Support Vector Machine Technique for the Short Term Prediction of Travel Time

Lelitha Vanajakshi, Member, IEEE, and Laurence R. Rilett

Abstract— A vast majority of urban transportation systems in North America are equipped with traffic surveillance systems that provide real time traffic information to traffic management centers. The information from these are processed and provided back to the travelers in real time. However, the travelers are interested to know not only the current traffic information, but also the future traffic conditions predicted based on the real time data. These predicted values inform the drivers on what they can expect when they make the trip. Travel time is one of the most popular variables which the users are interested to know. Travelers make decisions to bypass congested segments of the network, to change departure time or destination etc., based on this information. Hence it is important that the predicted values be as accurate as possible. A number of different forecasting methods have been proposed for travel time forecasting including historic method, real-time method, time series analysis, and artificial neural networks (ANN). This paper examines the use of a machine learning technique, namely support vector machines (SVM), for the short-term prediction of travel time. 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. A comparison of the performance of SVM with ANN, real time, and historic approach is carried out. Data from the TransGuide Traffic Management center in San Antonio, Texas, USA is used for the analysis. From the results it was found that SVM is a viable alternative for short-term prediction problems when the amount of data is less or noisy in nature.

Index Terms—Inductive loop detectors, Machine learning techniques, Support vector machines, Travel time prediction.

I. INTRODUCTION

TRAVEL time is a fundamental measure in transportation engineering that can be understood and communicated

Manuscript received December 1, 2006.

Lelitha Vanajakshi (corresponding author) was with Texas A&M University, College station, Texas, 77840, USA. The present address is Asst. Professor, Dept. of Civil Engg., IIT Madras, Chennai 600 036, India. (Phone: 91 44 2257 4291; e-mail: lelitha@iitm.ac.in).

Laurence Rilett was with Texas A&M University, College station, Texas, 77840, USA. He is now the Keith W. Klaasmeyer Chair in Engineering and Technology, University of Nebraska, Lincoln, NE, 68588-0531. (e-mail: lrilett2@unl.edu). by a wide variety of audience, including engineers, planners, administrators, and commuters. As a performance measure and decision-making variable, travel time is useful in many aspects of transportation planning, modeling, and decisionmaking applications. These applications include traffic performance monitoring, congestion management, travel demand modeling and forecasting, traffic simulation, air quality analysis, evaluation of travel demand, and traffic operations strategies. Travel time information is becoming increasingly important for a variety of real-time Some of the real-time transportation applications. applications include Advanced Traveler Information Systems (ATIS), Route Guidance Systems (RGS), etc., which are part of the Intelligent Transportation Systems (ITS). Thus, providing travelers with accurate and timely information to allow them to make decisions regarding route selection is one of the important applications that use travel time information in recent times.

Travel time data can be collected from the field either using direct methodologies such as probe vehicles, automatic vehicle identifiers (AVI) etc. or using point detector data, such as that obtained from inductance loop detectors (ILD). While less accurate, the latter is more popular because of the widespread deployment of point detectors in ITS applications across North America. Since ILD cannot measure the travel time directly, there is a necessity to estimate the travel time. There are different methods available to calculate travel time from loop detector data, the most popular among them being extrapolation of the point speed values. However, the accuracy of these speed-based methods reduces as the vehicle flow becomes larger. Other widely reported methods include statistical models and models based on the traffic flow theory, the majority of which are developed for either free-flow or congested-flow condition only.

Travelers, in general, will be interested in knowing the future travel time as it would inform them on what they could expect to encounter when they make the trip rather than what is happening right now or in the past. Intuitively, the performance of the application will be constrained if the current or historic traffic values are used, because by the time the user makes the trip the situation would have changed. Thus, there is a need for a methodology that will anticipate the values in the next few minutes and inform the travelers accordingly. There were different studies in the past on travel time prediction and different methods that can be used for this. Some of the more important methods include historical and real time algorithms [1], [2], regression, time-series and Kalman filtering models [3]-[8], and Artificial Neural Network (ANN) models [9], [10].

