

Bus Travel Time Predictions Using Additive Models

Matthías Kormáksson
IBM Research – Brazil
matkorm@br.ibm.com

Luciano Barbosa
IBM Research – Brazil
lucianoa@br.ibm.com

Marcos R. Vieira
IBM Research – Brazil
mvieira@br.ibm.com

Bianca Zadrozny
IBM Research – Brazil
biancaz@br.ibm.com

ABSTRACT

Many factors can affect the predictability of public bus services such as traffic, weather and local events. Other aspects, such as day of week or hour of day, may influence bus travel times as well, either directly or in conjunction with other variables. However, the exact nature of such relationships between travel times and predictor variables is, in most situations, not known. In this paper we develop a framework that allows for flexible modeling of bus travel times through the use of Additive Models. In particular, we model travel times as a sum of linear as well as nonlinear terms that are modeled as smooth functions of predictor variables. The proposed class of models provides a principled statistical framework that is highly flexible in terms of model building. The experimental results demonstrate uniformly superior performance of our best model as compared to previous prediction methods when applied to a very large GPS data set obtained from buses operating in the city of Rio de Janeiro.

1. INTRODUCTION

In this paper we are concerned with the problem of predicting bus travel/arrival times using GPS data from public buses. The main challenge in performing this task arises from the fact that GPS data only provide snapshots of bus locations at predefined (or in some cases irregular) time stamps. The observed GPS coordinates are therefore necessarily irregular in space as signal transmissions are not controlled with respect to bus locations. The difficulty of the problem is further increased when difference between time stamps is large.

The raw GPS data permit us to study the relationship between bus movements in time and space. However, other factors such as day of week, hour of day, and current traffic conditions may also influence travel times in some systematic way. The exact nature of such relationships between travel times and predictor variables is usually not known. Therefore, these factors need to be incorporated into prediction algorithms either indirectly through binned analyses or through direct modeling.

We propose to model travel times using Additive Models [8, 26], which provide a principled statistical framework for arrival time predictions. In particular, we model cumulative travel time as a smooth function of route location and further allow this functional relationship to vary smoothly across (clock) time. We also construct features that may

seemlessly be incorporated into the Additive Model, either as direct main effects or interaction effects in conjunction with other variables.

Previous approaches have used a mixture of statistical and machine learning algorithms for predicting bus travel times. [19, 20, 11, 18] based predictions of future travel times on historical averages, either through binned analysis, e.g., with respect to hour of day, or by taking averages over similar past trips. [22, 21, 17, 6] used Kalman filter or time series models to predict future travel times under the assumption of a direct relationship with previous travel times. The above approaches lack the ability to incorporate other features into the prediction algorithms in a model based manner.

As an alternative regression models provide a simple and highly interpretable framework for modeling travel time as a function of several features. However, [18, 9] all demonstrated that the above models lack the flexibility to deal with nonlinear features so often present in these types of data. Artificial Neural Network (ANN) models and Support Vector Regression (SVR) models address this problem in a principled manner and have gained recent popularity in predicting bus arrival times because of their ability to deal with complex and nonlinear relationships between variables [4, 5, 28, 3]. However, these methods suffer from slow learning process [1, 7, 3] and are difficult to interpret and implement unlike regression models.

A recurring problem in the above approaches is that they assume knowledge of travel times between fixed locations in space, in particular bus stops. Often times these data are available (e.g., Automatic Passenger Count (APC) data [17, 4]) and provide information about exact arrival, departure, and dwelling times at specified bus stops. In the absence of such data, interpolation is performed to infer these times at the route's bus stops [11, 18]. This is reasonable when difference between time stamps is small, say 20 seconds, but can lead to larger errors when difference is larger, say few minutes. Another problem arises for methods that account for temporal effects (e.g., Kalman filters) due to discretization made in the time dimension. This is again reasonable in the presence of high volumes of data, but may be problematic if data is sparse with irregularities in the time dimension.

The main advantage of Additive Models in this context is their ease of interpretability and flexibility in modeling complex non-linear relationships. Factors that are known (or

suspected) to affect traffic may be included in the model as traditional linear features, smooth functional effects, or interactions thereof. Additive Models do not require any discretization or interpolated observations, but rather are capable of handling directly the raw observed data. The only interpolation that applies is made when inferring the departure time from origin. However, a critical feature of our proposed solution is the inclusion of a (corrective) random intercept in the model that attempts to correct for this interpolation step thus redefining time zero for each bus. Experimental results show that the random intercept model uniformly dominates all other methods in all prediction scenarios.

To the best of our knowledge our proposed solution is the first method that: (1) models bus travel times directly using raw irregular GPS data; (2) models spatial and temporal effects through smooth functions thus avoiding any discretization; and (3) allows for flexible incorporation of additional traffic related features in a model based manner. The last point is an important one as it implies that our proposed framework may be used as a development framework for building more accurate travel time models through the incorporation of additional (perhaps city dependent) features.

The remainder of the paper is organized as follows: Section 2 provides a summary of the motivating GPS data; background on additive models is provided in Section 3; the proposed framework for predicting bus travel times is detailed in Section 4; experimental evaluation is provided in Section 5; related work is described in Section 6; and Section 7 concludes the paper.

2. PRELIMINARIES

In this section we start by describing the motivating data. We then explain how the data are normalized and introduce mathematical notation.

2.1 Motivating Data

The motivating data consist of GPS measurements collected from public buses in the city of Rio de Janeiro, Brazil, during the time period from September 26, 2013 to January 9, 2014. The complete data set contains information about more than 400 bus routes and 9000 buses. Each GPS data point contains information about the position of the bus (longitude, latitude), date and time stamp, bus ID, and route ID. In total there are more than 100 million location entries for the time period of this study. The time between consecutive GPS measurements ranges from anywhere under a minute to over 10 minutes, with an average of ≈ 4 minutes. A sequence of GPS coordinates of a given bus is called a *space-time trajectory* and provides information about bus movement in space and time.

