the data were in the afternoon  
 11 and at night where the different cities showed a less variance in the bike sharing usage (Figure 3).  
 12 At these times, models presented better results from the validation and illustrated a more logical  
 13 selection of variables that influenced ridership. Also, on weekends the rate of ridership distribution  
 14 was more similar between the cities than on the workdays (Figure 3) showing higher validation  
 15 results. Finally, it was found that the modeling results did not depend on the size of the cities but  
 16 on the similarity of the distribution of the rate of rentals.

17 **CONCLUSION**

18 To the best of the authors’ knowledge, this is the first study that analyzes factors affecting bike  
 19 sharing systems ridership on a local level in multiple cities. The resulting influencing factors are  
 20 not only based on one city but beyond the geographic boundaries, which will help to use the  
 21 resulting models to forecast the bike sharing usage in a different city. A data-driven method was  
 22 developed to analyze the influence of the built environment in the rentals of station-based bike  
 23 sharing systems in multiple cities. An original approach was considered by modeling different  
 24 cities with different sizes in two case studies 1) on a national level and 2) on an international level.

25 GBM with a logarithmic transformation of the dependent variables were found to vali-

1 date slightly better the dataset. Stepwise OLS and a logarithmic transformation of the dependent  
2 variables was found to select fewer variables than other models without decreasing the validation  
3 results significantly. In Germany, the most influencing variables selected were the city popula-  
4 tion, the distance from the stations to the city center, bakeries, memorials, car sharing stations,  
5 among others. Logical relationships between the variables with the historical bike sharing rentals  
6 over time intervals were displayed, such as higher arrivals on nights close to pubs, cinemas, and  
7 nightclubs; or the presence of bodies of water, parks or green areas on Sundays. On an interna-  
8 tional level, the distance to the marina and colleges played an important role. Different influencing  
9 factors were present between different seasons.

10 The data-driven method will help as a decision-making support tool to implement or ex-  
11 pand bike sharing systems. Transport planners will have in a relatively brief time an idea of hot  
12 and cold spots of arrivals and departures of bike sharing systems to help them set coverage areas  
13 and place stations. This method can also show the validity and increase the reliability of measures,  
14 policies, and shared mobility projects. Further research and implementation can be developed as  
15 improvements and expansion of the case of study ( including more variables, such as topography  
16 or population density) or correlations between factors influencing car and bike sharing.

## 17 **ACKNOWLEDGMENT**

18 This study was partially supported by “Hans Boeckler Stiftung” and the MO3 project granted  
19 by the International Graduate School of Science and Engineering of TUM. The authors also are  
20 thankful for the bike sharing systems Divvy Bikes, BIXI-Montreal and "Call a bike" for sharing  
21 the rentals data through an open portal.

## 22 **AUTHOR CONTRIBUTION STATEMENT**

23 The authors confirm the contribution to the paper as follows: study conception and design: D. Du-  
24 ran, E. Chaniotakis, C. Antoniou; data collection: D. Duran; analysis and interpretation of results:  
25 D. Duran, E. Chaniotakis, C. Antoniou; draft manuscript preparation: D. Duran, E. Chaniotakis,  
26 C. Antoniou. All authors reviewed the results and approved the final version of the manuscript.

## 27 **REFERENCES**

- 28 [1] Büttner, J. and T. Petersen, *Optimising Bike Sharing in European Cities: A Handbook*. OBIS,  
29 2011.
- 30 [2] Böckmann, M., *The Shared Economy: It is time to start caring about sharing; value creating*  
31 *factors in the shared economy. University of Twente, Faculty of Management and Govern-*  
32 *ance*, 2013.
- 33 [3] Shaheen, S., E. Martin, A. Cohen, and R. Finson, *Public Bikesharing in North America:*  
34 *Early Operator and User Understanding, MTI Report 11-19*. Mineta Transportation Institute,  
35 2012.
- 36 [4] Meddin, R. and P. DeMaio, *The bike-sharing world map*. [http://www. metrobike. net](http://www.metrobike.net), Ac-  
37 cessed on 21.06.2018, 2018.
- 38 [5] Firnkorn, J. and S. Shaheen, Generic time- and method-interdependencies of empirical  
39 impact-measurements: A generalizable model of adaptation-processes of carsharing-users'  
40 mobility-behavior over time. *Journal of Cleaner Production*, Vol. 113, 2016, pp. 897 – 909.
- 41 [6] Shaheen, S., S. Guzman, and H. Zhang, *Bikesharing in Europe, the Americas, and Asia: past,*

