

Forecasting Urban Travel Times in London

by Lucille Ablett

Submission date: 24-Mar-2017 02:42PM (UTC+0000)

Submission ID: 69354213

File name:

477659_Lucille_Ablett_Forecasting_Urban_Travel_Times_in_London_2799583_612364206.pdf (2.8M)

Word count: 8845

Character count: 45791

Forecasting of Urban Travel Times in London

Lucille Ablett and Christopher Baxter

24th March 2017

Word Count: 4987

Contents

1. Introduction	2
1.1. Data Description	2
2. Exploratory Spatio-Temporal Data Analysis	6
2.1. Temporal Characteristics	6
2.2. Spatio-Temporal Characteristics	8
2.3. Temporal Autocorrelation	10
2.4. Spatial Autocorrelation	14
3. Space-Time Forecasting	17
3.1. Autoregressive Integrated Moving Average (ARIMA) – Christopher Baxter.....	17
3.1.1. Method	17
3.1.2. Experimental Set-Up.....	17
3.1.3. Results	18
3.2. Support Vector Regression (SVR) – Lucille Ablett	23
3.2.1. Method and Experimental Set-Up	23
3.2.2. Results	24
4. Discussion	32
5. Conclusions and Further Work.....	33
References	34
Appendix A – Statistical Summary of Travel Times per Road Link	35
Appendix B – Support Vector Regression Models for All Links (Non-Spatial).....	36

Forecasting of Urban Travel Times in London

1. Introduction

This report covers the spatio-temporal analysis of urban travel times in London. The data analysed in this report is a subset of data collected for the London road network, as part of the London Congestion Analysis Project (LCAP) by Transport for London using automatic number plate recognition. The aim of this report is to explore the spatial and temporal variation in these travel times using various exploratory data analysis techniques, and to investigate two different forecasting methods: autoregressive integrated moving average (ARIMA) and support vector regression (SVR).

ARIMA is an autoregressive modelling method that uses linear regression to predict a value based on a series of one or more previous values. ARIMA incorporates a moving average model that uses the error in previously predicted values to improve future predictions. It relies on initial exploratory data analysis to identify temporal trends to be incorporated at the modelling stage (Cheng & Wang, 2011).

SVR, on the other hand, is a machine learning method developed from support vector machines (SVM). In addition to machine learning algorithms, SVR utilises kernel methods to map the data into a high dimensional feature space. This allows a linear algorithm to be used to solve what could potentially be a non-linear problem (Vapnik, 1999). Unlike artificial neural networks (ANN), SVR uses a convex quadratic optimisation algorithm and so cannot get stuck in local minima, therefore resulting in a globally optimal solution (Peng, et al., 2009). SVR also has a strong generalisation ability, making it good for handling noisy datasets (Zhang, et al., 2012). Since time series are commonly non-linear and prone to noise, SVR was selected as a forecasting method for this investigation.

SVR for time series analysis is relatively recent. In 2004, Vanajakshi & Rilett compared an SVR prediction model for traffic speed on freeways in San Antonio, Texas, to an equivalent ANN model, examining prediction accuracy for forecasts between 2 and 60 minutes. SVR was found to be comparable to ANN, with SVR performing better when the quality and quantity of the training data were poor. When compared to least squares and radial basis function networks for the prediction of chaotic time series, SVR was also found to generate more accurate predictions, particularly for longer term forecasts (Lau & Wu, 2008). In Taiwan, SVR was also found to significantly reduce the errors in predicted travel times, compared to previous methods, when utilised to examine travel times over three separate highway segments (Wu, et al., 2004).

1

1.1. Data Description

The dataset analysed comprises 8640 travel time readings for over 30 adjacent links passing through East Central London (Figure 2) between the hours of 06:00 and 21:00, from 1st – 30th January 2011. The data has been pre-processed into unit travel time, in seconds per meter, and aggregated to 5 minute intervals, resulting in 180 observations per day. These roads were selected for this investigation as they form two separate, yet intersecting, routes through London, and some additional adjoining roads.

A summary of the mean travel time data for all links analysed is given in Table 1. A full table of statistics for each individual road link is found in Appendix A.

Histograms and QQ plots (Figure 1) show a selection of distribution types observed within the dataset. The travel times for these links appear normally distributed when viewed using histograms, however QQ plots indicate significant deviation from normality, particularly at the tails where more extreme outliers exist in the dataset than would be expected for a normal distribution.

The range of travel times across the links is also variable. Whilst links 1877 and 446 show a relatively narrow spread of travel times, links 1604 and 2261 show a larger spread. This is observed in the large variation interquartile ranges between links, such as 0.047 and 0.150 for links 1877 and 1604 respectively.