We also had access to GTFS (General Transit Feed Specification¹) data, which contain general information about the bus routes, such as bus stop locations. In general, each route consists of two trips, one going from origin to destination and the second representing the return trip. The GTFS data contain a complete definition of each such trip as a sequence of latitude/longitude points tracing the streets of the route

¹<http://developers.google.com/transit/gtfs>

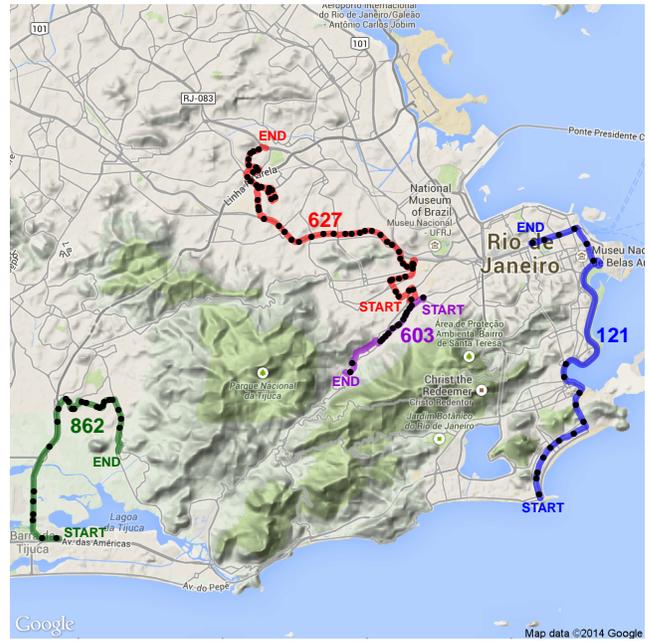


Figure 1: Map of the piecewise linear representations of the four routes analyzed in this paper. Route 603 (purple), 627 (red), 862 (green) and 121 (blue). Bus stops are marked by black points.

from origin to destination. In Figure 1 we display on map the bus stops and piecewise linear representations of routes 121 (running from Copacabana to City Center), 603 (running from Saens Pena to Usina), 627 (running from Saens Pena to Inhauma) and 862 (running from Jacarepagua to Barra da Tijuca).

Note that the observed data did not present itself without any challenges. For example, for each bus entry we only observe a general route ID for the round trip. The GPS data provide no further information about which direction the buses are travelling. However, by analyzing consecutive GPS measurements it is possible to infer the bus direction on the route. Other challenges involved erroneous or non-informative data entries. For example, for some buses the GPS measurements were observed far from the given routes and even at remote locations. We systematically removed all such non-informative entries in subsequent analyses.

2.2 Data Normalization

The GPS data in conjunction with the GTFS data provide us with the means to map GPS coordinates onto a 1-dimensional scale measuring distance from origin. For any given bus coordinate we project it onto the closest line segment of the corresponding route and then calculate its distance from origin along the piecewise segments.

By calculating differences between consecutive time stamps we may infer travel times of each bus between its observed locations. However, in order to analyze and compare travel times of buses, running at different hours, we need to normalize the time stamps onto a common cumulative time scale, i.e., we need to define a common time zero. This

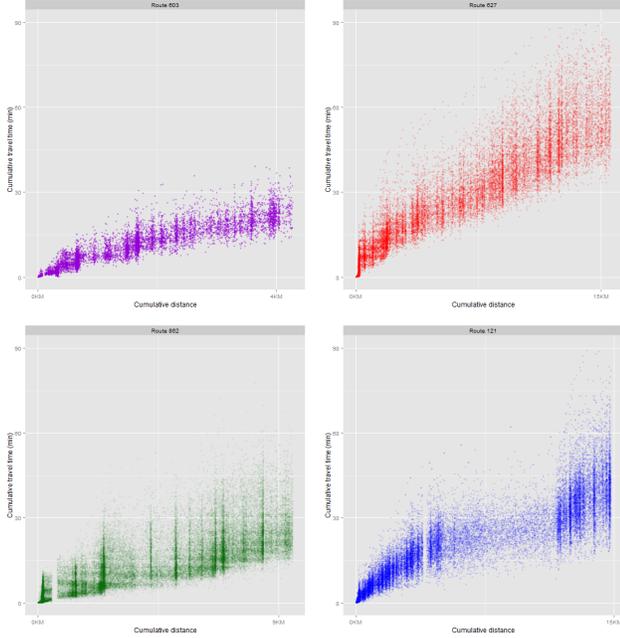


Figure 2: Cumulative space-time trajectories of the four bus routes analyzed in this paper. For each route a different transparency factor between 0 (completely transparent) and 1 (completely opaque) was chosen for plotting. This was done to normalize for varying numbers of data points per route and further make it possible to see where the bulk of the points lie. Note the distinct x-scales that reflect the different route lengths.

may be achieved by interpolating all the observed space-time trajectories at a common fixed point in space, e.g., origin, and defining that point as time zero. Space-time trajectories whose GPS coordinates have been mapped onto a cumulative distance scale and whose time stamps have been normalized to a common cumulative time scale are called *cumulative space-time trajectories*.

In Figure 2, we see the cumulative space-time trajectories of all buses (during the specified time period) running on the four routes analyzed in this paper. Note that the only interpolation made is at origin to define the common cumulative time scale. In all other aspects, the scatter plots represent raw measurements observed at irregular spatial locations.

2.3 Mathematical Notation

In general, we may normalize the time stamps at any arbitrary fixed point in space, in particular, at any of the route’s bus stops. Let $0 = p_0 < p_1 < \dots < p_K$ denote the distances of all bus stops of a given route from origin p_0 , where K denotes the number of on-route bus stops. *Cumulative space-time trajectories normalized at p_k* consist of cumulative distances $p_k \leq dist_{ijk} \leq p_K$, and corresponding cumulative travel times $T_{ijk} \geq 0$, $j = 1, \dots, m_{ik}$, where m_{ik} denotes the number of data points for bus i beyond p_k . The distances may either represent interpolated values at pre-specified fixed locations (e.g., subsequent bus stops) such as

in [11, 18], or raw GPS coordinates as in this paper. Both cumulative distances and cumulative times are defined from p_k onward such that $dist = 0$ and $T = 0$ at p_k . The cumulative time scale is inferred by interpolating two consecutive time stamps before and after p_k . We denote by $Traj(p_k)$ the set of thus normalized historical cumulative space-time trajectories.

3. THEORETICAL BACKGROUND

Additive models [8, 26] are linear models, which allow the linear predictor to not only depend on pure linear terms but also on a sum of unknown smooth functions of predictor variables. This class of models is particularly powerful when there is an evident smooth relationship between the response and predictor variables but exact parametric form can neither be theoretically nor intuitively inferred. However, we need to specify these functions in some meaningful way and determine the degree of smoothness. This section discusses both of these topics.

3.1 Penalized Spline Smoothing

Let us first consider one-dimensional functions through the simple scatterplot smoothing model

$$y_i = f(x_i) + \varepsilon_i, \quad (1)$$

$i = 1, \dots, n$. A common approach [26, 16], is to represent the function as $f(x) = \sum_{j=1}^q \beta_j \phi_j(x)$, where $\phi_j(x)$ are known basis functions and β_j are coefficients to be estimated. An intuitive example is the piecewise linear representation, involving basis functions $\phi_1(x) = 1$, $\phi_2(x) = x$, and $\phi_{j+2}(x) = (x - \tau_j)_+ \equiv \max(0, x - \tau_j)$, $j = 1, \dots, q - 2$, where τ_j are called knots that need to be chosen (e.g., equally spaced in x -domain). The exact choice of the number of knots and placement is not generally critical and is not the focus of this paper. In general the number should be chosen to be large enough to represent the underlying truth reasonably well, while at the same time maintaining computational efficiency. By letting $X = [1 \ x_i \ (x_i - \tau_1)_+ \ \dots \ (x_i - \tau_{q-2})_+]_{1 \leq i \leq n}$ the model function may now be written in matrix form as $f(x) = X\beta$. This representation is quite general and there exist several families of basis functions that fit into the above framework. For example, a simple cubic regression spline can be obtained by defining $\phi_{j+2}(x) = |x - \tau_j|^3$ instead of the truncated linear basis above.

