A Dynamic Prediction of Travel Time for Transit Vehicles in Brazil Using GPS Data.

by Zegeye Kebede Gurmu

A thesis submitted to the department of Civil Engineering & Management University of Twente in partial fulfillment of the requirements for the Degree of Master of Science in

Transportation Engineering and Management

Enschede, The Netherlands

2010

APPROVED BY:

Prof. Dr. Ir Eric Van Berkum Chair of Advisor Committee Dr.Tom Thomas

Abstract

GURMU, ZEGEYE KEBEDE. A Dynamic Prediction of Travel Time for Transit Vehicles in Brazil Using GPS Data.

The objective of this research is to develop a dynamic model that can provide accurate prediction for the Estimated Time of Arrival of a bus at a given bus stop using global positioning system (GPS) data which is collected in Macae, Brazil. The provision of timely and accurate travel time information of transit vehicles is valuable for both operators and passengers. It helps operators to monitor and manage their fleets in real time. It also allows passengers to plan their trips to minimize waiting times. Here, an artificial neural network (ANN) is developed for prediction due to its ability to solve complex non-linear relationships. The results obtained from the overall study are promising and the proposed ANN model can be used to implement an Advanced Public Transport System . The implementation of this system will improve the reliability of the public transport system, thus attracting more travelers to transit vehicles and helping relieve congestion. The performance of the proposed ANN model was compared with a historical average model under two criteria: overall precision and robustness. It was shown that the ANN outperformed the average approach.

Acknowledgements

It is a pleasure to thank the many people who made this thesis possible. First of all, I would like to thank Prof. Eric van Berkum, not only for his technical knowledge and insight, but also for his encouragement and provision of useful feedback during this research. I hold my utmost respect and sincere gratitude to my daily advisor Dr. Tom Thomas. This work would not have been possible without the support from him under whose guidance, I chose this topic. He has been abundantly helpful and has assisted me in numerous ways. I specially thank him for his infinite patience. The discussions I had with him were invaluable. I would like to say a big thanks to Ir. Warner Vonk for his encouragement , important discussions and feedback during this research. Through him , I would like to thank Rápido Macaense and APB Prodata for the provision of data for my research .

I thank the faculty and staff of CEM for their dedicated and kindness. Special thanks to Ir. Annet de Kiewit who has been providing me with all the necessary information concerning my study during my stay. I am usually amazed by your positive thinking, every time we speak.

I am grateful to all my friends for being the surrogate family during the two years I stayed in The Netherlands and for their continued love and support there after. On a different note, many people have been a part of my graduate education and I am highly grateful to all of them. My final words go to my great family,of course that includes Tsedi. I want to thank my family, whose love and guidance is with me in whatever I pursue. Special thanks goes to my brother Asfaw Kebede for supporting me in every aspect from the very beginning. I wish you a speedy recovery. Every day I pray for you.

Table of Contents

List of Tables
List of Figures
Chapter 1 Introduction
1.1 Background
1.2 Research Objectives
1.3 Thesis Organization
Chapter 2 Literature Review
2.1 Overview
2.1.1 Historical Data Based Models
2.1.1.1 Using Average Travel Time
2.1.1.2 Using Average Speed
2.1.2 Time Series Model \ldots 10
2.1.3 Regression Models
2.1.4 Kalman Filtering Model \ldots 12
2.1.5 Machine Learning Models $\ldots \ldots \ldots$
2.1.5.1 Artificial Neural Network Model
2.1.5.2 Support Vector Machines
2.1.6 Summary
Chapter 3 Data

3.1	Data Collection	20
3.2	Data Reduction	21
3.3	Preliminary Analysis	23
	3.3.1 Travel time analysis per section	25
	3.3.2 Running time per section	28
	3.3.3 Travel time over day time period	29
	3.3.4 Dwell time analysis	32
Chapte	er 4 Neural Network Development	36
4.1	Overview	36
4.2	The Network Architechture	37
4.3	Neural Network Development	38
	4.3.1 Training , Test and Validation Data Sets	40
	4.3.2 Activation Function	40
	4.3.3 Training and Learning Functions	41
4.4	The Prediction Formula	42
4.5	Concluding Remark	43
Chapte	er 5 Model Evaluation & Comparison	44
5.1	Model Performance	44
5.2	Comparison between Historical Average & ANN	46
Chapte	er 6 Conclusion	53

List of Tables

Table 3.1	Defined sections for the selected bus route	25
Table 5.1	Comparison of MAPE over the whole day for different section	52
Table 5.2	Comparison of prediction accuracies for section 2 to 5	52

List of Figures

Figure 1.1	Schematic overview of factors influencing the distribution of travel times	
	(Source-[Tu 2008])	3
Figure 3.1	Data collection scheme in the study area (Source: [Chung and Shalaby 2007])	21
Figure 3.2	Bus route LT11	22
Figure 3.3	correlation between reduced residuals of successive sections	26
Figure 3.4	Average variation vs. median travel time over sections and whole trajectory	
	$(y=\bigtriangleup T) \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots $	27
Figure 3.5	correlation between reduced residuals of running times $\ldots \ldots \ldots \ldots \ldots$	28
Figure 3.6	Travel time index as a function of time of day for the study area during	
	weekdays	29
Figure 3.7	Travel time index as a function of time of day during weekend $\hfill \ldots \ldots \ldots$	30
Figure 3.8	Travel time index for each section as function of time of the day	31
Figure 3.9	Median dwell times over each bus stops	32
Figure 3.10	Average dwell time on stop 34 over different hours of weekday \ldots .	33
Figure 3.11	Correlation between Dwell times of consecutive journeys at stop 34 in North	
	bound direction	34
Figure 3.12	Correlation between dwell times of consecutive stops $\ldots \ldots \ldots \ldots$	35
Figure 4.1	Input-Output Network Structure	38
Figure 4.2	Network training process	39
Figure 4.3	Schematic of a bus route with several stops	42

Figure 5.1	frequency of the number of stops away from the current bus location \ldots .	45
Figure 5.2	Observed travel time vs MAPE	46
Figure 5.3	MAPE /100 vs time of day section 2	47
Figure 5.4	$MAPE/100$ vs time of day for section $3 \dots $	47
Figure 5.5	MAPE /100 vs time of day for section 4 $\dots \dots \dots \dots \dots \dots \dots \dots \dots \dots \dots$	48
Figure 5.6	MAPE/100 vs time of day for section 5	48
Figure 5.7	MAPE/100 vs time of day for section 2 to section 3	49
Figure 5.8	MAPE/100 vs time of day for section 2 to section 4	50
Figure 5.9	MAPE/100 vs time of day for section 2 to section 5	51
Figure 5.10	$MAPE/100$ vs time of day for the whole trajectory $\ldots \ldots \ldots \ldots \ldots$	51

Chapter 1

Introduction

1.1 Background

Growing traffic congestion has posed threat to the quality of life of people in many countries over the past few decades. Congestion leads to a decrease in accessibility, travel time loss and air pollution. In developed countries still most of the people use private vehicles. Also in developing countries the level of vehicle ownership is rising at a faster rate. So far, different techniques for mitigation of congestion have been forwarded .One of them is to improve and expand public transport system [Kum and Israr 2007, Jamie et al 2009] . A good public transport is of increasing importance to maintain and improve quality of life by providing mobility and accessibility. Moreover, it helps to secure the environment, brings economic development and increases social cohesion. Different concepts for enhancing public transport have been suggested in the past. One of them is providing travelers with reliable travel information through the help of Advanced Public Transport System (APTS) ,which is one component Intelligent Transportation Systems (ITS) [Jill et al 2002, Jamie et al 2009, Vanajakshi et al 2009]. APTS applications include pre-trip and real-time passenger information systems, automatic vehicle location systems, timed transfers, bus arrival notification systems, and systems providing priority of passage to buses at signalized intersections. Other concepts include provision of comfort, improving stops and stations, transit oriented development etc.

Currently, by using global positioning systems, wireless communication systems, and other devices, passengers in most developed countries are able to get more information, for example, about when a transit vehicle or carpool will arrive.Travel time information is the most preferred information by travelers [Vanajakshi et al 2009, Pu and Lin 2008].However, this information can not be measured directly. As has been explained on [Jeong 2004],the provision of timely and accurate transit travel time information is important because it attracts additional ridership and increases the satisfaction of transit users,which will ultimately result in a decrease in congestion.Though many metropolitan areas in some parts of developed world are providing real-time travel time information on freeways and principal arterial, there are still difficulties in provision of accurate real-time travel time information on signalized urban streets because of the stochastic nature of urban traffic.This is even worse in developing countries considering the undisciplined traffic everywhere and lack of collected data, not to mention the technology. Therefore, to provide passengers with real-time travel time information,we need to develop a good algorithm that can predict travel time with reasonable accuracy.

[Van Lint et al 2003, Tu 2008] pointed out that travel times are the result of traffic flow operations, which in turn are attributed to the interaction between traffic demand and traffic supply characteristics. They have outlined these characteristics schematically on Figure 1.1. From the figure, it can be seen that the fluctuations in traffic demand and supply, which in most cases depend upon each other, result in travel time distribution. As far as these fluctuations continue to exist, we will frequently experience variability in, for example, link travel time irrespective of the type of vehicle, be passenger cars or transit vehicles. It is quite natural that high degree of variability will pose difficulty while predicting travel time information. From the traveler's perspective, a decrease in travel time variability reduces the uncertainty in decision-making about departure time and route choice as well as the anxiety and stress caused by such uncertainty. Therefore, knowledge of travel time variability is valuable for improving the reliability of traffic information services and increasing the accuracy of travel time predictions. This can theoretically be achieved by investigating the impacts of each the factors shown on Figure 1.1 on the distribution of travel time. However, due to limited resources in developing countries, it is not easy to get an information or data to study the impacts of all those factors. Therefore, prediction of travel times need to be done using the available data along with logical assumptions. For example, it can be assumed that travel demand at a certain time interval of the day is constant during weekdays etc.