For the Figure 1 links, the mean is located to the right of the distribution peak, indicating a positive skew. Links 1877 and 446 show moderate to strong positive skews of 5.60 and 7.43 respectively. In comparison, link 2261 shows a weak positive skew of 0.70. Although the peak of the distribution for link 1604 is clearly below that of the mean, the wide distribution and large number of slow travel times results in a relatively weak skew of 1.10.

A map of average travel times per road link (Figure 3), indicates that the slowest travel times are found along the road running north-south where interests the east-west road.

<i>Statistic</i>	<i>Value</i>
<i>Mean</i>	0.177
<i>Standard Deviation</i>	0.052
<i>Minimum Absolute Deviation</i>	0.037
<i>Median</i>	0.177
<i>Minimum</i>	0.083
<i>Maximum</i>	0.324
<i>Range</i>	0.241
<i>Skew</i>	0.579
<i>Kurtosis</i>	0.413

Table 1: Statistical summary of mean travel times per road link.

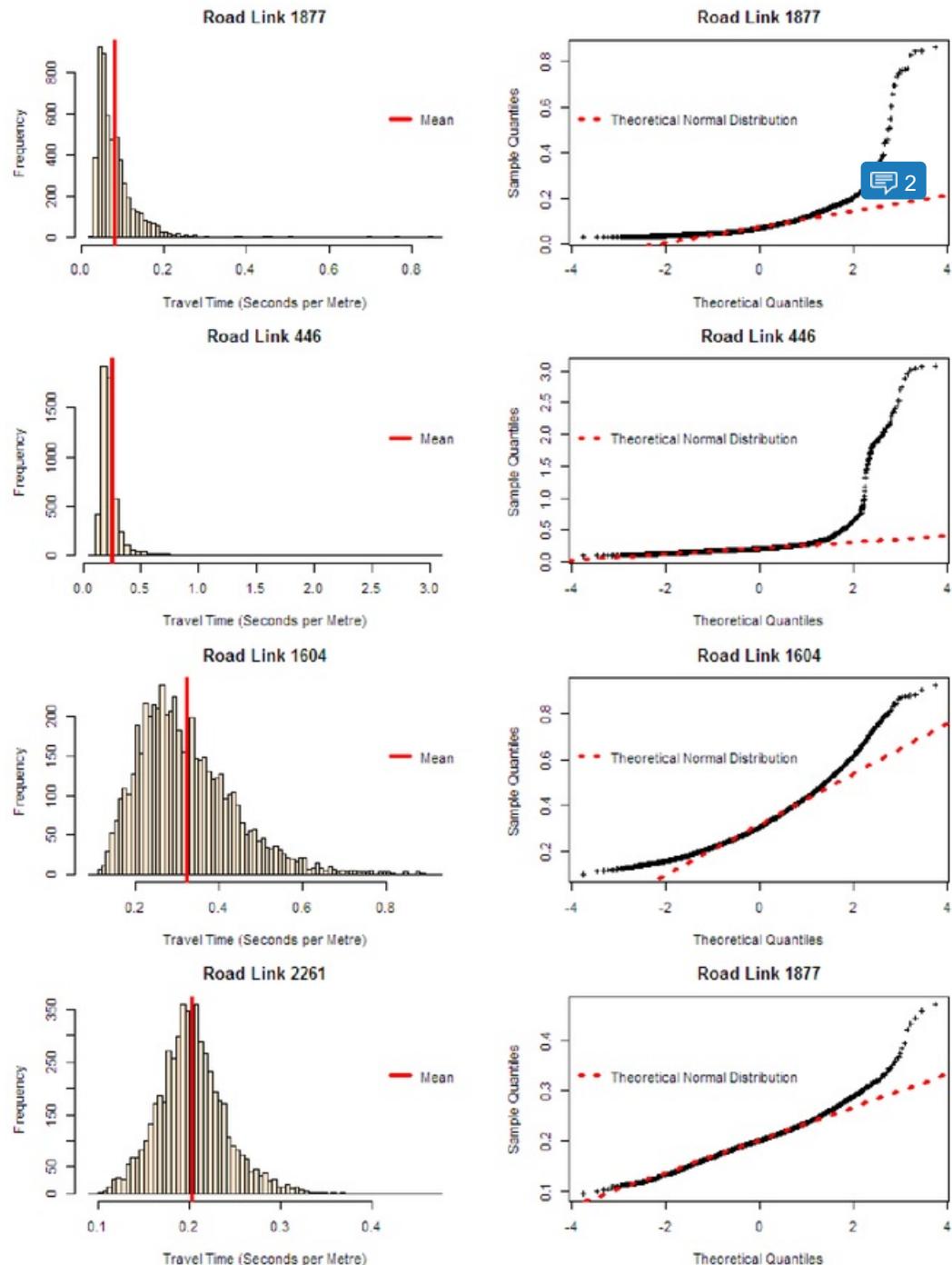


Figure 1: Histograms (left) and QQ plots (right) for four links showing the different distributions of travel times across the dataset. Although the histograms appear to show normal distribution, none of the travel time distributions shown above are truly normally distributed. This is evident from the QQ plots that indicate there is a strong case for non-normal distributions.