The above model may be estimated by least squares or by maximizing the loglikelihood function, $\ell(\beta)$, under a normality assumption on ε . However, in order to control the smoothness of the fit we need to work with the so called penalized loglikelihood

$$\ell_P = \ell(\beta) - \lambda \beta' D \beta, \quad (2)$$

where D is most often specified as $\text{diag}(0, 0, 1, \dots, 1)$ and λ is a smoothness parameter. If λ is chosen too large the resulting fit becomes closer and closer to a linear fit in the above case. On the other hand, choosing λ too small may lead to overfitting. In general the smoothness parameter may be estimated, for example, using Generalized Cross Validation (GCV).

3.2 Additive Models

The additive models that we consider in this paper have the form:

$$y_i = X_{0i}\beta_0 + f_1(x_{1i}) + f_2(x_{2i}) + f_3(x_{1i}, x_{2i}) + \varepsilon_i, \quad (3)$$

where y_i and ε_i are the response and error term respectively, $X_{0i}\beta_0$ represents purely linear terms in the model, and f_1 , f_2 , and f_3 represent smooth functions of the predictors x_1 and x_2 . We represent the one-dimensional functions as in the previous subsection: $f_1(x_1) = \sum_{j=1}^{q_1} \beta_{1j}\phi_j(x_1)$ and $f_2(x_2) = \sum_{k=1}^{q_2} \beta_{2k}\psi_k(x_2)$, where $\phi_j(x_1)$ and $\psi_k(x_2)$ are known (possibly distinct) basis functions. In this paper a tensor product basis [26, 16] is used to represent the bivariate term:

$$f_3(x_1, x_2) = \sum_{j=1}^{q_1} \sum_{k=1}^{q_2} \beta_{3jk}\phi_j(x_1)\psi_k(x_2). \quad (4)$$

Through a similar argument as in the previous subsection, each of the above functions f_i , $i = 1, 2, 3$, may be represented by $X_i\beta_i$, where the X_i matrices are appropriately specified in terms of the basis functions $\phi_j(\cdot)$, and $\psi_k(\cdot)$. The model terms may then be stacked in the traditional way: $X = [X_0 \ X_1 \ X_2 \ X_3]$, and $\beta = (\beta'_0, \beta'_1, \beta'_2, \beta'_3)'$, to obtain the linear model:

$$Y = X\beta + \varepsilon. \quad (5)$$

This demonstrates that with the appropriate specification of the smooth functions an additive model is simply a linear model whose smoothness of fit may be controlled by placing a penalty on the β terms. We may separately control the smoothness of each function by introducing function specific smoothness parameters λ_i . The penalized loglikelihood from (2) then generalizes naturally to:

$$\ell_P = \ell(\beta) - \sum_{i=1}^3 \lambda_i \beta'_i D_i \beta_i, \quad (6)$$

where D_i are specified similarly.

Estimation of the additive model may be performed by maximizing the penalized likelihood in (6) and estimating the smoothness parameters through GCV. Once the model has been estimated using training data, one can predict a new response in the usual manner.

3.3 Additive Model with Random Intercept

In this paper we also consider an additive model with a random intercept

$$y_i = b_{0i} + X_{0i}\beta_0 + f_1(x_{1i}) + f_2(x_{2i}) + f_3(x_{1i}, x_{2i}) + \varepsilon_i, \quad (7)$$

where $b_{0i} \sim N(0, \sigma_b^2)$ and $\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$. Note that the above model is not overparametrized as the b_{0i} are treated as random and not fixed. This model falls into the general class of Additive Mixed Models [26] (due to the mixed combination of random and fixed model terms) and we note that by specifying the smooth functions as before it may be represented in the matrix form:

$$Y = X\beta + Zb_0 + \varepsilon, \quad (8)$$

where Z is a single column matrix of ones.

The estimation of the above model is not straight forward and since space is limited we point to [26] for full theoretical

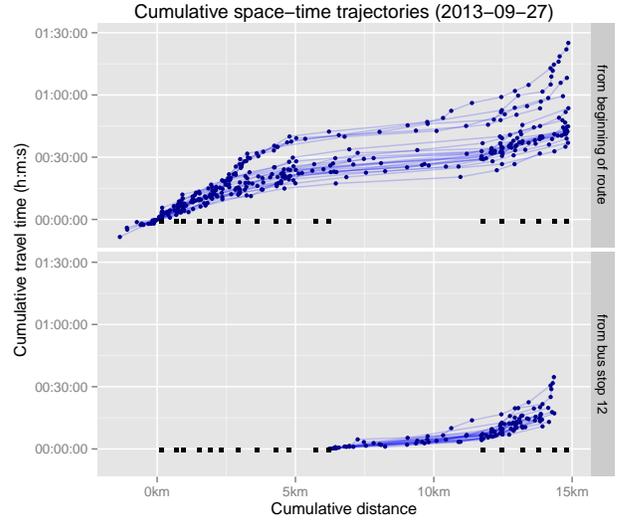


Figure 3: Cumulative space-time trajectories for route 121 from September 27, 2013. Upper panel demonstrates the trajectories from beginning of route and the lower panel from bus stop p_{12} onward. Bus stops are marked by black squares.

coverage. However, we note that through a correct likelihood specification an iterative maximization algorithm may be applied to obtain estimates of the parameters β , σ_b^2 , and σ_ε^2 . Given these estimates the prediction formula for the random effect vector is (see e.g. [16])

$$\hat{b}_0 = \sigma_u^2 Z' V^{-1} (y - X\beta), \quad (9)$$

where $V = \sigma_u^2 Z Z' + \sigma_\varepsilon^2 I$.

3.4 Computational Aspects

Additive models, such as (3) using penalized splines and tensor product smooths are implemented in a highly optimized R-package, *mgcv* which allows estimation of the model, [23, 24, 25, 27]. Additive mixed models, such as (8) are more computationally expensive than regular Additive Models, in particular when the number of random effects becomes large. However, the *mgcv*-package also has an optimized routine for estimation through a call to the *lme* function of the highly developed *nlme* R-package [12] that was specifically designed to estimate linear mixed models efficiently.