- 1 present, and future. *Transportation Research Record: Journal of the Transportation Research*  
2 *Board*, , No. 2143, 2010, pp. 159–167.
- 3 [7] DeMaio, P., Bike-sharing: History, impacts, models of provision, and future. *Journal of pub-*  
4 *lic transportation*, Vol. 12, No. 4, 2009, p. 3.
- 5 [8] Hamann, T. K. and S. Guldenberg, Overshare and Collapse: How Sustainable are Profit-  
6 Oriented Company-to-Peer Bike-Sharing Systems?, 2017.
- 7 [9] Nikitas, A., *Bike-sharing fiascoes and how to avoid them – an expert’s guide*.  
8 [https://theconversation.com/bike-sharing-fiascoes-and-how-to-avoid-them-an-experts-](https://theconversation.com/bike-sharing-fiascoes-and-how-to-avoid-them-an-experts-guide-84926)  
9 [guide-84926](https://theconversation.com/bike-sharing-fiascoes-and-how-to-avoid-them-an-experts-guide-84926) , 2017, accessed on: 17.07.2018.
- 10 [10] Chardon, C. M. D., G. Caruso, and I. Thomas, Bicycle sharing system ‘success’ determinants.  
11 *Transportation Research Part A: Policy and Practice*, Vol. 100, 2017, pp. 202 – 214.
- 12 [11] Zhao, J., W. Deng, and Y. Song, Ridership and effectiveness of bikesharing: The effects of  
13 urban features and system characteristics on daily use and turnover rate of public bikes in  
14 China. *Transport Policy*, Vol. 35, 2014, pp. 253 – 264.
- 15 [12] Faghih-Imani, A. and N. Eluru, Incorporating the impact of spatio-temporal interactions on  
16 bicycle sharing system demand: A case study of New York CitiBike system. *Journal of Trans-*  
17 *port Geography*, Vol. 54, 2016, pp. 218 – 227.
- 18 [13] El-Assi, W., M. Mahmoud, and K. Habib, Effects of built environment and weather on bike  
19 sharing demand: a station level analysis of commercial bike sharing in Toronto. *Transporta-*  
20 *tion*, Vol. 44, No. 3, 2017, pp. 589–613.
- 21 [14] Tran, T. D., N. Ovtracht, and B. F. d’Arcier, Modeling Bike Sharing System using Built  
22 Environment Factors. *Procedia CIRP*, Vol. 30, 2015, pp. 293 – 298, 7th Industrial Product-  
23 Service Systems Conference - PSS, industry transformation for sustainability and business.
- 24 [15] Faghih-Imani, A., R. Hampshire, L. Marla, and N. Eluru, An empirical analysis of bike shar-  
25 ing usage and rebalancing: Evidence from Barcelona and Seville. *Transportation Research*  
26 *Part A: Policy and Practice*, Vol. 97, No. Supplement C, 2017, pp. 177 – 191.
- 27 [16] Noland, R. B., M. J. Smart, and Z. Guo, Bikeshare trip generation in New York City. *Trans-*  
28 *portation Research Part A: Policy and Practice*, Vol. 94, 2016, pp. 164 – 181.
- 29 [17] Wang, X., G. Lindsey, J. Schoner, and A. Harrison, Modeling bike share station activity:  
30 Effects of nearby businesses and jobs on trips to and from stations. *Journal of Urban Planning*  
31 *and Development*, Vol. 142, No. 1, 2015, p. 04015001.
- 32 [18] Faghih-Imani, A., N. Eluru, A. M. El-Geneidy, M. Rabbat, and U. Haq, How land-use and ur-  
33 ban form impact bicycle flows: evidence from the bicycle-sharing system (BIXI) in Montreal.  
34 *Journal of Transport Geography*, Vol. 41, 2014, pp. 306 – 314.
- 35 [19] Mattson, J. and R. Godavarthy, Bike share in Fargo, North Dakota: Keys to success and  
36 factors affecting ridership. *Sustainable Cities and Society*, Vol. 34, 2017, pp. 174 – 182.
- 37 [20] Fishman, E., S. Washington, and N. Haworth, Bike Share: A Synthesis of the Literature.  
38 *Transport Reviews*, Vol. 33, No. 2, 2013, pp. 148–165, cited By 98.
- 39 [21] Schmöller, S. and K. Bogenberger, Analyzing External Factors on the Spatial and Temporal  
40 Demand of Car Sharing Systems. *Procedia - Social and Behavioral Sciences*, Vol. 111, 2014,  
41 pp. 8 – 17.
- 42 [22] Bishara, A. and J. Hittner, Testing the significance of a correlation with nonnormal data:  
43 comparison of Pearson, Spearman, transformation, and resampling approaches. *Psychologi-*  
44 *cal methods*, Vol. 17, No. 3, 2012, p. 399.

