

Review Paper for science article:  
**“A Simple and Effective Method for Predicting Travel Times on  
Freeways”**

Student name: Le Xu

## **Introduction**

In the traffic science article that published in the year 2004, both authors, Professor John Rice and Dr. Erik van Zwet from UC Berkeley, focus their attention on studying the travel time prediction. By applying real-world traffic data, and comparing three different kinds of methods, the authors discovered an empirical fact that the Linear Regression could be the simple and effective method for the prediction of future travel time. The predication can be done with the combination of the current travel time and historical data.

## **Data Sources**

The data that was being used in the analysis is from Performance Measurement System (PeMS) data, more specifically, it is the traffic flow and occupancy data from 116 single loop detectors along 48 miles of I-10 East in Los Angeles. The velocity was calculated based on the formula:

$$\text{velocity} = g \times \frac{\text{flow}}{\text{occupancy}}.$$

The authors are very carefully when defining the  $g$  value, which here is the unknown average length of the vehicles. The value of given  $g$  in the equation varies during the different time of the day and varies based on the different sensitivities of individual loop sensor.

## **Methods**

The methods of prediction in this paper are Linear Regression, Principal Components, and Nearest Neighbors. From those three methods, the main discovery of this paper is the empirical fact that there is a linear relationship between the current travel time and any future travel time when the starting point is the same.

At the very beginning paper, below formula was introduced by the Authors and it is to calculate the travel time from two loops, from a to b, which can be written as below (1):

$$T_d^*(a, b, t) = \sum_{i=a}^{b-1} \frac{2d_i}{V(d, i, t) + V(d, i + 1, t)}, \quad (1)$$

The assumption of this formula is that the  $T^*$  would be possible if there are no significant changes in traffic from loop point a to ending point b.

Historical mean travel time is also one of the most important components throughout this study. The reason is not only because we can compare the prediction to the historical data and the current data, but also due to that study the past data is the fundamental way to apply the Principal Components and Nearest Neighbors methods. As a result, historical mean travel time formula was brought out in the early beginning of the paper (2):

$$\mu_{TT}(t) = \frac{1}{|D|} \sum_{d \in D} TT_d(t). \quad (2)$$

Where  $t$  is the current time,  $D$  as the number of days, and  $TT$  represents the future travel time.

The review paper will walk you through those three learning models; however, the main focus would be on Linear Regression.

### ***1. Linear Regression:***

Firstly, the linear model was introduced as (3):

$$TT(t + \delta) = \alpha(t, \delta)T^*(t) + \beta(t, \delta) + \varepsilon. \quad (3)$$

The linear relationships between  $T^*$  (current status travel time) and  $TT$  (future travel time) with the given non-negative time lag  $\delta$  are formulated. By applying the data into the model, the problem became to solve the weighted function (weighted least squares), which is to minimize the below formula (4):

$$\sum_{\substack{d \in D \\ s \in T}} (TT_d(s) - \alpha(t, \delta) + \beta(t, \delta)T_d^*(t))^2 K(t + \delta - s), \quad (4)$$

by defining the value of  $\alpha(t, \delta)$ ,  $\beta(t, \delta)$ . The  $K$  is the Gaussian Density that has zero mean and customized variance that needs to be set specifically. So after the transformation, we are actually calculating the  $TT(\tau + \delta)$  as (5):

$$\widehat{TT}_e^{\alpha\beta}(\tau + \delta) = \hat{\alpha}(\tau, \delta)T_e^*(\tau) + \hat{\beta}(\tau, \delta). \quad (5)$$

the equation can be further transformed when we make  $\beta(\tau, \delta) = \beta'(\tau, \delta) \mu_{TT}(\tau, \delta)$ . Eventually, the equation represents a future travel time as the linear combination of the current travel time and historical mean travel time.

## 2. Principal Components

As a non-parametric method, the assumption of PC is that the travel times are independently and identically distributed on different days, and by a given day, the  $T^*$  and  $TT$  are multivariate normal. The mathematics that behinds the PC method is to find a few of the largest eigenvalues of the empirical covariance matrix of  $T^*$  and  $TT$  for a number of days, and the user of this method has to specify the number of the eigenvalues are retained.