The objective of this paper is to investigate the potential of Support Vector Machines (SVM) technique for the shortterm prediction of travel time. The performance of SVM method is compared with ANN, historic and real time methods. The analysis considered forecasts ranging from a few minutes ahead up to an hour into the future. Inductance loop detector data obtained from the TRANSGUIDE Traffic Management Center in San Antonio, Texas, USA was used to estimate the travel times and this information was input to the different prediction techniques.

In the following sections, a brief discussion of the historic method, real time method, ANN and SVM methods will be given followed by the details of the data used for the analysis. The implementation details and the results are discussed subsequently.

II. METHODS FOR TRAVEL TIME PREDICTION

Historic and real-time approaches are the two popular methods adopted in the field for forecasting travel times. Different studies in the past explored the use of other techniques such as regression, time-series and Kalman filtering models and ANN models for travel time prediction. Out of these, ANN has been extensively investigated by many researchers for the prediction of traffic parameters, and hence this study compares the performance of SVM with ANN in addition to real time and historic methods. A brief discussion of each of these is given in the following subsections.

The historic approach [1] is based on the assumption that the historic travel time profile can represent the traffic characteristics for a given time of the day. Thus, a historical average value will be used for predicting future values. An important component of the historic approach is the classification of many days of data into "representative" days that have similar profiles. This method can be valuable in the development of prediction models because they explain a substantial amount of the variation in traffic over days. However, for the same reason, the reliability of the prediction is limited because of its implicit assumption that the projection ratio remains constant.

However, commuters, in general, have a good idea about the average traffic conditions they will face on any given day. Consequently, they will be interested in conditions that are abnormal, that is, when the average values are not representative of the current or future traffic conditions. In the real time approach [2], it is assumed that the most recently obtained or estimated travel time information will hold into the future. This method can perform reasonably well for the prediction into immediate near future under traffic flow conditions without much variation [11]. Artificial Neural Network (ANN) with its learning capabilities [12] has been investigated by many researchers in the field of transportation engineering for the prediction of traffic parameters [9], [10]. In particular, multi-layer feed forward neural networks that utilize a back propagation algorithm have been applied successfully for forecasting traffic parameters [10],[13],[14]. Hence, in this study a multi-layer feed forward neural network with back propagation algorithm is used for comparison purposes. The back propagation algorithm neural network was coded in MATLAB. One hidden layer with 10 neurons was found to be the optimum.

SVM has been successfully applied to a number of applications ranging from particle identification to database marketing [15], [16]. In the present study, Support Vector Regression (SVR) was selected for the prediction of travel time. In SVR, the basic idea is to map the data into a highdimensional feature space via a non-linear mapping and do linear regression in this space. Thus, linear regression in a high dimensional (feature) space corresponds to non-linear regression in the low dimensional input space [15]. Further details of this technique can be seen in [17], [18] and [19]. The SVR model used a radial basis kernel function. SVM toolbox for MATLAB developed by Gunn was used for the present study [20].

III. DATA COLLECTION

The data for this study were collected from the TRANSGUIDE Traffic Management Center in San Antonio, Texas, USA [21]. The I-35 North freeway was selected, which is equipped with dual loop detectors at 0.5-mile intervals. The data are reported in 20-second intervals and includes flow, occupancy and speed. Also, this section was equipped with automatic vehicle identifiers, the data from which was utilized for validating the results in subsequent research [22].