- 1 [23] Box, G. E. P. and D. R. Cox, An Analysis of Transformations. *Journal of the Royal Statistical*  
2 *Society. Series B (Methodological)*, Vol. 26, No. 2, 1964, pp. 211–252.
- 3 [24] Chatterjee, S. and A. Hadi, *Regression analysis by example*. John Wiley and Sons, 2015.
- 4 [25] Tibshirani, R., Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Sta-*  
5 *tistical Society. Series B (Methodological)*, Vol. 58, No. 1, 1996, pp. 267–288.
- 6 [26] Friedman, J., Greedy function approximation: A gradient boosting machine. *Annals of Statis-*  
7 *tics*, Vol. 29, No. 5, 2001, pp. 1189–1232, cited By 2353.
- 8 [27] Lin, D., D. P. Foster, and L. H. Ungar, VIF Regression: A Fast Regression Algorithm for  
9 Large Data. *Journal of the American Statistical Association*, Vol. 106, No. 493, 2011, pp.  
10 232–247.
- 11 [28] Akaike, H., Maximum likelihood identification of Gaussian autoregressive moving average  
12 models. *Biometrika*, Vol. 60, No. 2, 1973, pp. 255–265.
- 13 [29] Schwarz, G., Estimating the dimension of a model. *The annals of statistics*, Vol. 6, No. 2,  
14 1978, pp. 461–464.
- 15 [30] Friedman, J., T. Hastie, and R. Tibshirani, *The elements of statistical learning*, Vol. 1.  
16 Springer series in statistics New York, 2001.
- 17 [31] Willing, C., K. Klemmer, T. Brandt, and D. Neumann, Moving in time and space – Location  
18 intelligence for carsharing decision support. *Decision Support Systems*, Vol. 99, 2017, pp. 75  
19 – 85, location Analytics and Decision Support.
- 20 [32] Deutsche Bahn AG, *Das smarte Leihfahrrad der Deutschen Bahn | Call a Bike*.  
21 <https://www.callabike-interaktiv.de/de>, 2017, accessed on: 31.10.2017.
- 22 [33] OpenStreetMap-contributors, *Planet dump retrieved from https://planet.osm.org* .  
23 <https://www.openstreetmap.org> , 2017, accessed on: 31.10.2017.
- 24 [34] Statistisches Bundesamt, *Bevölkerung und Erwerbstätigkeit. Bevölkerung mit Migrationsh-*  
25 *intergrund. Ergebnisse des Mikrozensus 2011. Fachserie 1, Reihe 2.2*, 2012, accessed on:  
26 31.10.2017.
- 27 [35] James, G., D. Witten, T. Hastie, and R. Tibshirani, *An introduction to statistical learning*,  
28 Vol. 112. Springer, 2013.
- 29 [36] BIXI-MONTREAL, *Open Data - BIXI Montréal*. [www.bixi.com/en/open-data](http://www.bixi.com/en/open-data), 2017, ac-  
30 cessed on 17.12.2017.
- 31 [37] Bikes, D., *Divvy Data*. <https://www.divvybikes.com/system-data>, 2018, accessed on:  
32 15.01.2018.
- 33 [38] Deutsche Bahn (DB), *Call A Bike - Open-Data-Portal – Deutsche Bahn Datenportal*.  
34 <http://data.deutschebahn.com/dataset/data-call-a-bike>, 2017, accessed on: 31.10.2017.