## 3. Nearest Neighbors

Compare to the Principal Component method, Nearest Neighbor takes fewer assumption. The name of the method is quite self-explanatory. The nearest neighbor is aiming to find a day in the past that is most close or similar to the present day in some degree. Variable  $M$  was introduced as the distance between days. In order to minimize the distance  $M$ , the authors proposed two ways to calculate  $M$ , one is the speed difference, another is the travel time difference.

## Results

In the study, the authors calculated the estimated prediction error of the root mean squared (RMS) from the linear regression, with a time lag of 0 and 60 minutes respectively. It was satisfactory to the authors that the prediction on current travel time performs well, and the linear regression method shows the best result out of all three methods.

Up to the time lag of 60 minutes, the RMS error of the regression predictor stays within 10 minutes. In the authors' point of view, the result overall is rather impressive given the data from the rush hour. The Fig 7 below shows the RMS error for all three prediction methods with a time lag of 60 minutes.

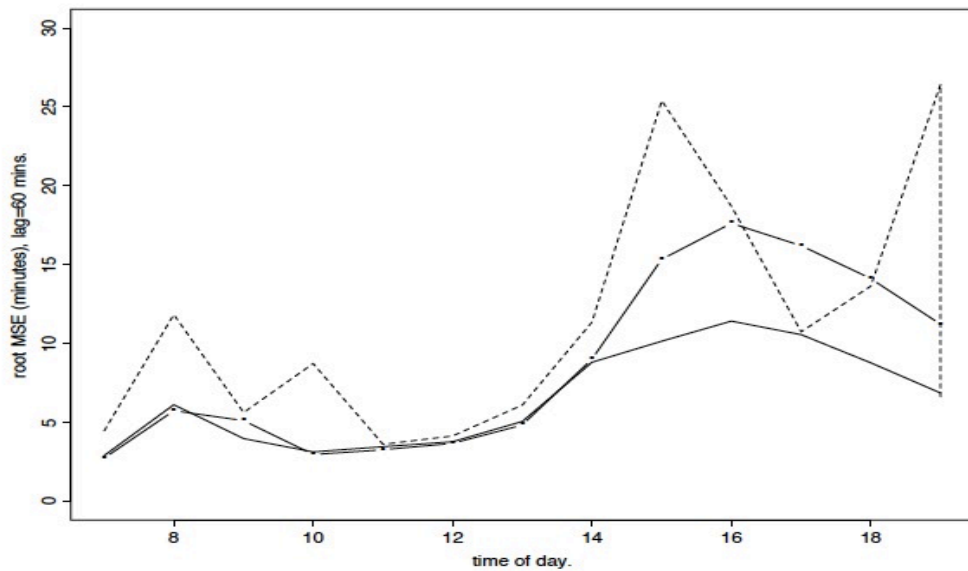


Fig. 7. Estimated RMSE, lag=60 minutes. Historical mean (— · —), current status (— —) and linear regression (—).

From the graph, we can also tell, as the lag increases more than 15 minutes, the historical mean starts to change abruptly and it is evident that linear regression is the relative stable method for prediction error.

## More Thoughts

There are three more worth-mention findings that I would like to point out. The first one is, based on the authors that the prediction method could be performed in a real time manor. Since the calculation of the smoothing variables could be expensive, it can be done offline. The second one is if we were asked to calculate travel time within a network, it is possible to calculate the travel time separately of each edge and united them together. The last thing is if we were given the arrival time, by regressing the arrival time on  $T^*$ , we can also predict the travel time based on arrival time, and then we can answer the question such as when to departure in order to arrive at the certain time frame.

To further study the methods and predication models on predication travel time, I would like to approach the study in two directions, the first one is I will use the same data to apply into other predication methods such as Support Vector Machine or other black box machine learning techniques to compare the resulted RMSE, the other one is I will use different dataset to see if the previously studied methods still holds universally.