Figure 1.1: Schematic overview of factors influencing the distribution of travel times (Source-[Tu 2008])

The need for the model or technique to predict transit travel time using AVL data is increasing. Although quite reasonable number of researches on this topic have been conducted [Lin and Zeng 1999, Patnaik et al 2004, Jeong 2004], it has been shown that more research on this topic is still required. This thesis mainly focuses on the development of algorithms for prediction of travel time for buses in Brazil using available GPS data in order to provide reliable travel time information for transit users.

1.2 Research Objectives

As indicated earlier the increase in the transit ridership and the satisfaction of transit users can be achieved by the deployment of pre-trip and real-time traveler information systems. Despite the fact that Bus arrival time is important information for passengers, providing it is not an easy task. For example, [Patnaik et al 2004] described that bus arrivals at stops in urban networks are difficult to estimate because travel times on links, dwell times at stops, and delays both at signalized and non-signalized intersections fluctuate both spatially and temporally due to the impact of the factors shown on Figure 1.1 on travel time distribution. In order to consider for such impacts, sound algorithms which could accurately predict bus arrival times are needed. However, developing such models while considering the effects of the aforementioned factors that influence travel time is a complex task, specially when one considers data limitation.

The main objective of this thesis is to develop and test a model that can provide accurate prediction for the Estimated Time of Arrival of a bus at a given bus stop using GPS data. The scope of this thesis encompasses:

- Analysis of the travel time variability and computation of the historical mean travel time. This could be taken as a baseline prediction value on which prediction algorithms will be based.
- Analysis of systematic and random variations of travel times.
- Looking for temporal and spatial correlations between travel times in order to improve the baseline prediction algorithm further. For example temporal correlations could be looked between travel times of successive journeys and spatial correlation between travel times of successive route sections for same journeys.
- Analysis of the impact of dwell time at stops on the prediction model.

• Testing of the prediction algorithms .Once the algorithms have been developed, they will be tested using real sets of data to see their performance.

1.3 Thesis Organization

This section describe the contents of each chapter. The thesis consists of 6 chapters including the introductory chapter.

Chapter 2 presents an extensive literature review on prediction models used for forecasting traffic states, i. e. , traffic flow and travel time. Special attention is given to their application for bus travel time prediction.

Chapter 3 presents the data collection scheme, data reduction and preliminary data analysis done to see the characteristics of the data.It discusses different definitions and notions of travel time.

Chapter 4 outlines different concepts about the proposed model , i.e ,the artificial neural network model and how it is developed.

Chapter 5 presents evaluation of the performance of the proposed model and the comparison with a historical average model.

Chapter 6 concludes the whole research and gives some research recommendations for the future

Chapter 2

Literature Review

2.1 Overview

A variety of prediction models for forecasting traffic states such as travel time and traffic flow have been developed over the years. The five most widely used models include historical data based models [Williams and Hoel 2003, Jeong 2004], time series model [Thomas et al 2010, Al-Deek et al 1998], regression models [Jeong 2004, Ramakrishna et al 2006], Kalman filtering model [Chien et al 2002, Shalaby and Farhan 2004] and machine learning models [Bin et al 2006, Yasdi 1999, Jeong 2004]. [Kirby et al 1997, Zheng et al 2006], however, discussed that no single predictor had yet been developed that presented itself to be universally accepted as the best, and at all times, an effective traffic state forecasting model for real-time traffic operation. Recently, hybrid models, e.g. a combination of Kalman filtering and neural network [Chien et al 2002, Chen et al 2004], a combination of time series and kalman filtering [Thomas et al 2010], also drew much attention. Of course, different classifications of prediction methods were suggested by researchers. For example [Sun et al 2007] classified them as models based on historical data, multilinear regression models and artificial neural network models; [Hoogendoorn and Van Lint 2008] classified them in to three as naive methods that use the real (historic) traffic state values, model based methods that use traffic simulation models and data-driven models which use statistical techniques and correlate actual traffic state values to available data whereas [Chien et al 2002] classified them as univariate forecasting models, multivariate forecasting models and artificial neural networks. However, we will stick to the aforementioned five widely used models in the discussion presented on the next subsections. Special attention is given to their application to travel time prediction, mostly for transit vehicles.

2.1.1 Historical Data Based Models

This type of prediction model gives the current and future bus travel time from the historical travel time of previous journeys on the same time period. The current traffic condition is assumed to remain stationary. [Williams and Hoel 2003] pointed out that the phenomenon that traffic conditions follow nominally consistent daily and weekly patterns leads to an expectation that historical averages of the conditions at a particular time and day of the week will provide a reasonable forecast of future conditions at the same time of day and day of the week. Therefore, these models are reliable only when the traffic pattern in the area of interest is relatively stable, e.g. rural areas. [Lin and Zeng 1999] developed four simple GPS data-based arrival time estimation algorithms based on historical data to provide real time information for Blacksburg which is one of the rural areas in USA where traffic pattern is relatively stable. The four algorithms were developed with different assumptions on input data. The authors claimed that their algorithms outperform several algorithms. Bus location data, schedule information, the difference between schedule and actual arrival time (schedule adherence), and waiting time at time-check stops were used as a main input in developing the models. Since the algorithms were developed basically for buses in the rural areas where congestion is minimum, the authors did not consider the effect of traffic congestion and the dwell time at bus stations. Interestingly, they evaluated the performance of their four algorithms in terms of overall precision, robustness, and stability. The overall precision was calculated to determine the average deviation of the predicted arrival from the actual arrival time. The robustness measure was computed to determine if an algorithm would occasionally give a prediction that is far off the actual arrival time. The stability measure was to check if the prediction given by an algorithm fluctuates from time to time. Using a real life data, they tested their algorithms and based on the aforementioned criteria. the fourth algorithm, that utilized all the four inputs mentioned earlier, showed the best performance.

Historical data based models are usually used in combination with the most recently observed traffic data to provide real-time travel information. Though a variety version of historical data based models have been developed for prediction of travel time, for the sake of simplicity, they are classified here broadly as models that use average travel time and models that use average speed in order to give the estimated value of travel times.

2.1.1.1 Using Average Travel Time

These models use the historical average travel time directly or in combination with other inputs in some way to give bus arrival time. In most researches they were developed for comparison purpose [Jeong 2004, Ramakrishna et al 2006, Shalaby and Farhan 2004, Vanajakshi et al 2009]. And in almost all those researches they were outperformed by the respective proposed main algorithms. However, they were also shown to outperform some of the models, e.g. multilinear regression models [Jeong 2004, Ramakrishna et al 2006]. [Chung and Shalaby 2007] developed an expected time of arrival (ETA) statistical model using explanatory variables. The proposed model predicted arrival time from the input of two categories: the last several days' historical data and the current day's operational conditions. An operational strategy was additionally incorporated into the model to reduce the risk that an overestimated arrival time can result in missing the bus. In developing the model, the most notable constraint was the size of historical data to calibrate the model. Unlike transit vehicles, school buses have one run per route per day and their schedules are revised every school year. That necessitated that the ETA model should be based on a method applicable to the relatively small size of historical data. Because the buses were conventionally operated along fixed routes and stops and according to published schedules, their ETA model assumed that the travel time between two stops can be explained by the historical trends of bus travel times and other independent correlated variables. For the operational conditions, the study incorporated schedule adherence and weather condition but did not consider for dwell time due to stable demand. The authors evaluated the performance of the model using data collected from real-world operations of school buses on which a global positioning system-based automatic vehicle location (AVL) system was installed. It was mentioned that the proposed model consistently showed lower levels of prediction error than moving average and regression approaches. With the operational strategy, the model provided a sufficiently reliable service in which approximately 99%–100% of students do not miss the bus, with the tolerable wait time of 162–177 seconds. As has been discussed earlier, average travel time can also be combined with real-time travel time information to develop a dynamic travel time prediction algorithm, e.g. [Chien and Kuchipudi 2003].

2.1.1.2 Using Average Speed

These type of models use the average speed of vehicles over certain links to predict travel times. They are specially applicable to predict travel time using data collected by GPS technology as the distance taken over links can be calculated using the position information. Commonly, these models make use of map matching techniques, which can be established on a Geographic Information System (GIS) Software to estimate the vehicle position and travel times. [Weigang et al 2002] developed a model to estimate bus arrival times at bus stops using GPS information, which was implemented in the SITCUO-Information System for Urban Bus Transportation, in Brasilia. The model consisted of a main algorithm that calculates the estimated arrival time and two sub algorithms to determine the position and the speed of the bus en route. First, the bus route is divided into a number of short, straight lines, sub routes then these lines were modeled as first-degree equation in a plane. When the GPS equipment in the bus transmits its position, speed and other related information to the control center, it is unlikely that this position would coincide with any point on the straight line graphs. Thus, the actual position was mapped to a point on the graph to get the position of the bus. Second, when using the speed information from the GPS, the arrival time at the stop point will be infinite if the vehicle is stationary. In order to solve this problem, they made use of historical bus travel speed along the route segment and current speed of the bus derived from the GPS data. With the improved method and empirical calibration, the results from the developed model were found to be satisfactory in the implementation and experiment. They found the mean error between output results from system and the actual position of bus is less than 8 %. They indicated that these errors could be reduced by increasing the number of lines representing the bus route. The same kind of model has later been developed by [Sun et al 2007] with slight modification. The proposed prediction algorithm combined real-time location data from global positioning system receivers with average speeds of individual route segments, taking into account historical travel speed as well as temporal and spatial variations of traffic conditions. Recalling the estimated average speed to a station, proposed by [Weigang et al 2002], would depend primarily on its historical average speed along the route as the bus is far away from the station, [Sun et al 2007] argued that the current speed of a bus is usually a more important factor influencing how fast the bus will travel over the distance to the station of interest. Their algorithm basically included of two components. The first component consisted of real-time bus tracking model with the purpose of processing the GPS data, projecting them onto the electronic map and then obtaining the distance to each bus station. The second component was a bus arrival time prediction model used to estimate the time to downstream bus station in real time on the basis of the output of the first component and various other factors. The system was implemented as a finite state machine to ensure its regularity, stability, and robustness under a wide range of operating conditions. A case study on a real bus route was conducted to evaluate the performance of the proposed system in terms of prediction accuracy. The results indicated that the proposed system was capable of achieving satisfactory accuracy in predicting bus arrival times and perfect performance in predicting travel direction. However, it was observed that their model less performed during peak hour than off peak hours. This is due to the variation in traffic condition and thus speed increases as the level of congestion increases. The performance of the algorithm was also compared with the algorithm proposed by [Weigang et al 2002] in terms of prediction accuracy and showed an improvement.