The data were analyzed over a continuous 24-hour period for five consecutive days starting from February 10th Monday to February 14th Friday of 2003. The travel time from link 1 between stations 159.500 and 159.998 on all the five days was analyzed first. After extracting the specific detector data from the whole set, an extensive data reduction and quality control process were carried out to identify and correct any discrepancies in the data. The quality control and analysis included checking for missing data and threshold checking on the speed, volume and occupancy observations, individually as well as in combination. After the data reduction and quality control, the data were aggregated into 2-minute intervals in order to reduce the quantum of data to be processed while still capturing most of the trends in the varying traffic conditions. Thus, an original data file for a 24-hour time period having 4300 records was reduced to 720 records after aggregation. The next step was the estimation of the corresponding travel time from this data. This estimation can be carried out using a variety of different techniques available such as by dividing the known distance by the speed reported by the detectors, or by using methods based on statistical or traffic flow theory based models.

A travel time estimation procedure was developed in the present study which can be used for both peak, off-peak, and transition period traffic flow conditions. The methodology is based on the characteristics of the stochastic vehicle counting process and the principle of conservation of vehicles. The model estimates speed and travel time as a function of time directly from flow measurements. The methodology is based on the traffic flow theory and uses flow, occupancy, and speed data from the detectors as input. The details of the methodology can be found in [22], [23] and is not detailed here since it is beyond the scope of the present paper. This estimated travel time was used as input for the prediction problem in this study. However, this prediction methodology would work equally well with travel time collected directly or estimated using any other methods also.

Four days data was used for training and one day data was used for testing the results. The training was carried out based on data from February 10 to 13th 2003 (Monday to Thursday). The data from February 14th Friday was kept for validation. Fig. 1 shows the travel time distribution for all the five days as a function of time. It can be seen that on Tuesday, February 11, 2003, the data is showing lesser magnitude throughout the day compared to all the other days. Also, it can be seen that on Wednesday, February 12 the peak value in the travel time is small compared to other days. The mean absolute difference (MAD), as given in (1), was calculated between each day's data with Friday data, which is the value to be predicted. The MAD came to be 3.85, 7.85, 4.87, and 3.99 for Monday, Tuesday, Wednesday, and Thursday data respectively.

$$MAD = \sum \frac{|actual - predicted|}{N} \tag{1}$$

Thus, one week data was used for the analysis with Friday data for testing and Monday to Thursday data for training. It is well known that predicting Friday data from previous Fridays' data will be a better option since the pattern of input and output data will be more similar. If such a data is used the results will be better due to more similar input output pattern. However, the focus of this investigation is to explore the possible use of SVM as a technique for predicting travel time rather than the best possible data set to be used as input for prediction. Thus this study checks the performance in a worst case scenario. If the results in this study are encouraging, it can be used with different data set such as same day's data from different weeks where the data has less variation.

IV. RESULTS

Travel time was predicted into future time steps using

historic method, real time method, ANN method and SVM method and the results are compared. The analysis considered prediction times ranging from 2-minute ahead up to an hour ahead.

First, the 2-minute aggregated data was normalized based on the range of the travel time values. The input and output data was selected as the travel time for the 5 previous time step values and the travel time for the next time step value respectively. Because the data was grouped in 2-minute intervals, five time steps correspond to a 10-minute interval. Thus, the prediction was based on the previous 10-minute travel time values. The model then predicts the next 2minute travel time as shown in (2).

 $T(k+\Delta t) = f(T(k-4 \Delta t), T(k-3 \Delta t), T(k-2 \Delta t), T(k-\Delta t), T(k))$ (2) where,

 $\Delta t = time interval,$

T = travel time, and

k = current time interval.

The prediction was subsequently carried out to 4 minutes, 6 minutes, *etc.*, up to an hour ahead. The results of 2-minute ahead prediction using all the four methods are detailed first. The results obtained using the historic method, which assumes that the historic average represent the future travel time, is shown in Fig. 2. The travel time of Friday, predicted using the other four days data, based on the historic method and the corresponding actual travel time for the 24 hour period are plotted. The MAPE, as given in (3), is calculated between the predicted travel time and the actual travel time for the 24 hour period and came to be 9.36%.

$$MAPE = \frac{\sum \frac{|actual - estimated|}{actual}}{Number of observations} \times 100$$
(3)

Fig. 3 shows the same 2-minute ahead predicted travel time using the real time method, which assumes that the current travel time is going to continue to the future time step. As expected, the predicted travel time leads the actual travel time by the 2-minutes prediction interval. The corresponding MAPE for the whole 24 hour period came to be 9.66%.