4. PROPOSED SOLUTION

In this section we present additive models for analyzing historical cumulative space-time trajectories such as those observed in Figure 3. In the upper and lower panel, respectively, we see examples of historical trajectories, $Traj(p_k)$, that have been normalized at p_k for $k = 0$ (origin) and $k = 12$ (bus stop p_{12}). We note that the cumulative travel time variance beyond bus stop p_{12} is reduced dramatically when normalized at p_{12} as compared to at p_0 . Therefore, we propose to train additive models on each of the historical trajectories, $Traj(p_k)$, for bus stops $k = 0, \dots, K - 1$, where $K - 1$ corresponds to the second to last bus stop on route. The objective is then to base future travel time pre-

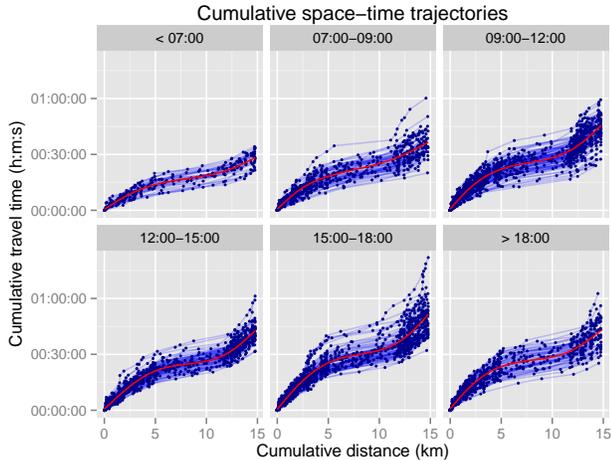


Figure 4: Cumulative space-time trajectories of route 121, stratified by hour, during the time period September 26 - October 10, 2013. The dots represent raw measurements and the blue interpolated curves represent each bus trajectory for illustration purposes. A smooth mean curve for each time category is depicted in red.

dictions of a bus close to bus stop p_k on the corresponding additive model trained on $Traj(p_k)$.

In order to make our presentation more coherent, we model and analyze the bus trajectories of bus route 121 (Copacabana-Center) for the first two weeks of our observed time period. We analyze trajectories starting from origin, $Traj(p_0)$ and for ease of notation we omit the k -subindex of Subsection 2.3. However, we note that all discussions generalize to trajectories starting from any given bus stop along the route, $Traj(p_k)$. For the first two weeks of our study we observed $n = 385$ trajectories for route 121 with on average $m_i = 13$ measurements per bus ride. Through statistical reasoning, we construct three models whose performances are compared to previous approaches (Section 5). All numerical summaries in this section apply to this data set.

4.1 Basic Additive Model for Travel Times

In Figure 4, we see the cumulative space-time trajectories of all bus rides of route 121 during the specified time period. The trajectories are stratified by hour and a smooth mean curve is fitted through each scatterplot to illustrate travel time trends. We note that morning travel time duration peaks between 9am and noon (morning rush hour). Then a slight reduction in travel times is observed between noon and 3pm, followed by an afternoon rush hour. We also note that there is not only a difference in total travel times across hours, but also in the shapes of the mean curves. This figure inspires the following model of bus travel time, T_{ij} , as function of distance from origin, $dist_{ij}$, and time of departure, $time_i$:

Model 1: Basic Additive Model (BAM)

$$T_{ij} = \beta_0 + f_1(dist_{ij}) + f_2(time_i) + f_3(dist_{ij}, time_i) + \varepsilon_{ij},$$

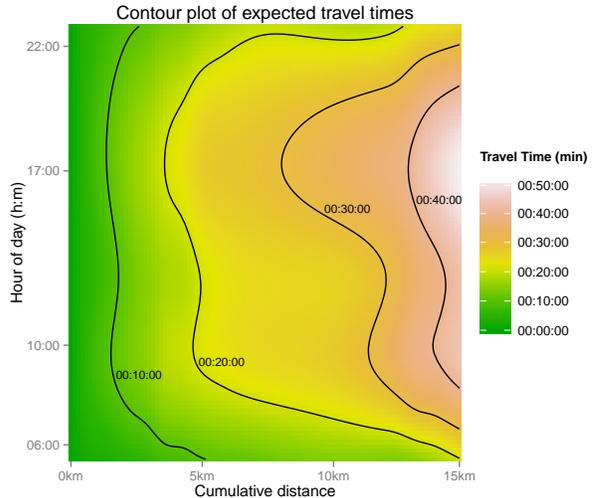


Figure 5: A contour plot of estimated travel times as a smooth function of time of day and cumulative distance from origin.

$i = 1, \dots, n$, and $j = 1, \dots, m_i$, where β_0 , and ε_{ij} represent an overall model mean and error term, respectively. In what follows we assume the random error terms are mean zero and normally distributed. The terms f_1, f_2, f_3 denote unknown smooth functions designed to capture functional relationships such as those observed in Figure 4. The f_1, f_2 -terms can be thought of as smooth main effects of $dist$ and $time$ on T , respectively, whereas the f_3 -term represents an interaction effect of the two variables. The interaction allows the functional relationship between T and $dist$ to change with $time$, as observed in Figure 4.

The functions f_1, f_2 and f_3 were represented by cubic regression splines and tensor product smooths (see Section 3). We placed one knot at each bus stop between origin and destination to capture smooth transitions from one station to the next. We placed 5 equally spaced knots in the time space, which was large enough to capture the two rush hour trends in the morning and afternoon, respectively. Larger number of time-knots did not seem to affect the fit of the model.

We estimated the Basic Additive Model using the *mgcv* R package. The numbers showed that each of the functional effects f_1, f_2 , and f_3 was deemed statistically significant by the F-test (p-values $< 10^{-16}$) and the overall adjusted R^2 of the model was 0.903. To illustrate the smooth relationship between the two variables and travel time, Figure 5 shows a contour plot of estimated travel time with cumulative distance from origin on x-axis and time of day on y-axis. We can see that at 10am it takes the bus on average approximately 30min to travel around 12km, while at 5pm it only travels around 8km in half an hour. We also see that the two rush hour peaks, at approximately 10am and 5pm, are more amplified at 12km than at 2km, as exemplified by the 30min and 10min contour lines, respectively. These observations demonstrate the importance of including the interaction term f_3 in the model.

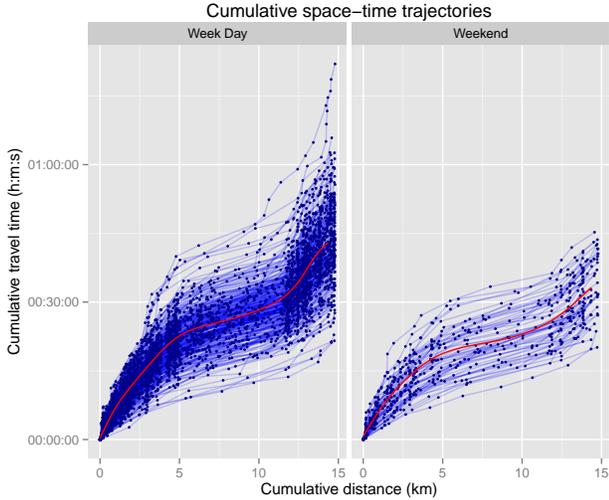


Figure 6: Cumulative space-time trajectories of route 121 during week days (left) and weekends (right). The dots represent raw measurements and the blue interpolated curves represent each bus trajectory for illustration purposes. Smooth mean curves are depicted in red.