Generally, historical data based models require an extensive set of historical data, which may not be available in practice, especially when the traffic pattern varies significantly. These models are not suitable for large cities where both travel time and dwell time experience large variations. Their accuracy largely rely on the similarity between the real time and historic traffic patterns.

2.1.2 Time Series Model

These models assume that the exogenous factors acting upon the dynamical system either remain constant, or can be measured and accounted for in the model, if they vary in time. In terms of traffic, they assume that the historical traffic patterns will remain the same in the future. As has been indicated on [Chien et al 2002], the accuracy of time series models is a function of the similarity between the real-time and historical traffic patterns. Variation in historical travel time data or changes in the relationship between historical and real-time travel time data could significantly cause inaccuracy in the prediction results. To the author's knowledge, these models have not been used for prediction of bus travel time so far. However, they have been used and indicated to be effective for link travel time and traffic volume predictions either alone or in combination with other models, e.g. Kalman filtering [Al-Deek et al 1998, Thomas et al 2010].

2.1.3 Regression Models

These models predict and explain a dependent variable with a mathematical function formed by a set of independent variables [Chien et al 2002]. Unlike historical data based prediction models, these are able to work satisfactorily under unstable traffic condition. Regression models usually measure the simultaneous effects of various factors, which are independent between one and another, affecting the dependent variable. [Patnaik et al 2004] proposed a set of multilinear regression models to estimate bus arrival times using the data collected by automatic passenger counter (APC). They used distance, number of stops, dwell times, boarding and alighting passengers and weather descriptors as independent variables. They indicated that the models could be used to estimate bus arrival time at downstream stops. However, this approach is reliable when such equations can be established. [Jeong 2004] and [Ramakrishna et al 2006] also developed multilinear regression models using different sets of inputs. Both studies indicated that regression models are outperformed by other models. One great advantage of multilinear regression model is that it reveals which inputs are less or more important for prediction. For example, [Patnaik et al 2004] discovered that weather was not an important input for their model. Also [Ramakrishna et al 2006] found out that bus stop dwell times from the origin of the route to the current bus stop in minutes and intersection delays from the origin of the route to the current bus stop in minutes are less important inputs. In general, the applicability of the regression models is limited because variables in transportation systems are highly inter-correlated [Chien et al 2002].

2.1.4 Kalman Filtering Model

Kalman filtering models have elegant mathematical representations (e.g. linear state-space equation) and the potential to adequately accommodate traffic fluctuations with their time-dependent parameters (e.g. Kalman gain) [Chien et al 2002]. It has been used extensively for predicting bus arrival time [Chen et al 2004, Vanajakshi et al 2009, Wall and Dailey 1999, Chien et al 2002, Yang 2005] and many more. Its basic function is to provide estimates of the current state of the system. But it also serves as the basis for predicting future values or for improving estimates of variables at earlier times, i.e., it has the capacity to filter noise [Kalman 1960, Lesniak et al 2009, Thomas et al 2010]. [Yang 2005, Wall and Dailey 1999] presented a short term transit vehicle arrival times prediction algorithm by combining real-time AVL data with an historical data source in Seattle, Washington. Their algorithm consists of two components: tracking and prediction. They used a Kalman filter model to track a vehicle location and statistical estimation for prediction of bus arrival time purpose. As has been tried to mention above, the model relied on the real-time location data and historical statistics of the remaining time to arrival. That is, it assumed that other variables possibly influencing the arrival time as mentioned on [Tu 2008] were implicitly included in the statistics. Therefore, they did not explicitly deal with dwell time as an independent variable. It was mentioned in the literature that some empirical results had shown that the proposed algorithm was flexible enough to function in adverse conditions and was able to produce predictions that could be useful to the rider. It was found that they could predict bus arrival time with less than 12% error (i.e. when the predicted bus arrival time is 15 minutes, 70 percent of the buses will arrive in between 13 and 17 minutes). The algorithm was implemented as a web application finally to provide the predicted arrival times to users. [Shalaby and Farhan 2004] developed a bus travel time prediction model using the Kalman filtering technique. They used downtown Toronto data collected with four buses equipped with AVL and automatic passenger counter (APC). They used five-weekday data in May 2001. Four days of data were used for learning and developing models, and one-day data were used for testing. They developed two Kalman filtering algorithms to predict running times and dwell times separately. However, when they developed a historical average model, a regression model, and a time lag recurrent neural network model, they included dwell times in link travel time. They

defined a link as the distance between two time check point stops and each link included between 2 and 8 bus stops. Consequently, they predicted dwell time only at time check points, not at every stop. To develop a dynamic, real-time model, they updated the predicted time of bus arrival and departure at time check points. Of the 27 stops on the route, their model was updated at only the six time check points. They claimed that Kalman filtering techniques outperformed the historical models, regression models, and time lag recurrent neural network models in terms of accuracy, demonstrating the dynamic ability to update itself based on new data that reflected the changing characteristics of the transit-operating environment.

[Chien and Kuchipudi 2003] developed a travel time prediction model for vehicles with real-time data and historic data. Here also Kalman filtering algorithm was employed for travel time prediction because of its ability to continuously update the state variable with changing observation. Their study, however, concentrated on a comparison of the path-based and link-based travel time values. Results revealed that during peak hours, the historic path based data used for travel time prediction were better than link based data due to smaller travel time variance and larger sample size. The advantage of using historic data over the link-based model is procurable, allowing prediction at any given time, but at the expense of prediction accuracy under congestion situations. [Yang 2005] focused on the traffic characteristics after special events (e.g. conventions, concerts, football match) and predicted the travel time after graduation ceremony using recursive discrete time Kalman filtering as a case study. GPS equipped test vehicles were used for data collection and the predicted travel times at a given instant of time was determined from observed and predicted travel times at the previous time instant. The performance of the model was quantified using mean absolute relative error. The prediction error was found to be around 17.6%. This value was acceptable by traffic engineers given the fact of many uncertainties (e.g. weather, traffic condition, signal timing) associated with such event. [Vanajakshi et al 2009] developed an algorithm based on Kalman filtering algorithm under heterogeneous traffic conditions on urban roadways in the city of Chennai, India. Their motivation was that they believed most studies used data collected from homogeneous lane-disciplined traffic, either directly from the field or indirectly through simulation models. The unique feature of their algorithm is that the discretization had been performed over space rather than

over time unlike the aforementioned Kalman filtering models. It was mentioned that this feature could be used in reflecting the effects of events such as accidents that had taken place in the previous subsection of the route on the travel time predicted for the given subsection. The results obtained from the overall study were promising. Their algorithm outperformed the average approach over 7 days out of 10 days. In general, Kalman filtering algorithms give promising results on providing a dynamic travel time estimation which other most models lack.

2.1.5 Machine Learning Models

Machine learning methods present some advantages with respect to statistical methods: they are able to deal with complex relationships between predictors that can arise within large amounts of data, are able to process non-linear relationships between predictors and are able to process complex and noise data [Ricknagel 2001]. These models can be used for prediction of travel time, without explicitly addressing the (physical) traffic processes [Hoogendoorn and Van Lint 2008]. However, they are location-specific solutions, requiring significant efforts in input- and model selection for each specific application, via for instance correlation analysis, or genetic algorithms or trial-and-error procedures. Results obtained for one location are (typically) not transferable to the next, due to location-specific circumstances (geometry, traffic control, etc.). Artificial Neural Network (ANN) and Support Vector Regression methods are presented under these categories.

2.1.5.1 Artificial Neural Network Model

ANNs have been recently gaining popularity in predicting bus arrival time because of their ability to solve complex non-linear relationships [Chen et al 2004, Ramakrishna et al 2006, Jeong 2004, Chien et al 2002, Park et al 2004]. ANNs, motivated by emulating the intelligent data processing ability of human brains, are constructed with multiple layers of processing units, named artificial neurons. The neurons contain activation functions (linear or nonlinear) and are highly interconnected with one another by synaptic weights. Information can be processed in a forward or feedback direction through fully or partially connected topologies. Meanwhile, the synaptic weights can be adjusted to map the input-output relationship for the analyzed system automatically through a learning process [Hagal et al 1996]. A detailed discussion about artificial neural network is presented on Chapter 4. Unlike the aforementioned models, ANNs can be developed without specifying the form of the function, while the restrictions on the multicollinearity of the explanatory variables can be neglected. [Chien et al 2002] developed an enhanced Artificial Neural Network model to predict dynamic bus arrival time. The so called Back-Propagation algorithm was used. Their motivation was that due to long learning process of ANN, it is usually hard to apply ANNs on-line. Consequently, an adjustment factor to modify travel time prediction with new input of real-time data was developed. They generated traffic volume and passenger demand that AVL can not collect, using Corridor Simulation model (CORSIM) to use them as inputs. For an actual implementation they assumed they could obtain similar data from Automatic Passenger Counters (APC) and AVL systems. Therefore, Automatic Passenger Counters (APC) needs to be deployed in addition to AVL systems if their model is to be implemented practically. In the study, dwell time and scheduled data were not considered. They checked the performance of their model using simulation result. They claimed their model can accurately perform well for both single and multiple stops. On another study, [Chen et al 2004] developed a methodology for predicting bus arrival time using data collected by Automatic Passenger Counter (APC). Their model consisted of an Artificial Neural Network (ANN) model to predict bus travel time between time points and a kalman filter based dynamic algorithm to adjust the predicted arrival time using bus location information. The ANN was trained with four input variables, day-of-week, time-of-day, weather and segment; and produced a baseline estimate of the travel time. The dynamic algorithm then combined the most recent information on bus location with the baseline estimate to predict arrival times at downstream time points. The algorithm not only explicitly considered variables influencing the travel time but also updated it using the real-time APC data. The authors indicated that their model was powerful in modelling variations in bus-arrival times along the service route. It was observed that the dynamic algorithm performed better than the corresponding ANN model because it incorporated the latest bus-arrival information into the prediction. The ANN model also performed better than the timetable. [Jeong and Rilett 2004] also proposed an ANN model for predicting bus arrival times and demonstrated its superior performance as compared with the historical data based and multilinear regression models. Historical data based

model gave superior results, as compared to the multiple linear regression. The authors have tested 12 training and 14 learning functions and the best functions were chosen for the prediction purpose. The advantage of their models was that traffic congestion, schedule adherence and dwell times at stops were considered as inputs for the prediction. [Ramakrishna et al 2006] developed a Multiple Linear Regression (MLR) model and an Artificial Neural Network (ANN) model for prediction of bus travel times using GPS-based data. These models were applied to a case study bus route in Chennai city, India. It was indicated that Artificial Neural Network model performed better than Multiple Linear Regression model.