Fig. 4 and 5 show the predicted travel time using ANN method and SVM method. It can be seen that the travel time predicted by both SVM and ANN were able to follow the trends in the actual data better than the historic and real time methods. The MAPE values of ANN and SVM were 8.64% and 7.38% respectively.

An enlarged view of the actual travel time values and the corresponding predicted values for a 2-hour evening peak and off-peak for the 2-minute ahead prediction using all the four methods is shown in Fig. 6 for illustration. This figure clearly illustrates the historic method performing very poorly for the prediction of peak period. Also it can be seen that in the case of the real time method the predicted travel time leads the actual travel time by the 2-minutes prediction

interval. And the SVM and ANN following the trends in the actual travel time can also be seen.

The next step was to extend the prediction further ahead to 4-minutes, 6-minutes etc. up to an hour ahead into future. The result obtained in each prediction is compared with the actual value and the MAPE value is calculated. The prediction was carried out for the full 24-hour data. The training data was varied from one day to four day data and each result is detailed below to show the effect of the quantity and quality of data on the prediction result.

Fig. 7 shows the error in prediction when a single day data (Monday) was used for training the network and Friday travel time was predicted. MAPE values are shown for 2minute ahead prediction up to one hour ahead prediction. The MAPE for the historical method, the real time method, the ANN and SVM methods are shown in this figure. It can be seen that the historic method outperformed real time method through out the prediction. SVM performed better than historic only up to 6 minutes of prediction and ANN performed better up to 10 minutes of prediction ahead. Thus, historic method outperformed other methods in this case after 10 minutes of prediction time ahead, which can be explained based on Fig. 2. In Fig. 2, both the training data (Monday) and testing data (Friday) had the same pattern with an MAD of 3.85, which makes historical method the best method for prediction. It can also be observed that ANN performed better than the SVM in this case.

Fig. 8 shows the MAPE values when 2 days data were used for training (Monday and Tuesday) and when the 24hour Friday data was predicted. Comparing Fig. 8 with Fig. 7, it is seen that there is an increase in prediction error using historic method from 9.3% to 14.8% when the training data was changed from Monday data alone to Monday and Tuesday data together. This is due to the fact that the Tuesday travel time data differed in magnitude when compared to Monday and Friday data. It can be seen that the Monday and Friday data having very similar trends throughout with an MAD of 3.85, whereas the MAD between Tuesday and Friday data came to be 7.84. This difference of Tuesday data makes the training data different from testing data, reducing the performance of historic method. The reduced performance of ANN also is due to the same reason. The SVM method out-performed all other methods in this case and that historic method failed throughout the one hour prediction period compared to other methods.

Fig. 9 shows similar result when 3 days data was used for training (Monday, Tuesday and Wednesday) and the Friday data was predicted. It can be seen that with more data being added to the training set, the effect of Tuesday data is getting reduced. As in the previous case, here also the SVM performed better than all the other methods. Up to 35 minutes of prediction ahead, the other methods performed better than historic method.

Fig. 10 shows similar result when 4 days data was used

for training (Monday, Tuesday, Wednesday and Thursday) and the Friday data was predicted. In this case the historic method outperformed other methods after about 30 minutes of prediction ahead. SVM performed better than ANN and real time methods throughout the prediction.

It can be seen that, as more and more data is added to the training set, the influence of the Tuesday data reduces and this is reflected in the reduction in error for the historic and ANN methods. A comparison of SVM and ANN shows that their performances are comparable to each other with SVM having slight advantage over ANN in this case. Both ANN and SVM performed better than real time and historic method showing the ability of these methods to capture the variability in travel time in a better way.