4.2 Extended Additive Model with Additional Features

It is well known that traffic patterns in cities are different on a weekday as compared to the weekend (see Figure 6). This phenomenon can easily and quite flexibly be incorporated into our basic additive model above. Let $weekend_i$ denote an indicator variable that determines whether bus ride i occurred on a weekday or on the weekend. Then by adding the linear term $\beta_1 \cdot weekend_i$ into the model we account for differences in overall mean travel times between weekdays and weekends. However, there is an evident interaction of the weekend factor with distance from origin as can be seen in Figure 6, where the mean difference between weekday and weekend travel times increases as a function of distance. Therefore, in addition to the main effect $\beta_1 \cdot weekend_i$, we propose replacing the functional term $f_1(dist_{ij})$ in Model 1 by the interaction term $f_1(dist_{ij}, weekend_i)$. This term in fact generates two different smooths, one for weekday and the other for weekend trajectories.

Another feature that intuitively seems likely to correlate well with travel time of a given bus is the travel time of the last bus in front of it. We therefore define the feature T_{ij}^{last} to be the cumulative travel time at $dist_{ij}$ of the last bus that passed before time of departure of bus i . It is important to point out here that it is unlikely that the last bus will transmit a GPS signal at the same locations $dist_{ij}$ as bus i . Therefore, interpolation of the cumulative space-time trajectory of the last bus is performed at $dist_{ij}$ to construct the feature T_{ij}^{last} . We observed that T_{ij}^{last} had a strong linear relationship with the observed travel times, T_{ij} , with a sample correlation of $r \approx 0.92$.

The following model extends Model 1 to include the features

discussed above:

Model 2: Extended Additive Model (EAM)

$$T_{ij} = \beta_0 + \beta_1 \cdot weekend_i + f_1(dist_{ij}, weekend_i) + \beta_2 \cdot T_{ij}^{last} + f_2(time_i) + f_3(dist_{ij}, time_i) + \varepsilon_{ij}.$$

We fitted the above model to our data set and observed that all effects, including $weekend_i$, the interaction term $f_1(dist_{ij}, weekend_i)$, and the linear predictor T_{ij}^{last} were highly significant (p-values $< 10^{-16}$). The adjusted R^2 increased from the previous model to 0.919.

4.3 Additive Mixed Model

We recall that in Section 2.3 we normalized all the space-time trajectories to a common cumulative time scale. Since the actual times of departure are not known, an approximation was made by taking two consecutive time stamps before and after origin and defining $T = 0$ as the interpolated time-stamp at origin. However, note that this introduces an error in the form of a vertical trajectory shift, due to incorrect specification of time of departure. One bus driver may, for example, take a short break at origin, while another may depart soon after arriving. This error is unpredictable in our context and can amplify when time resolution is poor such as in our data, where GPS coordinates are only transmitted on average every 4 minutes.

In order to correct for the misspecification of time of departure, we propose an additive mixed model (see section 3.3) that includes a (corrective) random intercept term b_{0i} for each and every bus ride $i = 1, \dots, n$:

Model 3: Additive Mixed Model (AMM)

$$T_{ij} = \beta_0 + b_{0i} + \beta_1 \cdot weekend_i + f_1(dist_{ij}, weekend_i) + \beta_2 \cdot T_{ij}^{last} + f_2(time_i) + f_3(dist_{ij}, time_i) + \varepsilon_{ij},$$

where $b_{0i} \sim N(0, \sigma_b^2)$.

We fitted the above model and found that the random intercept term b_{0i} was indeed highly significant (p-value $< 10^{-16}$) and the adjusted R^2 value increased significantly to 0.968. The standard deviation σ_b was estimated to be 3 minutes, which indicates that the estimated (interpolated) time of departure indeed requires adjustment.

5. EXPERIMENTS

5.1 Experimental Data

We performed a prediction analysis on four bus routes in the city of Rio de Janeiro: 603, 627, 862, and 121 (see Figure 1). These routes are located in distinct regions of the city and further have different lengths, number of bus stops, and frequency of bus rides, see Table 1. Further these routes demonstrate distinct traffic patterns as can be seen in Figure 2. We have made all the data sets available in [13], which we believe this will stimulate further research in the area.

We note that since GPS locations are transmitted on average every 4 minutes the density of points in the spatial dimension of Figure 2 provides some insight into traffic behavior at different route segments. We see, for example, a lighter

Table 1: Route data summary

Route	# trajectories	# stops	Length (in km)
603	1,276	15	4
627	1,325	54	15
862	7,882	24	10
121	2,515	18	15

blue section in the middle of route 121 (due to fewer observations), which represents an inner city thruway less prone to traffic congestions. On the other hand, darker sections represent locations where traffic may experience regular stops or delays, such as traffic signals or frequently congested road segments. Locations with no data points represent either tunnels or regions with poor reception. We note that although both routes 121 and 627 have the same length their cumulative space-time trajectories are quite different. These four distinct routes represent a wide range of prediction scenarios we want to cover in our experiments.

5.2 Experimental Setup

The total number of trajectories in each route is presented in Table 1. We randomly selected 14 days after November 1st 2013 as our test data. This guaranteed at least 30 days of historical data for each test date. For each bus running on any of these 14 days we performed travel time predictions using three sets of historical data involving all bus rides in the last 10, 20, and 30 days, respectively. This was done to get a sense of whether the size of the historical data set has an influence on the accuracy of the tested models.

Travel time predictions were made for each bus in test set from every bus stop until end of route to reflect the real world problem of predicting bus arrivals from any on-route location onward. More precisely, for each bus stop p_k we recorded for bus i in test set the first observed bus entry $dist_{i1k}$ after p_k . We then made travel time predictions at all remaining (observed) points $dist_{ijk}$, $j = 2, \dots, m_{ik}$. Since the data at $dist_{ijk}$ represent the raw data whose cumulative travel times T_{ijk} (from bus stop p_k) are known we could thus calculate and analyze prediction errors. In order to get a sense of how error changes as a function of distance from the bus stop beyond which predictions were made we recorded the prediction distances $|dist_{ijk} - p_k|$.

5.3 Evaluation Measures

To evaluate overall performance of each method for a given bus route we calculated the *mean absolute relative error*, defined as $(1/N) \sum_{ij} |T_{ij} - \hat{T}_{ij}|/T_{ij}$, where N denotes total number of predictions made. Since the distributions of the relative errors was right skewed in all cases a median could have been used instead of mean. However, as the mean is less robust to outliers it may also provide insight about worst case errors. Overall conclusions were not affected by replacing the mean with median. We performed a non-parametric paired Wilcoxon test to compare the overall performances between methods.

Since error was greater at later parts of route, we also analyzed the distributions of absolute errors stratified by predic-

tion distances, $|dist_{ijk} - p_k|$. The distance space was binned into one kilometer bins $[0,1)$, $[1,2)$, $[2,3)$, ... etc. Visual comparison of distributions was performed using boxplots and a 95th percentile curve; see Figure 7. Since absolute errors were right skewed for each method we performed a two-sided non-parametric paired Wilcoxon test to compare methods within each distance bin.