In general, ANN models have the ability to capture the complex non-linear relationship between travel time and the independent variables. These models have been proved to be effective for the provision of satisfactory bus arrival time information. They could be very useful in prediction when it is difficult or even impossible to mathematically formulate the relationship between the input and output. Though the learning and testing process is inherently delicate and is slow to converge to the optimal solution [Hagal et al 1996], it is still possible to do an off-line training and adapting ANNs to real-time condition if the inputs are chosen carefully.

2.1.5.2 Support Vector Machines

Support vector machines (SVMs) are a set of related supervised learning methods used for classification and regression. While other machine learning techniques, such as ANN, have been extensively studied, the reported applications of SVM in the field of transportation engineering are very few. Support vector machine and support vector regression (SVR) have demonstrated their success in time-series analysis and statistical learning [Chun-Hsin et al 2003]. Since support vector machines have greater generalization ability and guarantee global minima for given training data, it is believed that support vector regression will perform well for time series analysis. [Bin et al 2006] proposed Support Vector Machine (SVM) as a new neural network algorithm to predict the bus arrival time. They pointed out that unlike the traditional ANN, SVM is not amenable to the overfitting problem, and it could be trained through a linear optimization process. This study predicted the arrival time based on the travel time of a current segment and the latest travel time of the next segment. The authors built separate models according to the time-of-day and weather conditions. The developed model was tested using off-line data of a transit route and exhibited advantages over an ANN based model methods. These models have been developed for prediction of travel time on highways, e.g. [Chun-Hsin et al 2003]. They compared their proposed SVR predictor to other baseline predictors, the results showed that the SVR predictor can reduce significantly both relative mean errors and root mean squared errors of predicted travel times. However, [Bin et al 2006] indicated that when SVM is applied for solving large-size problems, a large amount of computation time will be involved. In addition, the methods for selecting input variables and identifying the parameters should be further researched.

2.1.6 Summary

The above extensive review highlights the efforts made so far to predict travel times, mostly for transit vehicles. It is worth mentioning that there were other few new type of models that could not be classified under any of the five commonly used travel time prediction methods but gives satisfactory results. For example, [Kidwell 2001] presented a computer algorithm for predicting the arrival times of public fixed-route buses at their stops, based on real-time observations of the vehicles' geographic positions. The algorithm worked by dividing each transit route into a number of zones. Each zone was arbitrarily drawn to approximate how far a bus would likely travel in about $1\frac{1}{2}$ to 3 minutes. The predictions were based on the most recent observation of a bus passing through each zone. The algorithm tracked a bus' entrance into and exit from a zone, and how much time it spent in the zone. This transit time through a zone was stored in a table. Every time a bus passes through a zone, the length of time spent in that zone was used to update the time table. It was reported that the predictions made by this algorithm were more accurate than the published schedules because they were based on where the bus is right now, and how long it took to cross the distance between the user and the bus on the last run of the route. The algorithm has been tested and proved valid by feeding it a few hours' worth of historical data and then testing its predictions against later data which was not given to the algorithm. The reviewed studies have shown that most arrival time prediction models are based on historical arrival patterns and/or other explanatory variables correlated with the arrival time. The explanatory variables used in the previous studies include historical arrival time (or travel time), schedule adherence, weather condition, time-of-day, day-of-week, dwell time, number of stops, distance between stops and road -network condition. The collection and transmission of data on such variables have been largely made possible using the emerging technologies of wireless communication, AVL (e.g. GPS), APC, and other sensing technologies. The effect of congestion was treated in most models differently. For example, some have used traffic properties like volume and speed from simulation results, some clustered their data into different time periods and some left it at all because their models were based on black box approach. Historical based models were used in areas where congestion is minimum because the models assumed traffic patterns are cyclical. However, it could be argued that it is also be possible to observe such patterns in areas where congestion is sever. This can be investigated from extensive historical data analysis by looking into the distribution of travel time over time of day or day of week and so on. Many of the studies evaluated the performance of their models in terms of accuracy precision. Only few researchers looked into the robustness and the stability of their algorithms. And because some of the researchers used data collected for few days, some for few weeks and others for months and also the traffic characteristics of their test beds were quite different, it is not possible to order the models in accordance with their performances. Even it is difficult to compare the same kind of model developed by different researchers for the reasons stated. However in general, Kalman filtering techniques and ANN, which were used mostly in urban areas, have been shown to outperform historical average and multilinear regression models in many studies using same resolution of data and test beds. Kalman filtering models can be applied on-line while the bus trip is in progress due to its simplicity in calculation. But these models are effective in predicting travel time one or two time periods ahead, and they deteriorate with multiple time steps. While other models are dependent on cyclical traffic data patterns or need independence between dependent and independent variables, ANNs do not require that variables be uncorrelated and/or that they have a cyclic pattern. Multilinear regression models were seen to perform least, though they had been established based on a lot of explanatory variables. New prediction concepts like map matching and vector support regression models were also found to give promising results.

Even the latter was claimed to outperform ANN model. In conclusion, in areas where the are stable demand and similar traffic pattern, historical based models are able to give satisfactory bus arrival time information, so no need to go for complex prediction models. From the review, to get bus arrival time information that considers traffic condition in real-time, map-matching and Kalman algorithms are recommended. However, it can be argued that ANN and Vector support regression models can also be used for providing a better real-time travel time information considering that they are more effective in handling non-linear relationships between the factors that influence the distribution of travel time presented on Figure 1.1, if trained well off-line. And this thesis focuses mainly on development of ANN model for providing transit vehicle users with a real-time travel time information in Brazil.

Chapter 3

Data

3.1 Data Collection

Data were collected using Automatic Vehicle Location (AVL) systems. In these systems, GPS receivers are usually interfaced with a GSM modems and placed in the buses. They basically record point locations in latitude-longitude pairs, speeds of the buses, date and time. Arrival and departure time records at each bus stops are the most important ones. The data were collected from November 2008 to May 2009 for different buses in Macae, Brazil. Data collection scheme is presented on Figure 3.1below.



Figure 3.1: Data collection scheme in the study area (Source: [Chung and Shalaby 2007])

3.2 Data Reduction

Bus line LT11 has been chosen for the case study because it has more number of records as compared to the other bus lines. Records at some stops are very small and hence have been excluded. The route has been shown on Figure 3.2after removing stops with insufficient number of records. The picture was obtained by plotting the coordinates of the GPS data points from one of the files superimposed on a map of Macae. The fleet monitoring system is still relatively new in the considered area ,so the GPS records did not contain a trip identity. The passing times at bus stops or in between were also stored arbitrary in the data-base. It was therefore not possible to identify a bus trip in a straightforward way. The fact that bus stops have the same identity for both directions makes trip identification even more difficult. It is thus also not possible to identify the direction of the bus in a straightforward way. It may be not practical to assign trip identities to the records (although it some countries it is done by the bus driver who needs to reset the GPS machine whenever he starts a new trip). Therefore, trip identification could become easier when all bus stops at least have their own identity number, i.e. that the bus stops in opposite directions have different IDs.

In the trip identification process, first the records were sorted according to the vehicle ID. Thus, per vehicle the trips were assigned. Then, the records were sorted by successive time stamps per



Figure 3.2: Bus route LT11

vehicle (year, month, day and time of the day). Therefore, a trip was recognized when for successive records the bus stops are in the correct order. When the order is broken, or when the days or vehicle IDs are different, a new trip is assigned. The direction of the bus follows from the order of bus stops. Also when at a stop the difference between the departure and arrival time is very large, it means most likely the bus returns back from that stop. In this rather rough assignment process, the differences between passing times were not considered. For example, when this difference is very large, it is less likely that the respective records belong to the same trip. However, in the data analysis, extraordinary time differences were considered. When traffic accidents or any other unique event occur, it is most likely that there exist an observation or a set of observations which appear to be inconsistent with the remainder set of data.For example, when an incident occurs, it could take longer travel time. These observations are called outliers. Outliers were detected using the **3**

sigma formula. Also wrong measurements were observed. So, it is very important to detect these observations and deal with them as soon as possible. For example, at a stop the time of arrival was greater than the time of departure and the difference gave a negative value.