Theoretically the accuracy of the SVM prediction does not depend on the amount of data used once the support vectors are selected. This performance of SVM can be explained based on the inherent nature of the SVM training process. Once SVM chooses the data points, which can represent the input data (support vectors), its performance is more or less independent of the amount of training data. Hence, if the support vectors selected from the training data are not affected, its performance may not get affected by the amount of training data. However, in the case of ANN, the network can learn more about the data as the amount of training data increases and this change the results for the better.

Hence, in scenarios where the training data has variations and the availability of data is limited, SVM will be a better choice than ANN. In this study also, the results agreed with this and showed SVM performing better when the data was having more variations between the training set and testing set. In cases where large amount of data is available and the data do not have much variation, ANN was found to be an equally good predictive algorithm as SVM. Overall based on the results obtained for the particular data set under analysis, one can conclude that SVM is a better choice for travel time prediction problems if the amount of training data is less, or when the training data has lot of variations. Analysis of data on other links and other days showed similar results (22).

V. SUMMARY AND CONCLUSIONS

This paper investigated the usefulness of SVM for the short term prediction of travel time. A comparison was carried out between the performance of SVM and other popular methods such as ANN, historic method and real time method. The ANN model used is a multi-layer feed forward neural network and the SVM model used was a support vector regression with radial basis kernel function. The analysis considered forecasts ranging from 2 minutes ahead up to an hour into the future. Up to four days data was used as training set, with the previous 10 minutes data as input and next 2 minute data as output. One full day data was left for cross validation to evaluate the prediction errors.

The results of this study showed that current traffic

conditions are good predictors while long-range predictions need the use of historical data. The ANN and SVM methods performed better for some range into future (up to around 30 minutes ahead). Also, both these methods have good dynamic response and show better performance compared to the traditional models.

Comparison between ANN and SVM showed that, the performance of both SVM and ANN are comparable to each other. SVM becomes a better choice for the short-term prediction of travel time, if the amount of training data is less, or when the training data has more variations compared with the testing data. Also, it was found that the influence of the amount of training data used is more on the ANN method than on the SVM method. Overall, it was found that SVR is a viable alternative to ANN for short-term prediction problems when the amount of data is less or when the training data was not a good representative sample of the testing data.

REFERENCES

- C. Hoffman, and J. Janko, "Travel Time As A Basis Of The LISB Guidance Strategy", In Proc. Of IEEE Road Traffic Control Conference, IEEE, New York, pp. 6-10, 1990.
- [2] P. Thakuriah, A. Sen, J. Li, N. Liu, F. S. Koppelman, and C. Bhat, "Data Needs For Short Term Link Travel Time Prediction", *Advance working paper series number 19*, Urban Transportation Center, University of Illinois, Chicago, 1992.
- [3] J. Rice, and E. van Zwet, "A Simple and Effective Method for Predicting Travel Times on Freeways", *IEEE intelligent Transportation systems conf. Proc.*, pp. 227-232, 2001.
- [4] J. Kwon, B. Coifman, and P. Bickel, "Day-To-Day Travel Time Trends And Travel Time Prediction From Loop Detector Data", *Transportation Research Record: Journal of the Transportation Research Board, No. 1717*, TRB, National Research Council, Washington, D.C., pp. 120-129, 2000.
- [5] A. Sen, N. Liu, P. Thakuriah, and J. Li, "Short-Term Forecasting Of Link Travel Times: A Preliminary Proposal", *ADVANCE Working Paper Series*, Number 7, 1991.
- [6] M. Saito, and T. Watanabe, "Prediction and Dissemination System for Travel Time Utilizing Vehicle Detectors", *Proc. of the 2nd world congress on Intelligent Transp. Systems*, Yokohama, Japan, 1995.
- [7] S. Chien, X. Liu, and K. Ozbay, "Predicting Travel Times For The South Jersey Real-Time Motorist Information System", CD-ROM. Transportation Research Board, National Research Council, Washington, D.C., 2003.
- [8] C. Nanthawichit, T. Nakatsuji, and H. Suzuki, "Application Of Probe Vehicle Data For Real Time Traffic State Estimation And Short Term Travel Time Prediction On A Freeway", CD-ROM. Transportation Research Board, National Research Council, Washington, D.C., 2003.
- [9] J. W. C. Van Lint, S. P. Hoogendoorn, and H. J. van Zuylen, "Freeway Travel Time Prediction with State Space Neural Networks", CD-ROM. Transportation Research Board, National Research Council, Washington, D.C., 2002.
- [10] D. Park, and L. R. Rilett, "Forecasting Freeway Link Ravel Times With A Multi-Layer Feed Forward Neural Network", *Computer Aided Civil And Infra Structure Engineering*, vol. 14, pp. 358 – 367, 1999.
- [11] S. Shbaklo, C. Bhat, F. Koppelman, J. Li, P. Thakuriah, A. Sen, and N. Rouphail, "Short-Term Travel Time Prediction", *ADVANCE Project Rep., TRF-TT-01*, Illinois University Transp. Research Consortium, 1992.
- [12] S. Haykin, Neural Networks: A Comprehensive Foundation, Prentice Hall, 1999.
- [13] J. Mc Fadden, W. T. Yang, and S. R. Durrans, "Application Of Artificial Neural Networks To Predict Speeds On Two-Lane Rural Highways", CD-ROM. Transportation Research Board, National Research Council, Washington, D.C., 2001.