To account for multiple testing, p-values were recorded for each comparison and then adjusted using the Benjamini-Hochberg method [2]. Statistical significance was determined if adjusted p-values were < 0.05 .

5.4 Implemented Methods

Additive Models: No model selection or parameter tuning was performed during the training. Instead for each and every training set we estimated the exact same three models as defined in Section 4. Once estimation had been performed the estimated model parameters, $\hat{\beta}$, along with a complete set of test features was plugged into the Additive Model formula (5) to obtain travel time predictions at subsequent route locations. For AMM, in order to estimate the random effect b_{0i} of (7) for a new trajectory i in the test set at least one observed travel time is needed. Since predictions are always made given the current location of the bus, the first observation, T_{i1} , may be used for that purpose. By plugging this value in for y in the formula (9) we obtain an estimate of b_{0i} . Then the formula (8) may be used in conjunction with a complete set of test features to obtain travel time predictions at subsequent route locations. A minor implementation detail we want to point out involves predictions beyond bus stops very close to route destination. In this case the training data $Traj(p_k)$ can become scarce and full spline function representation as defined in our three proposed models in Section 4 may lead to overfitting. Therefore, in these cases, we replaced the smooth functions with the more simple linear model terms: $\alpha_1 dist_{ij} + \alpha_2 dist_{ij} \cdot weekend_i + \alpha_3 time_{ij} + \alpha_4 dist_{ij} \cdot time_{ij}$.

Support Vector Machine (SVM): Bin et al. [3] used support vector machine regression to predict the arrival time of the next bus. They divided the bus trajectories in segments and then used as features the travel time of current bus at previous segment and the latest travel time of a previous bus in the next segment to predict the travel time for the next segment. Since we are not only interested in predicting the travel time of the next segment but all subsequent segments until end of route, we added to the training data the latest travel times at all subsequent segments. Similar to [3], we used a linear kernel and the implementation was performed using the R package “e1071”.

Kernel Regression: Sinn et al. [18] proposed an instance-based method that uses weighted averages of historical trajectories to make predictions. Trajectories with similar behaviour up to the current bus location are given more weight. Weights are defined by a gaussian kernel: $\exp(-\|x - y\|^2/b)$, where x and y are cumulative space-time trajectories, and b is the bandwidth of the kernel². For further details, we refer the reader to [18].

²Similar to [18], we set $b = 1$ in our experiments

Table 2: Mean Absolute Relative Error

Route	# days	Method				
		BAM	EAM	AMM	Kernel	SVM
603	10	19.9%	19.7%	18.4%	21.3%	64.4%
	20	20.1%	19.8%	18.5%	21.3%	64.7%
	30	19.8%	19.6%	18.3%	21.3%	64.8%
627	10	16.3%	14.7%	13.8%	18.1%	28.8%
	20	15.2%	14.2%	13.4%	17.3%	30.0%
	30	15.1%	14.0%	13.2%	17.1%	29.4%
862	10	22.1%	19.5%	18.0%	23.8%	26.4%
	20	22.5%	19.3%	18.0%	23.6%	26.8%
	30	22.2%	19.3%	17.9%	23.4%	25.6%
121	10	23.1%	20.9%	19.2%	23.9%	41.5%
	20	22.9%	20.7%	19.1%	23.6%	41.4%
	30	22.7%	20.3%	18.9%	23.4%	41.2%

Both approaches, SVM and Kernel Regression, perform predictions only at predefined route segments. However, since GPS data consist of irregular points in space, both of these methods relied on interpolation at predefined route locations, such as bus stops in [3]. We therefore performed interpolation at all bus stops of the route, which in fact resulted in a consensus in training data across all methods. To be more precise, for predictions from bus stop p_k onward, all approaches used as training data the historical space-time trajectories $Traj(p_k)$. The key difference is that for SVM and Kernel Regression the cumulative distances $dist_{ijk}$, underlying $Traj(p_k)$, coincide exactly with subsequent bus stops beyond p_k , whereas in our approach they correspond directly with the raw GPS measurements.

5.5 Experimental Results

In Table 2 we see the mean absolute relative errors for each method. The first thing to note is that our Additive Models (BAM, EAM, and AMM) outperformed the Kernel Regression and SVM in all scenarios. SVM’s overall performance was notably worse than any of the other methods. The main comparisons of interest are thus between the Kernel Regression approach and each one of our Additive Models. The Wilcoxon paired test revealed statistically significant differences between the Kernel Regression and all our proposed three Additive Models, in all scenarios. Further, the Wilcoxon paired test revealed that in all scenarios the AMM outperformed all other methods. Another observation from Table 2 is that the size of the training data does not seem to affect performance of any of the 5 methods.

To give a more detailed view of the results, we show in Figure 7 boxplots of absolute prediction errors aggregated across all training data sets, 10, 20, and 30 days, and further stratified by route and prediction distance. The boxplots are displayed for all methods except for SVM as their performance was greatly inferior for larger distances and only interfered with visualization. It should be noted that the SVM approach of [3] was specifically designed to predict only the travel time at next route segment and therefore inferior performance at larger distances may be expected. However, even at the smaller bins the SVM approach was outperformed by all other methods.

As expected, the error increases with distance from bus stops beyond which the predictions were made. We note that all distributions were right skewed with several outliers (as defined by the ends of the boxplot whiskers) mostly due to heavily delayed buses. These outliers are not displayed as they interfere with visualization and do not reveal any significant trends beyond those seen in the boxplots. However, in order to get a sense of this “outlier effect” we plotted the 95th percentiles (dashed lines) along with the boxplots. These lines give us a sense of “worst case” scenario performance of each method.

Although perhaps not visually striking in all distance bins, the AMM statistically outperformed all methods for all routes and in all distance bins (except for the 14km bin on route 627 where no difference existed between EAM and AMM). In the first distance bin, $[0, 1)$, the Kernel Regression method outperformed both BAM and EAM at all routes except for route 603 (where no statistical difference existed). However, in all other distance bins the two Additive Models statistically outperformed the Kernel Regression. Thus, on the whole, the visualization and stratified analysis confirmed the performance order observed in Table 2.

The fact that Kernel Regression outperformed BAM, and EAM in the first distance bin suggests that the Additive Models tend to put more priority on minimizing error in later parts, when it is in fact larger, at the expense of short term predictions. Perhaps this may be fixed by placing more knots at the beginning of the route or through additional features. However, as we discussed before, AMM performed statistically better than all other methods in the first distance bin. This suggests that the random corrective intercept term plays an important roll in rescuing the incorrectly specified cumulative time scales as obtained by interpolation.

It is interesting to note that route 121 showed the highest worst-case scenario across all routes, as demonstrated by the 95th percentile curves. This is further reflected in the highest relative error in Table 2. This fact can perhaps be explained by the fact that the destination of route 121 lies in the heart of the city center.