3.3 Preliminary Analysis

Definitions

Travel time between any two stops i & j for a journey 'p' initiated at a certain time interval 't' of a day 'd' is the is defined as the difference between the arrival time at stop 'j' and departure time at stop 'i'.i.e,

$$T_{i-j}^{tdp} = T_{Aj}^{tdp} - T_{Di}^{tdp}$$

Running time between any two stops i & j for a journey 'p' initiated at a certain time interval 't' of a day 'd' is the is defined as the difference between the arrival time at stop 'j' and arrival time at stop 'i'.i.e,

$$RT_{i-j}^{tdp} = T_{Aj}^{tdp} - T_{Ai}^{tdp}$$

Dwell time between any stop 'i' for a journey 'p' initiated at a certain time interval 't' of a day 'd' is the is defined as the difference between the departure time at stop 'i' and arrival time at stop 'i'.i.e,

$$DT_i^{tdp} = T_{Di}^{tdp} - T_{Ai}^{tdp}$$

where,

 T_{i-j}^{tdp} is the travel time between stop j & i for a journey p initiated at time interval t of day d T_{Di}^{dpt} is the departure time of a bust at stop i

 $T_{Aj}^{dpt} \mbox{is arrival time of a bus at stop j}$

 RT_{i-j}^{tdp} is the running time of a bus between stop i & j for a journey p initiated at time interval t of day d

 DT_i^{tdp} is the dwell time of a bus at stop i

In order to calculate the travel time, the travel times were pooled together in 30 minute intervals (e.g.6:00-6:30, 6:30-7:00.etc) and the average value of each interval was determined. This means that all data sets need to be clustered by time period, because transit vehicles have different departure time by time period and also there are no bus schedules in the study area, resulting in different travel patterns. Usually the average travel time from historical data can be taken as a baseline prediction model. It has been indicated that it is possible to improve prediction models from average values by considering errors or variations and correlations between different values of the variable under consideration [Thomas et al 2010]. For example, if the previous bus was slower during the whole trip, it is likely that the next bus will be also slower than normal if only the traffic situation is considered. If the bus is also slower during the first section of the trip, it is likely that the bus will also be slower on the second section of the trip. Therefore, we say that travel time of successive buses or of successive trip sections may be correlated. If this correlations are found, the information about the travel time of a previous bus or trip section can be used to update travel time prediction of the next bus or trip section. The same explanation holds true for correlations between dwell times and running times of successive journeys. Thus, different correlations between travel times of successive trip sections, between dwell times of successive journeys and stops and different running times of successive sections were investigated. Using these correlations, it is possible to predict travel time in a better way than the average travel time calculated from historic journeys.

Intuitively, the travel time of a bus is variable. Because there are traffic lights, congestion or other traffic conditions, we see random variability among successive travel times. Also we normally anticipate high travel demand during rush hours, so we expect different travel times during peak and off-peak periods that would result in systematic variation. The distribution of average travel time over different hours of the day was investigated to study the systematic and random variations.Systematic error which always occurs (with the same value) when we use the instrument (the AVL systems in this case) in the same way, and random error which may vary from observation to observation. The systematic error is sometimes called statistical bias. We control it by using very carefully standardized procedures; in this case for example by using a more precise GPS receivers like differential global positioning systems (DGPS) receivers. The random error (or random variation) is due to factors which we can not usually control. It may be too expensive or we may be too ignorant of these factors to control them each time we measure. It may even be that whatever we are trying to measure is changing. The results of the above two investigations, i.e. correlations and variations are presented in the next subsections.

3.3.1 Travel time analysis per section

Trajectory	South-North stop ID	North-South stop ID
whole trajectory	Stop 0 to 34	Stop 34 to 0
Section 1	Stop 0 to 8	Stop 34 to 26
Section2	Stop 8 to 16	Stop 26 to 21
Section 3	Stop 16 to 21	Stop 21 to 16
Section 4	Stop 21 to 26	Stop 16 to 8
Section 5	Stop 26 to 34	Stop 8 to 0

The whole traject shown on Figure 3.2has been divided in to five sections first.

Table 3.1: Defined sections for the selected bus route

For these sections the average travel times were calculated and translated in to residuals, i.e. the difference between an observed travel time and the average travel time (per interval). When we divide this residual by the square root of the average travel time (per interval) we get a reduced residual. This can indicate by how many percent an observed travel time differs from the average travel time. The reduced residual (r_{i-j}) of a travel time between any two stops i & j for a journey 'p' initiated at a certain time interval 't' of a day 'd' is given by:

$$r_{i-j} = \frac{T_{i-j}^{tdp} - T_{i-j-avg}^{tdp}}{\sqrt{T_{i-j-avg}^{tdp}}}$$

where $T_{i-j-avg}^{tdp}$ is the average travel time between stop j & i for a journey p initiated at time interval t of day d



Figure 3.3: correlation between reduced residuals of successive sections

Figure 3.3above shows a correlation between reduced residuals of successive sections. A correlation coefficient of 0.42 has been found between these relative residuals .Thus, if the bus is faster than average on one section there is a reasonable probability that it will also be faster on another section. Suppose the reduced residuals on trajectory i and i+1 of two consecutive sections, one in x-axis and one in y-axis, are r_x and r_y respectively. Hence, the standard variation in $\frac{(r_x - r_y)}{\sqrt{2}}$ equals to 5.2. Therefore, if we would take all this correlation into account then according to this result, the random variation would be about: $5.2\sqrt{T_{i-j}^{tdp}}$. This means that if there is no systematic variation, for a journey of 15 minutes (900 seconds), the variation in travel time will be $5.2 * \sqrt{900} = 2.6$ minutes.

The uncertainty range for this thesis was chosen by taking the 10th and 90th percentile of the distribution of travel times. Thus 80% of all travel times lie within the uncertainty range ,i.e. the smallest 10% and the largest 10% lie outside the uncertainty range. This uncertainty range could be used to define the total variation in travel time at a certain time interval of the day. Thus, the variation in travel time ΔT is is defined as half the width of the uncertainty range and larger than average travel times have a positive error. Smaller than average travel times have a negative error.

It is given by:

$$\triangle T = 0.5(90^{th}_{percentile} - 10^{th}_{percentile})$$

The variance of the total variation is the sum of the variances in the systematic and random variations. This can be expressed as:

$$\triangle T = \sqrt{\{var(random) + var(systematic)\}}$$

Figure 3.4 presents this variation versus the median travel time (T) over the five sections and over the whole trajectory. It can be seen from the figure that the solid line function fits the data reasonably well. This function is given by:

$$\triangle T = \sqrt{5.2 * T^2 + (0.23 * T)^2}$$



Figure 3.4: Average variation vs. median travel time over sections and whole trajectory ($y = \Delta T$)

Therefore summarizing the above two graphs, it can be concluded that the random variation or minimum prediction error is $5.2\sqrt{T_{i-j}^{tdp}}$. Thus, for a travel time of 15 minutes, for example, the

minimum prediction error is 2.6 minutes. But also a larger systematic variation aroud 23 % has been observed and this obviously will increase the total prediction error to 4 minutes. This could in principle be reduced by taking day-to-day and hour-to-hour variations into account in addition to the spatial and temporal correlations that exist between travel times.

3.3.2 Running time per section

Here only the running travel time is considered and checked if the correlation has been improved or not.



Figure 3.5: correlation between reduced residuals of running times

Figure 3.5 shows the correlation between reduced residuals of successive section considering only the running times, i.e excluding dwell times. Here the correlation was changed only from 0.42 to 0.40 and thus a better correlation can not be derived from running travel times only. This is because there is a high possibility for a bus to stop at signalized intersections and also congestion may occur.

3.3.3 Travel time over day time period

If journey 'p' is just before rush hour peak period, it is likely journey 'p+1' will take longer period. Therefore, it is appropriate to use average travel time in a certain-interval as a baseline prediction. Thus if a bus will leave between 9:00 and 9:30 hours, the expected travel time could be the average travel time of all buses that left between 9:00 and 9:30 in the past. Consequently, for our travel time analysis also we pooled the travel times together in 30 minute intervals and took average value.Figure 3.6shows the travel time index as function of time of the day, i.e. the ratio of the average travel time per weekday and the average travel time over all days. The different lines correspond with different workdays, i.e. Monday to Friday. The upper panel demonstrates the line LT11 in the North bound direction (form Stop 0 to Stop 34) and the lower panel shows the South bound direction (form Stop 34 to Stop 0).



Figure 3.6: Travel time index as a function of time of day for the study area during weekdays

From these data it can be seen that there is a significant variation in (average) travel time over different times of the day. In the evening rush hour (17-18h), the travel times are about 30% higher

than the average (over all time periods and all days, including weekends). During the morning rush hour (around 7.00 - 7.30h), travel times are about 20% higher than average for the south bound direction. Travel time variations during the day are thus significant and should be taken into account. Another observation from the figure is that the differences between working days are small except for the evening peak in the south bound direction. On Fridays (yellow line), the travel times are slightly higher in afternoon and evening. The figure moreover indicates that the evening rush hours starts earlier on Friday, at least for the south bound direction.



Figure 3.7: Travel time index as a function of time of day during weekend Upper panel:North bound direction,Lower panel:south bound direction

Figure 3.7shows the results for the travel time index as function of time of the day, i.e. the ratio of the average travel time for the weekend days (Saturday and Sunday) and the average travel time for these days. Here, the variation in (average) travel time during the day is less significant. The differences with working days, see Figure 3.6, are however significant. The travel times are significantly lower, and no rush hour peaks can be identified. On Sundays traffic is light and travel

times are the smallest.



Figure 3.8: Travel time index for each section as function of time of the day Upper panel:North bound direction,Lower panel:south bound direction

Figure 3.8 presents the travel time index for each section as function of time of the day, i.e. the ratio of the average travel time per working day and the average travel time over all days. The Sections are presented below on table 3.1. Here it is assumed there is no significant difference in travel time between different days of the week, so we took the whole weekdays together. The different lines correspond with the different sections.

The travel times over the sections also show the rush hour peaks. There are however differences between sections. The last two south bound sections, from city center southwards (light blue and yellow lines) show significant higher travel times (up to 60%) than average during the evening rush hour. These sections correspond with the first two sections (green and blue lines respectively) in the opposite direction (north bound) that show the highest travel times compared to the other sections during the morning rush hour. It can be observed that in the third section the index is higher in most of the cases, especially during off peak hours. This might be because usually the buses go to the central station which is located a few miles away from the stop stations for the buses under consideration.

3.3.4 Dwell time analysis

Intuitively it would be expected that dwell time like the average travel time would also be different by the time period. To account for these differences, data were clustered by time of the week and time of the day. The median dwell time for all stops of bus line LT11 is plotted in Figure3.9. It shows the median dwell time on the thirty-five bus stops of the bus line number in both directions aggregating all days together. It has been observed buses do not stop at every bus stop but sometimes just continue the journey. This is expected to be due to the fact that there was no passenger on that specific stop . A stop time of zero was not included in the analysis which can significantly change the plot pattern. However, if the bus has to stop due to travel demand – a passenger - the graph indicates the amount of time a bus must wait in order let passengers embark, disembark and continue the trip.