- [14] S. H. Huang, and B. Ran, "An application of neural network on traffic speed prediction under adverse weather condition", CD-ROM. Transportation Research Board, National Research Council, Washington, D.C., 2003.
- [15] V. Kecman, Learning and Soft Computing: Support Vector Machines, Neural Networks, And Fuzzy Logic Models, The MIT press, Cambridge, Massachusetts, London, England, 2001.
- [16] C. Campbell, "Kernel Methods: A Survey Of Current Techniques", *Neuro Computing*, Vol. 48, 2002, pp. 63-84.
- [17] V. Vapnik, *The Nature of Statistical Learning Theory*, Springer Verlag, New York, 1995.
- [18] N. Cristianini, and J. Shawe-Taylor, An Introduction to Support Vector Machines, Cambridge University Press, Cambridge, UK, 1999.
- [19] A.J. Smola, and B. Schölkopf, "A tutorial on Support Vector Regression", Technical Report NC2-TR-1998-030, NeuroCOLT2, 1998.
- [20] S. R. Gunn (Nov. 23, 2006), "Support Vector Machines for Classification and Regression, Available: <u>http://www.ecs.soton.ac.uk/~srg/ publications/pdf/SVM.pdf.</u>
- [21] Transguide Model Deployment Initiative Design Report and transguide Technical Paper (Nov. 1, 2006), Texas Department of Transportation, Available: <u>http://www.transguide.dot.state.tx.us/</u> PublicInfo/papers.php.
- [22] L. Vanajakshi, "Estimation and Prediction Of Travel Time From Loop Detector Data For Intelligent Transportation Systems Applications", Ph.D. Dissertation, Dept. of Civil Eng., Texas A&M University, 2004.
- [23] L. Vanajakshi, and L. R. Rilett, "Travel Time Estimation from Loop Detector Data", *ITS Safety and Security Conference*, Florida, 2004.



Fig. 1 Travel time distribution on the study dates.



Fig. 2 Travel time predicted by historic method.







Fig. 4 Travel time predicted by ANN method.



Fig. 5 Travel time predicted by SVM method.



Fig. 6 Comparison of the predicted values during peak period using different methods for a 2minute ahead prediction.



Fig. 7 MAPE for prediction using one-day data for training.



Fig. 8 MAPE for prediction using two-day data for training.



Fig. 9 MAPE for prediction using three-day data for training.



Fig.10 MAPE for prediction using four-day data for training.