6. RELATED WORK

A list of related works on bus arrival time prediction may be found in [29, 14, 1]. In this section we present an overview of the main methods, but refer to [29, 14, 1] for a more exhaustive list of references. The discussion is divided into categories based on the type of models in question.

6.1 Historical Data-Based Models

The models falling into this category base predictions of future travel times on historical averages [19, 20, 11, 18]. The proposed algorithm in [19] combined real-time GPS coordinates and current bus speed with historical average speeds of individual route segments. The methods proposed in [20, 11, 18] were all based on averages of similar past bus trajectories. These methods work best when current bus has traveled some distance and its trajectory until current location has sufficient data points that can be compared to historical trajectories. In [19, 11] the analyses were stratified by hour of day by defining time bins. Our proposed

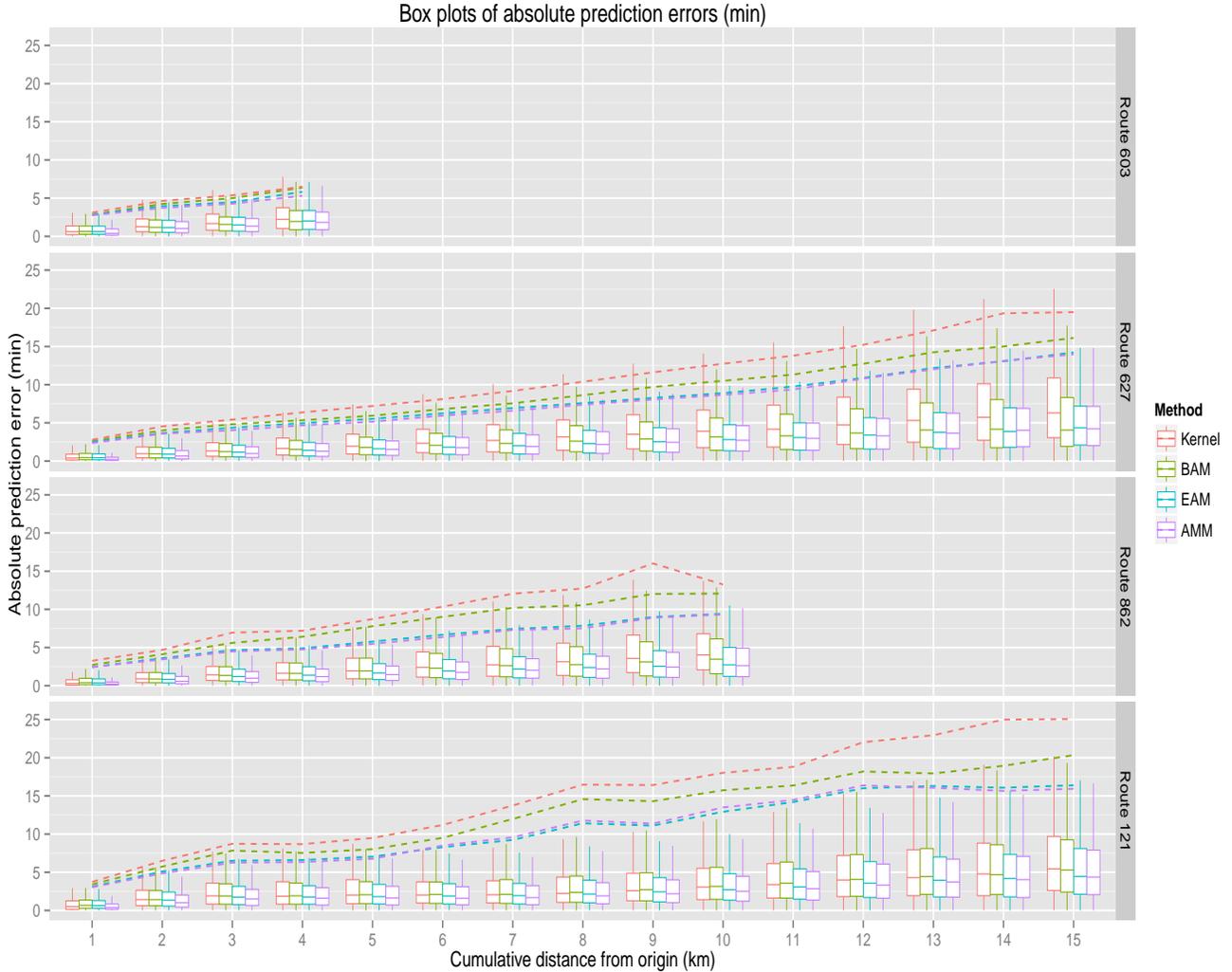


Figure 7: Boxplots of absolute prediction errors for routes 603, 627, 862, and 121, ordered by overall performance of methods. The dashed lines represent 95th percentiles of corresponding absolute errors.

BAM falls into this class of models but in addition models temporal effects as a smooth function as opposed to using categorical binning.

6.2 Regression Models

Regression models predict and explain a response variable through a function of predictor variables. [9] and [15] developed multiple linear regression models using different sets of predictors and both studies indicated that regression models are outperformed by other models. Further, the Kernel Regression method (analyzed in the Experiments section) was demonstrated to have superior performance over regression [18]. However, a great advantage of regression models is that they reveal which predictors have a significant effect on the response. Further, they provide a principled statistical framework for handling features and are highly interpretable. Our proposed EAM and AMM methods enjoy all the benefits of regression models but in addition allow for flexible modeling of nonlinear features through smooth

functions.

6.3 Kalman Filter Models

Kalman filters [10] and other time series models have been proposed for predicting bus arrival times [22, 21, 17, 6]. For the bus prediction problem the most common implementation of Kalman filter involves the assumption that travel time on a given route segment depends on a previously observed travel time at the same route segment [22, 17, 6]. However, [21] took a different approach and assumed that travel time on a given route segment depends on the travel time of a previous route segment. This approach resembles the SVM approach that was implemented for comparison purposes in the experiments section, see subsection below. All of the above methods rely on discretization of either time or space. [22] developed a seasonal autoregressive moving average process for short-term traffic forecasts. [6] treated the average travel time of tagged vehicles in a given time interval as the true value to predict the travel time in the next

time period. In addition to using the travel time in current time interval, [17] also used the last three-day historical data of actual running times in the next time period to predict the next running time. The main limitation of a Kalman filter in the context of our data is the irregularity of observations. Large parts of the data contained time periods where no bus was observed and those time periods of missing data were generally different across different days. Therefore, a clear implementation strategy (e.g. time discretization) that covers all prediction scenarios would require some additional work. However, as noted in [1] Kalman filter give promising results on providing a dynamic travel time estimation. We note that both of our proposed EAM and AMM methods have a Kalman filter flavor as they include the last bus travel time as feature.