Figure 3.9: Median dwell times over each bus stops

This is expected to be due to the fact that there was no passenger on that specific stop . A stop time of zero was not included in the analysis which can significantly change the plot pattern. However, if the bus has to stop due to travel demand (a passenger), the graph indicates the amount of time a bus must wait in order let passengers embark, disembark and continue the trip. Figure 3.9 shows that dwell time has a significant impact on the total travel time and is therefore an important parameter for real time traveler information systems in order to make correct travel time predictions.



Figure 3.10: Average dwell time on stop 34 over different hours of weekday

Figure 3.10 shows the average distribution of dwell time over different hours of the day at stop 34. This stop is chosen as it has large number of records as compared to the other stops. Here also the zero dwell times were not considered in the calculation. Dot in the graphs shows the data is not sufficient. As can be seen from the graph, it can be said averagely the dwell time is with an acceptable range, i.e., less than 1 minute. It was also possible to see evening peak dwell times are large, especially in the south bound direction.



Figure 3.11: Correlation between Dwell times of consecutive journeys at stop 34 in North bound direction

Figure 3.11 shows the correlation between dwell times of consecutive journeys at stop 34. Here it is assumed that trips initiated with in 20 minutes are consecutive trips. As can be shown, few data points were found and the correlation is really weak which is about 0.15. This correlation will not that much help in improving a travel time prediction. For other stops also similar results were obtained, an of course even with far less number of data points.



Figure 3.12: Correlation between dwell times of consecutive stops

The next correlation considered was between dwell times of consecutive stops. This is presented on Figure 3.12. Here also it can be seen that the correlation is very weak which is about 0.12. The last correlation investigated was between dwell times of different consecutive sections. This also resulted in a weaker correlation. Therefore, these different correlation analysis results show that it is complicated to derive a certain pattern out of dwell times and hence the author chooses not consider running and dwell times individually.

In general travel time variations during the day are thus significant and should be taken into account. It can be concluded that there is a significant variation in (average) travel time during the day. In the evening rush hour (17-18h), the travel times are about 30% higher than the average (over all time periods and all days, including weekends). It was not possible to get satisfactory linear relationships to improve the baseline prediction between travel times, running times and dwell times. Hence, it was decided to use a non linear models such as artificial neural network models because these models are able to get any kind of relationship that exist between travel times either temporal or spatial relationships.

Chapter 4

Neural Network Development

4.1 Overview

ANNs, motivated by emulating the intelligent data processing ability of human brains, are constructed with multiple layers of processing units, named artificial neurons. The neurons contain activation functions (linear or nonlinear) and are highly interconnected with one another by synaptic weights. Information can be processed in a forward or feedback direction through fully or partially connected topologies. Meanwhile, the synaptic weights can be adjusted to map the input-output relationship for the analyzed system automatically through a learning process [Hagal et al 1996]. As explained on the second chapter, ANNs have been recently gaining popularity in predicting bus arrival time because of their ability to solve complex non-linear relationships [Chen et al 2004, Ramakrishna et al 2006, Jeong 2004, Chien et al 2002, Park et al 2004]. ANNs learn from examples and capture subtle functional relationships among the data even if the underlying relationships are unknown or hard to explain. Thus ANNs are well suited for problems whose solutions require knowledge that is difficult to specify. Another advantage of ANNs is, they can generalize. After learning the data fed into them (a sample or example), ANNs can often correctly infer the unseen part of a population even if the example data contain noisy information. To gain the maximum benefit from neural network, there should be enough data or observations[Zhang et al 1998].

4.2 The Network Architechture

Many different ANNs have been proposed in the past few decades for forecasting purposes. The most popular connected multilayer perceptron (MLP) neural network architecture was chosen in this study as it can approximate almost any function if there are enough neurons in the hidden layers, i.e it has a very good capability of arbitrary input-output mapping. And it is also easy to implement. The ANN architecture is typically composed of a set of nodes and connections arranged in layers. In this thesis three layers have been used: an input layer, a hidden layer, and an output layer. The first layer is an input layer when external information is received. The last layer is an output layer where problem solution is obtained. Usually one or two hidden layers are used in between the two aforementioned layers to predict reasonably well. The actual processing in the network occurs in the nodes of hidden layer and the output layer. The input layer is simply at which the data vector is fed into the network. It then feeds into the hidden layer which in turn feeds into the output layer. The connections are typically formed by connecting each of the nodes in a given layer to all of the neurons in the next layers. The hidden layer generates weight of these connections and bias parameter during the training process. It is the hidden nodes in the hidden layer that allow the neural network to detect the feature to capture the pattern in the data, to perform non linear mapping between input and output variables. A single hidden layer has been proved to be sufficient for ANNs to approximate any non linear functions [Zhang et al 1998]. A suitable number of nodes in the hidden layer was determined by experiment; a typical number was about 15. A fully connected MLP with one hidden layer is presented on Figure 4.1



Figure 4.1: Input-Output Network Structure

4.3 Neural Network Development

Though the basic training procedure of ANNs is the same, the accuracy of the result is greatly dependent upon the type of input /output combinations. And all the ANNs developed for prediction of travel time so far differ in their input-output combinations. This thesis also presents a unique input-output combination for the prediction purpose. Four input variables have been used in combination with one output variable. The input variables are the time interval of the day (X_1) , the stop number that is currently being passed (X_2) , the time the bus takes to travel up to this stop after the trip has been initiated (X_3) and the stop number where we want to know the travel time information (X_4) . These variables contribute to the variation of bus travel time. It is worth noting that the choice of input variables is limited to the available data. The output is the travel time the bus will take to reach the required location from the currently passed stop location (Y). All the variables were coded so that they could be understood by Matlab. Since variable X_3 indicates a real time travel time information, our neural network is considered as a dynamic prediction model.

As has been indicated earlier, commonly neural networks are trained so that a particular input leads to a specific target output. Such a situation is presented below on Figure 4.2 . As can be seen from the figure, the network is adjusted, based on a comparison of the output and the target, until the network output matches the target. Typically many such input/target pairs are used. This kind of learning to train a network is called supervised learning. The supervised training methods are commonly used, but other networks can be obtained from unsupervised training techniques or from direct design methods. Unsupervised networks can be used, for instance, to identify groups of data.



Figure 4.2: Network training process

The input variables are fed into the network in such a way that $Y = f(X_1, X_2, X_3, X_4)$. This can be written formally as Y = f'(wX + b). The scalar inputs X are transmitted through a connection that multiplies their strength by the scalar weights w, to form the product wX, again scalar. Here the weighted input wX is the only argument of the transfer or activation function f, which produces the scalar output Y. The neurons also have a scalar bias, b. The bias can be viewed as simply being added to the product wX by the summing function or as shifting the function f to the left by an amount b. The bias is much like a weight, except that it has a constant input of 1.

4.3.1 Training, Test and Validation Data Sets

A training and a test sample are typically required for building an ANN forecaster. The training sample is used for ANN model development and the test sample is adopted for evaluating the forecasting ability of the model or to measure its performance. In this thesis a third one called the validation sample is also utilized to avoid the over-training problem or to determine the stopping point of the training process. The symptom of over-training is that the network performs well with data in the training set, while its performance over the test data set (those "unseen" by the network) starts to deteriorate. It is common to use one test set for both validation and testing purposes particularly with small data sets. The first issue to deal with during ANN development is the division of the data into the training and test sets. Although there is no general procedure to do this p, several factors such as the data type and the size of the available data should be considered in the division. The whole data set was first sorted by week number. Then the first 70 % of the data set was taken as a training set where as the next 30% of the data set as a testing set. This division has been used by most researchers. Out of the training set, 20 % of the data set was takes as a validation test. It is worth mentioning that different percentage combinations of training, test and validation sets has been investigated and it was possible to get a minimum mean square error using the above combinations.

Usually, the the need to normalize the input/target data sets before training. This is because the contribution of an input will depend heavily on its variability relative to other inputs. If one input has a range of 0 to 1, while another input has a range of 0 to 2,000,000, then the contribution of the first input will be swamped by the second input. So it is essential to rescale the inputs so that their variability reflects their importance, or at least is not in inverse relation to their importance. The choice of range to which inputs and targets are normalized depends mainly on the activation function f of output nodes, with typically [0, 1] for logistic function and [-1, 1] for hyperbolic tangent function.

4.3.2 Activation Function

Activation functions for the hidden units are needed to introduce non linearity into the network. It determines a nonlinear relationship between inputs and outputs of a node and a network. The sigmoidal functions such as logistic and hyperbolic tangent functions (tanh) are the most common choices. Functions such as tanh or arctan that produce both positive and negative values tend to yield faster training than functions that produce only positive values such as logistic in practice. For this thesis the inputs and targets are scaled to [-1,1] using the tanh function. This function is given by:

$$f(x) = \frac{(exp(x) - exp(-x))}{(exp(x) + exp(-x))}$$

4.3.3 Training and Learning Functions

Training and learning functions are mathematical procedures used to automatically adjust the network's weights and biases [Mathworks]. The training function dictates a global algorithm that affects all the weights and biases of a given network. The learning function can be applied to individual weights and biases within a network. The training procedure chosen was the most commonly used back-propagation algorithm, whose learning process is shown in Figure 4.2. The back propagation neural network which is arguably the most popular algorithm for transportation use. The objective of the training process is to find a weight matrix that minimizes the mean squared error (MSE), which is defined as the average squared error between the ANN predictions and the actual values of travel times. MATLAB offers a number of training and learning functions. Considering the scope of this research only two training functions and one learning function have been used. Out of twelve training functions, Jeong 2004] found that the Bayesian Regularization training function and the Levenberg-Marquardt Back propagation training function outperformed the other ten training functions. Here in this thesis, it was found that the Levenberg-Marquardt Back propagation training function, which is commonly used by most researchers, outperformed the Bayesian Regularization training function. [Jeong 2004] also showed that there was no significant differences in the results from the fourteen tested learning functions. Gradient descent with momentum weight and bias learning function has been used in this study because it is available in the used MATLAB version.



Figure 4.3: Schematic of a bus route with several stops

4.4 The Prediction Formula

Consider a bus route shown on figure 4.3. Suppose a journey is initiated from stop 'O' at a certain time interval 't' and the bus is currently at location 'i' after passing stop 'C'. It is required to provide a travel time information for a person at stop 'S'. The inputs will be the time interval 't', the coded number for stop 'C', the coded number for stop 'S' and the travel time taken from stop 'O' to 'C'. The output will be the travel time the bus will take from stop 'C' to stop 'S'.Therefore, the travel time information from the current bus location 'i' to stop 'S' can be calculated as:

$$T_{S-i}^t = T_{S-C}^t - T_{i-C}^t$$

Where,

 T_{S-i}^t is the travel time between the current location and the stop under consideration for a trip initiated at time interval t of the day

 T_{S-C}^t is the travel time between stop C and S for a journey initiated at time interval t of the day. It is obtained from the ANN and is equal to the output Y.

 T_{i-C}^t is the difference between the current time and the departure time at stop 'C'

4.5 Concluding Remark

Constructing a program for Neural Network is not a difficult task once one understands the concept behind. Basically, it was only several steps of algorithms that are easily followed even by novice practitioners. However, preparing the network for training is a difficult task since the network dealing with a large amount of data. Another problem is when to stop the training. Over training could cause memorization where the network might simply memorize the data patterns and might fail to recognize other set of patterns. Thus, early stopping is recommended to ensure that the network learn accordingly.

Chapter 5

Model Evaluation & Comparison

5.1 Model Performance

Now after the prediction models have been developed, it is necessary to evaluate them in terms of prediction accuracy. Here the ANN model has been evaluated for its performance and compared with that of historical average travel time model. The Mean Absolute Percentage Error (MAPE) was used as the measure of model performance. The MAPE is shown in equation below. It represents the average percentage difference between the observed value (in this case observed arrival times at a bus stop) and the predicted value (in this case predicted arrival times at a bus stop).

$$MAPE = \frac{1}{n} \sum_{i}^{n} \frac{|Y_p - Yobs|}{Y_p} * 100\%$$

Where,

 Y_p is Predicted travel time from current bus stop to target bus stop

 Y_{obs} is Observed travel time from current bus stop to target bus stop

 \boldsymbol{n} is the number of test sets

It was generally found that prediction of travel time in the study area could be given with an MAPE value of 18.3 % using the ANN model. This figure resulted because of larger values of prediction errors above this average value. It was discovered that the ANN model gave a better

prediction of travel time when the station where the travel time information required is located at least 5 stops away from the current bus location. This is illustrated on Figure 5.1, below. The figure shows the frequency of the number of stops away from the current bus location that gave an MAPE value above 18.3% of the test data set.



Figure 5.1: frequency of the number of stops away from the current bus location

Around one fourth (25 %) of the test data set resulted in MAPE higher than the total MAPE i.e 18.3 %. The above figure, thus ,shows that only few test data sets could be observed with a higher MAPE when the number of stops between the current bus location and the stop under consideration becomes larger. On one hand, for an observed travel times with values between 20 to 50 minutes, prediction can be given with less than 10% error. In other words, it is possible to provide a real time travel time information for a person who is 20 to 50 minutes away from the current bus location with only 10 % error. This is illustrated on figure 5.2below. On the other hand, when the observed travel time becomes larger, i.e. greater that 50 minutes, the error in prediction started to increase. This could be due to the fact that when a bus stays longer in a trip, there is a probability that it stops at many of the stops specially when demand is high during peak hours. Therefore, a person who is 20 to 50 minutes away from the current bus location receives the optimal travel time information. The error in prediction for smaller travel times, for example between 5 to 10 minutes, were found to be higher even for the ANN. This indicates that for a short distance trips there are higher travel time variability as one might expect due to a number of factors. For example a bus may wait for 1 minute at a traffic light on a link where the bus usually takes 2 minutes to cross it. So the error for this situation could be exaggerated and may reach up to 40% which in turn will increase the overall MAPE. However, it should be acknowledged that this kind of travel time information is less important as it is with in an acceptable range of waiting time, specially considering the longer waiting time experience of the people in most developing countries.



Figure 5.2: Observed travel time vs MAPE

Therefore, considering the provision of travel time is new in Brazil, the overall performance result of the model is quite promising.

5.2 Comparison between Historical Average & ANN

MAPE from both models were compared for different sections and for the whole traject. The sections have already been defined in the third chapter on Table 3.1 .There are 36 time interval per day each representing 30 minutes. For example, the first time interval represent the time between 6:00 to 6:30 am and so on. The MAPE for both models were drawn against time of day for different sections and combination of them.



Figure 5.3: MAPE /100 vs time of day section 2 $\,$



Figure 5.4: MAPE/100 vs time of day for section 3 $\,$



Figure 5.5: MAPE /100 vs time of day for section 4



Figure 5.6: MAPE/100 vs time of day for section 5

The MAPEs for the ANN model and historical average of individual sections, from section 2 to section 5 are shown in the above four consecutive figures. As can be seen from the figures, it is difficult to define a certain pattern on how the MAPE vary over time of day. Overall, the prediction of the proposed ANN algorithm performed better than the average approach in terms of prediction

accuracy. However, it was observed that at certain time intervals of the day the average approach also performs well but the difference is not much. For instance, considering the above four sections, though the ANN perfomed better during off- peak and peak hours (specially during evening peak),the average approach also sometimes perfomed well during morning peak hours and sometimes during off peak hours. The reason for these comparable results is that the travel times over the individual four sections are approximately between 8 to 18 minutes.



Figure 5.7: MAPE/100 vs time of day for section 2 to section 3



Figure 5.8: MAPE/100 vs time of day for section 2 to section 4

For sections with travel times between 20 to 50 minutes, the ANN outperformed the average approach on more than 70 % of the time intervals of the day. This has been illustrated on figures 5.7 and 5.8 above.

Usually a low MAPE value is desirable for an algorithm. However, an algorithm with a low value of MAPE may occasionally yield a prediction with a large deviation. This is undesirable since it may divert passengers away from the bus stop and eventually cause them to miss the bus. Therefore, it is important to define a second measure that will be used to detect this behavior. It examines the robustness of an algorithm such that its maximum deviation is within a certain range. Here, the robustness measure R_o is defined as: $R_0 = max\{MAPE\}$ of a section at a certain time interval of the day. As can be seen from all figures, i.e figure 5.3 through figure 5.10, the maximum MAPE of the average approach is greater that the corresponding MAPE of the ANN. Hence, the proposed ANN is more robust than the average approach.



Figure 5.9: MAPE/100 vs time of day for section 2 to section 5



Figure 5.10: MAPE/100 vs time of day for the whole trajectory

As the observed travel time increases more than 50 minutes like in the case of figure 5.9 and 5.10, the MAPEs of both models become comparable. But the MAPEs of these bigger sections are less than individual sections, see table 5.1.

Section Name	MAPE (ANN)%	MAPE (Historical Average)
Section 2	18.7%	19.6%
Section 3	13.7%	15.6%
Section4	13.9%	19%
Section 5	12.7%	14.4%
Section 2 to 3	8.8%	10.4%
Section 2 to 4	7.2%	10.3%
Section 2 to 5	8.8%	9.2%
Whole section	8.6%	8.9%

Table 5.1: Comparison of MAPE over the whole day for different section

Another general prediction accuracy measure can be defined to see if the ANN has the capacity to filter noise by combining the prediction errors of the individual sections and comparing them with a bigger section that consists of the individual sections. For example, the prediction accuracy of a section from the begining of section 2 to the end of section 5 , represented as Pac_{2-5} , can be calculated using the prediction errors of section 2 through section 5. Let for *m* number of buses, the observed travel times of the sections are represented by $To_2, To_3, To_4 and To_5$ and the corresponding predicted travel times are represented by Tp_2, Tp_3, Tp_4 and Tp_5 . Then the prediction accuracy from section 2 to section 5 Pac_{2-5} is given by:

$$Pac_{2-5} = \sum_{j=1}^{m} \frac{|Tp_2 + Tp_3 + Tp_4 + Tp_5 - To_2 - To_3 - To_4 - To_5|}{Tp_2 + Tp_3 + Tp_4 + Tp_5} \frac{*100}{m}$$

Using the above formula the prediction accuracies for both the ANN model and historical average model were calculated and the result is presented below on table 5.2.

	MAPE of whole section 2 to 5	$Pac_{2-5}(\sec 2 + \sec 3 + \sec 4 + \sec 5)$
ANN model	8.8 %	7.4%
Average model	9.2%	9.1%

Table 5.2: Comparison of prediction accuracies for section 2 to 5

As can be seen from the above table the MAPE of the ANN has decreased when different sections were combined to calculate the prediction accuracy which indicates that indeed the ANN has the ability to filter noises. However, the MAPE of the average model has remained almost the same. Overall, the ANN outperformed the average approach in terms of both prediction accuracy and robustness.

Chapter 6

Conclusion

This research explored mainly the use of an artificial neural network for travel time prediction of transit vehicles under undisciplined traffic conditions given GPS data. It is part of the background work required to implement Advanced Public Transport System (APTS) in Macae, Brazil and no attempt has been made before to the author's knowledge. The requirements for such an APTS system include real-time data collection and methodologies for quick travel time prediction.Before the ANN development, spatial and temporal correlations between travel times, running times and dwell times were investigated. The common variations in travel times ,i.e systematic and random variations have also been studied.