6.4 Artificial Neural Network Models

Artificial Neural Network (ANN) models have gained recent popularity in predicting bus arrival times because of their ability to deal with complex and nonlinear relationships between variables [15, 4, 5]. [15] developed an ANN model for prediction of bus travel times using GPS-based data and demonstrated superior performance over multiple linear regression. [4] developed an ANN model that further applied a dynamic Kalman filter algorithm to adjust predictions using bus location information. In order for the models in [4, 5] to be practically implementable Automatic Passenger Count data need to be available in addition to the GPS data [1]. Additive Models in general share the ability of ANNs to flexibly deal with nonlinear relationships. However, they are further easily interpretable like regression models and do not suffer from slow learning process as reported for ANNs [1, 7]. As implementation of ANN involves delicate setup of construction parameters (i.e., input variables, hidden layers, etc.) and none of the ANNs above were directly applicable to our setting we did not include ANNs in our experimental comparison. However, we did implement an SVM, discussed in the next subsection, which is a method that shares some functionalities with ANNs.

6.5 Support Vector Regression Models

SVM and Support Vector Regression (SVR) have demonstrated their success in time-series analysis and statistical learning [28, 3]. [28] compared their SVR algorithm to baseline predictors for prediction of travel time on highways and demonstrated superior performance. [3] proposed SVM for travel time predictions and pointed out that unlike the traditional ANN, their method is not amenable to the overfitting problem. However, they also indicated that when SVM is applied for solving large problems the computation time becomes a problem.

7. CONCLUSIONS

In this paper we discussed the problem of predicting travel times of public buses based on GPS data. We proposed Additive Models as a flexible and a statistically principled framework for model building. We modelled cumulative travel time as a sum of linear terms and smooth functions of predictor variables. We showed that by including a random intercept in the model we were able to correct for an interpolation error incurred when normalizing space-time trajectories onto a cumulative time scale. We demonstrated

on a large real-world GPS data that our proposed Additive Models achieved superior performance as compared to other existing prediction methods.

8. REFERENCES

- [1] M. Altinkaya and M. Zontul. Urban bus arrival time prediction: A review of computational models. *Int'l Journal of Recent Technology and Engineering*, 2:164–169, 2013.
- [2] Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society, Series B*, 57:289–300, 1995.
- [3] Y. Bin, Y. Zhongzhen, and Y. Baozhen. Bus arrival time prediction using support vector machines. *Journal of Intelligent Transportation Systems*, 10(4):151–158, 2006.
- [4] M. Chen, X. Liu, J. Xia, and S. I. Chien. A dynamic bus-arrival time prediction model based on APC data. *Computer-Aided Civil and Infrastructure Engineering*, 19(5):364–376, 2004.
- [5] S. I.-J. Chien, Y. Ding, and C. Wei. Dynamic bus arrival time prediction with artificial neural networks. *Journal of Transportation Engineering*, 128(5):429–438, 2002.
- [6] S. I. J. Chien and C. M. Kuchipudi. Dynamic travel time prediction with real-time and historic data. *Journal of Transportation Engineering*, 129(6):608–616, 2003.
- [7] M. T. Hagan, H. B. Demuth, and M. Beale. *Neural Network Design*. PWS, Boston, 1996.
- [8] T. Hastie and R. Tibshirani. *Generalized Additive Models*. Chapman & Hall, 1990.
- [9] R. H. Jeong. The prediction of bus arrival time using automatic vehicle location systems data. *A Ph.D. Dissertation at Texas A&M University*, 2004.
- [10] R. E. Kalman. A new approach to linear filtering and prediction problems. *Journal of Basic Engineering*, 82(1):35–45, 1960.
- [11] W.-C. Lee, W. Si, L.-J. Chen, and M. C. Chen. Http: a new framework for bus travel time prediction based on historical trajectories. In *Proc. of the ACM SIGSPATIAL*, pages 279–288, 2012.
- [12] J. Pinheiro, D. Bates, S. DebRoy, D. Sarkar, and R Core Team. *nlme: Linear and Nonlinear Mixed Effects Models*, 2014. R package version 3.1-117.
- [13] Project Webpage. <http://goo.gl/HIuSPX>, 2014.
- [14] R. Rajbhandari. *Bus Arrival Time Prediction Using Stochastic Time Series and Markov Chains*. New Jersey Institute of Technology, 2005.
- [15] Y. Ramakrishna, P. Ramakrishna, and R. Sivanandan. Bus travel time prediction using GPS data. *Proceedings Map India*, 2006.
- [16] D. Ruppert, M. P. Wand, and R. J. Carroll. *Semiparametric Regression*. Cambridge Series in Statistical and Probabilistic Mathematics, 2003.
- [17] A. Shalaby and A. Farhan. Prediction model of bus arrival and departure times using AVL and APC data. *Journal of Public Transportation*, 7(1):41–62, 2004.
- [18] M. Sinn, J. W. Yoon, F. Calabrese, and E. Bouillet. Predicting arrival times of buses using real-time GPS

- measurements. In *IEEE Int'l Conf. on Intelligent Transportation Systems*, pages 1227–1232, 2012.
- [19] D. Sun, H. Luo, L. Fu, W. Liu, X. Liao, and M. Zhao. Predicting bus arrival time on the basis of global positioning system data. *Transp. Res. Rec.: Journal of the Transportation Research Board*, 2034(1):62–72, 2007.
- [20] D. Tiesyte and C. S. Jensen. Similarity-based prediction of travel times for vehicles traveling on known routes. In *Proc. of the ACM SIGSPATIAL*, page 14, 2008.
- [21] L. Vanajakshi, S. Subramanian, and R. Sivanandan. Travel time prediction under heterogeneous traffic conditions using global positioning system data from buses. *Intelligent Transport Systems*, 3(1):1–9, 2009.
- [22] B. Williams and L. Hoel. Modeling and forecasting vehicle traffic flow as a seasonal arima process: Theoretical basis and empirical results. *Journal of Transportation Engineering*, 129(6):664–672, 2003.
- [23] S. N. Wood. Modelling and smoothing parameter estimation with multiple quadratic penalties. *Journal of the Royal Statistical Society (B)*, 62(2):413–428, 2000.
- [24] S. N. Wood. Thin-plate regression splines. *Journal of the Royal Statistical Society (B)*, 65(1):95–114, 2003.
- [25] S. N. Wood. Stable and efficient multiple smoothing parameter estimation for generalized additive models. *Journal of the American Statistical Association*, 99(467):673–686, 2004.
- [26] S. N. Wood. *Generalized Additive Models: an introduction with R*. Chapman & Hall/CRC, 2006.
- [27] S. N. Wood. Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *Journal of the Royal Statistical Society (B)*, 73(1):3–36, 2011.
- [28] C.-H. Wu, J.-M. Ho, and D.-T. Lee. Travel-time prediction with support vector regression. *IEEE Trans. on Intelligent Transportation Systems*, 5(4):276–281, 2004.
- [29] B. Yu, W. H. K. Lam, and M. L. Tam. Bus arrival time prediction at bus stop with multiple routes. *Transportation Research Part C*, 19:1157–1170, 2011.