The research used GPS data which is quite a new data collection scheme in the area and an algorithm based on an artificial neural network has been developed for predicting the travel time of transit vehicles between any two bus stops under consideration. As one might expect, the traffic conditions in developing countries are different with heterogeneity lack of lane discipline. Therefore, the prediction algorithm needs more care during development as compared to short-term travel time prediction that used homogeneous data in most previous reported studies on . The lack of historic data and permanent data collection schemes add to the difficulties. Therefore, this research is one of the first approaches to predict travel time under a mixed traffic condition, taking an urban route in Macae, Brazil, as a case study. The predicted travel time between subsections of the route under

consideration were compared with the measured data. The performance of the model is also tested and compared with a historical average approach, where the predicted travel time is taken to be the average of the travel times of previous buses that traveled between any two stops under consideration. Prediction accuracy or overall precision and robustness were used as a performance measures. The overall precision measure determines the average deviation of the predicted travel time from the observed travel time. The robustness measure determines if an algorithm will occasionally give a prediction that is far off the actual arrival time. The ANN outperformed the average approach in both performance measures. The results obtained from the overall study are promising and the proposed ANN model can be used to implement an Advanced Public Transport System to predict the arrival time at bus stops in the study area where there are undisciplined traffic flow, especially considering the fleeting monitoring system is new in the area. The implementation of this system will improve the reliability of the public transport system, thus attracting more travelers to buses and helping relieve congestion in Macae, Brazil. As has been tried to indicate earlier, the fact that the buses follow irregular schedule and lack of collected traffic data such as flow, density and weather data contributed to the prediction error. With an appropriate data collection scheme, it can be claimed that the results obtained can be further improved.

Bibliography

- [Lesniak et al 2009] Andrzej Lesniak, Tomasz Danken, Marek Wojdyla (2009), "Application of Kalman Filter to Noise Reduction inMultichannel Data", S C H E D A E I N F O R M A T I C A E, Vol 17/18,Poland.
- [Al-Deek et al 1998] Al-Deek, H., M. D'Angelo and M. Wang (1998). Travel Time Prediction with Non-Linear Time Series. Proceedings of the ASCE 1998 5th International Conference on Applications of Advanced Technologies in Transportation, Newport Beach, CA, pp.317-324.
- [Bin et al 2006] Bin, Yu, Zhongzhen, Y., and Baozhen, Y. (2006). Bus arrival time prediction using support vector machines. Journal of Intelligent Transportation Systems, 10(4), 151–158.
- [Chen et al 2004] Chen M, Liu X B , Xia J X (2004). A dynamic bus arrival time prediction model based on APC data . ComputerAided Civil and Infrastructure Engineering , 19 :364 -376.
- [Chien et al 2002] Chien, S.I.J., Ding, Y. and Wei, C. (2002). "Dynamic Bus Arrival Time Prediction with Artificial Neural Networks." Journal of Transportation Engineering, Volume 128, Number 5, pp. 429-438.
- [Chien and Kuchipudi 2003] Chien, S.I.J. and Kuchipudi, C.M. (2003). "Dynamic Travel Time Prediction with Real-Time and Historic Data." Journal of Transportation Engineering, Volume 129, Number 6, pp. 608-616.

- [Chun-Hsin et al 2003] Chun-Hsin, W., Chia-Chen, W., Da-Chun, S., Ming-Hua, C. & Jan-Ming,
 H. (2003), Travel time prediction with support vector regression, in 'Proceedings of the
 2003 IEEE Conference on Intelligent Transportation Systems', IEEE, Shanghai, China.
- [Kum and Israr 2007] Dewan, K.Kum and Ahmad Israr (2007), Carpooling: A Step To Reduce Congestion. Engineering Letters 14:1,EL_14_1_12,Advanced online publication, India
- [Chung and Shalaby 2007] E.-H. Chung and A. Shalaby.(2007) "Expected time of arrival model for school bus transit using real-time global positioning system-based automatic vehicle location data". Journal of Intelligent Transportation Systems, 11(4):157—167
- [Hagal et al 1996] Hagan, M. T., Demuth, H. B., and Beale, M. (1996). "Neural network design", PWS, Boston.
- [Higatani et al 2009] Higatani, Akito, Kitazawa, Toshihiko, Tanabe, Jun, Suga, Yoshiki, Sekhar, Ravi and Asakura, Yasuo (2009)'Empirical Analysis of Travel Time Reliability Measures in Hanshin Expressway Network', Journal of Intelligent Transportation Systems,13:1,28 — 38
- [Hoogendoorn and Van Lint 2008] Hoogendoorn, S. & Van Lints, H (2008). Reistijdvoorspellingen en reistijdbetrouwbaarheid. NM Magazine, 2008, 3.
- [Jill et al 2002] Hough Jill.A, Crystal Bahe, Mary Lou Murphy and Jennifer Swenson (2002), "Intelligent Transportation systems:Helping Public Transport Welfare To Work Initiatives", Upper Great Plains Transportation Institute, North Dakota State University
- [Jamie et al 2009] Houghton Jamie, John Reiners and Colin Lim (2009), "Intelligent Transport System:How cities can improve mobility", IBM Institute for Business Value, USA
- [Van Lint et al 2003] J.W.C. van Lint, S.P. Hoogendoorn and H.J. van Zuylen (2003), "Accurate freeway travel time prediction with state-space neural networks under missing data",Delft University of Technology, The Netherlands

- [Jeong 2004] Jeong R.H.(2004), "The Prediction of Bus Arrival time Using Automatic Vehicle LocationSystems Data", A Ph.D. Dissertation at Texas A&M University
- [Kidwell 2001] Kidwell, B.(2001)" Predicting Transit Vehicle Arrival Times". GeoGraphics. Laboratory, Bridgewater State College, Bridgewater, Mass
- [Jeong and Rilett 2004] Jeong, R., and Rilett. L. (2004). "The Prediction of Bus Arrival Time using AVL data." 83rd Annual General Meeting, Transportation Research Board, National Research Council, Washington D.C., USA.
- [Vanajakshi et al 2009] L. Vanajakshi, S.C. Subramanian, and R. Sivanandan.(2009) Travel time prediction under heterogeneous traffic conditions using global positioning system data from buses,IET Intell. Transp. Syst. 3, p1 -9
- [Weigang et al 2002] Li Weigang; Koendjbiharie, W.; de M Juca, R.C.; Yamashita, Y.; Maciver, A. (2002)" Algorithms for estimating bus arrival times using GPS data", Intelligent Transportation Systems, 2002. Proceedings. The IEEE 5th International Conference on Volume, Issue ,Page(s): 868 - 873
- [Lin and Zeng 1999] Lin, W.H. and Zeng, J. (1999). "Experimental Study of Real-Time Bus Arrival Time Prediction with GPS Data." Transportation Research Record, 1666, pp. 101-109.
- [Kalman 1960] Kalman R.E. A new approach to linear filtering and prediction problems, Transactions of the ASME–Journal of Basic Engineering, 82 (Series D), 1960, pp.35–45.
- [Kirby et al 1997] Kirby, H., Dougherty, M., and Watson, S. (1997). "Should we use neural networks or statistical models for short term motorway traffic forecasting?" Int. J. Forecasting, 13(1), 43–50.
- [Tam and Lam 2006] Mei Lam Tam, William H. K. Lam (2006), "Real-Time Travel Time Estimation Using Automatic Vehicle Identification Data in Hong Kong". ICHIT 2006 352-361

- [Patnaik et al 2004] Patnaik.J, S.Chein,and A.Bladihas (2004) "Estimation of Bus Arrival Times. Using APC Data". Journal of Public Transportation, Vol. 7, No. 1, pp. 1–20. 7
- [Pu and Lin 2008] Pu, W., J. Lin (2008) "Urban Street Travel Time Prediction Using Real Time Bus Tracker Data". Transport Chicago 2008, Chicago, IL, June 6, 2008
- [Ramakrishna et al 2006] Ramakrishna Y, Ramakrishna P, Sivanandan R (2006), "Bus Travel Time Prediction Using GPS Data", proceedings, Map india 2006
- [Ricknagel 2001] Recknagel, F., 2001. Applications of machine learning to ecological modelling. Ecol. Model. 146, 303–310.
- [Shalaby and Farhan 2004] Shalaby, A., and A. Farhan.(2004) "Bus Travel Time Prediction for Dynamic Operations Control and Passenger Information Systems". CD-ROM. 82nd Annual Meeting of the Transportation Research Board, National Research Council, Washington D.C.
- [Sun et al 2007] Sun, D. Luo, H. Fu, L. Liu, W. Liao, X. Zhao, M.(2007)," Predicting Bus Arrival Time on the Basis of Global Positioning System Data", TRANSPORTATION RESEARCH RECORD, Washington, DC; National Academy.
- [Park et al 2004] T. Park, S. Lee, and Y.-J. Moon (2004). Real time estimation of bus arrival time under mobile environment. In Computational Science and Its Applications - ICCSA 2004, International Conference, volume 3043 of Lecture Notes in Computer Science, pages 1088–1096.
- [Thomas et al 2010] T. Thomas, W.A.M. Weijermars and E.C. Van Berkum (2010) .Predictions of Urban Volumes in Single Time Series, IEEE Transactions on intelligent transportation systems, vol. 11, no1, pp. 71-80
- [Tu 2008] Tu, H. (2008). "Monitoring Travel Time Reliability on Freeways". Transportation and Planning. Delft,. Delft University of Technology. PhD series

- [Yasdi 1999] Yasdi, R., 1999. Prediction of road traffic using a neural network approach. Neural Computing and Applications 8, pp. 135–142
- [Wall and Dailey 1999] Wall, Z., and D. J. Dailey (1999). "An Algorithm for Predicting the Arrival Time of Mass Transit Vehicles Using Automatic Vehicle Location Data". CD-ROM. 78th Annual Meeting of the Transportation Research Board, National Research Council, Washington D.C
- [Williams and Hoel 2003] Williams, B. and Hoel, L. (2003). Modeling and forescating vehicle traffic flow as a seasonal arima process: Theoretical basis and empirical results. Journal of Transportation Engineering, 129(6):664–672.
- [Yang 2005] Yang, J.-S. (2005). Travel time prediction using the GPS test vehicle and Kalman filtering techniques. In American Control Conference, Portland, Oregon, USA.
- [Zheng et al 2006] Zheng, W., Lee, D.-H., and Shi, Q. (2006) Short-Term Freeway Traffic Flow Prediction: Bayesian Combined Neural Network Approach, Journal of Transportation Engineering, 132(2), 114–121.
- [Zhang et al 1998] Zhang, G., Patuwo, B. E., & Hu, M. Y. (1998). Forecasting with artificial neural networks: The state of the art. International Journal of Forecasting, 14 (1), 35–62.

[Mathworks] The Mathworks Inc, MATLAB Ver 6.1 software, Massachusetts, USA