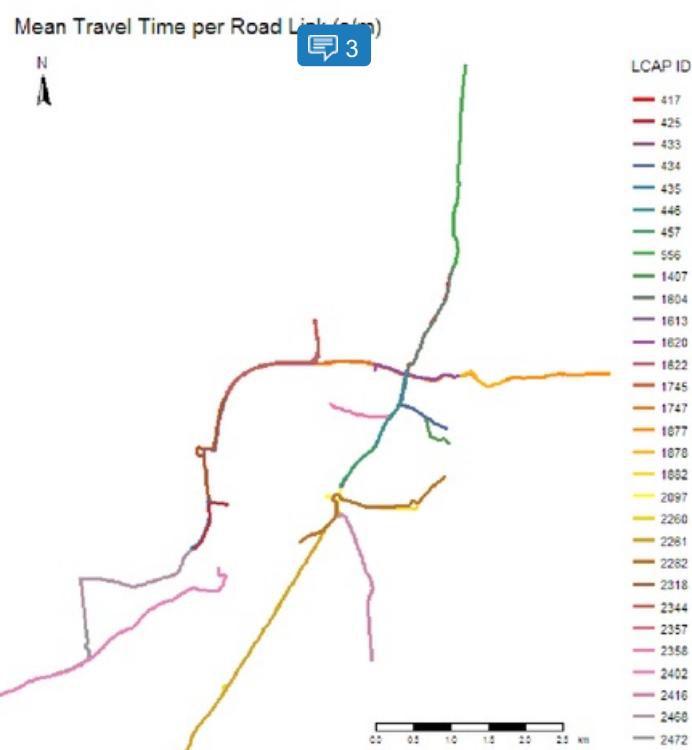


Figure 2: IDs for all links analysed in this report.

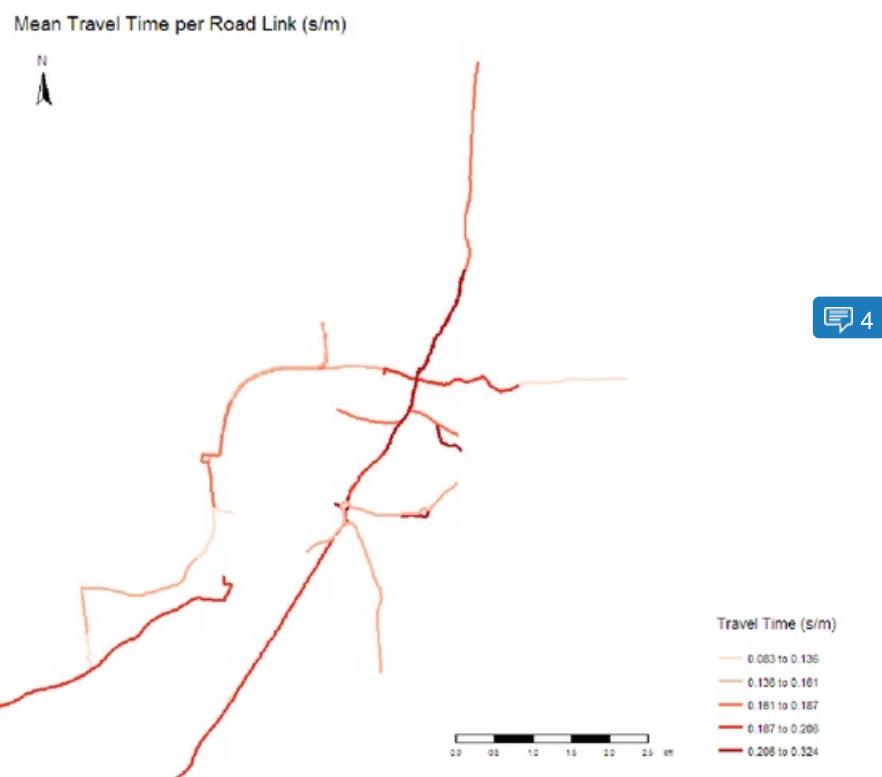


Figure 3: Mean travel times for all links in seconds per meter.

2. Exploratory Spatio-Temporal Data Analysis

2.1. Temporal Characteristics

Analysis of mean travel times across all links (Figure 4) shows a distinct pattern of high (slow) and low (fast) travel times, both daily and weekly. Travel times at the beginning of the month show less variation than in the second half of the month. It should be noted that Monday 3rd January 2011 was observed as a holiday. Figure 4 also shows clear distinctions in travel time trends between weekdays and weekends, as the slowest travel times are during the weekdays. During the weekends, the daily trends are still evident, but variation in travel time is significantly reduced.

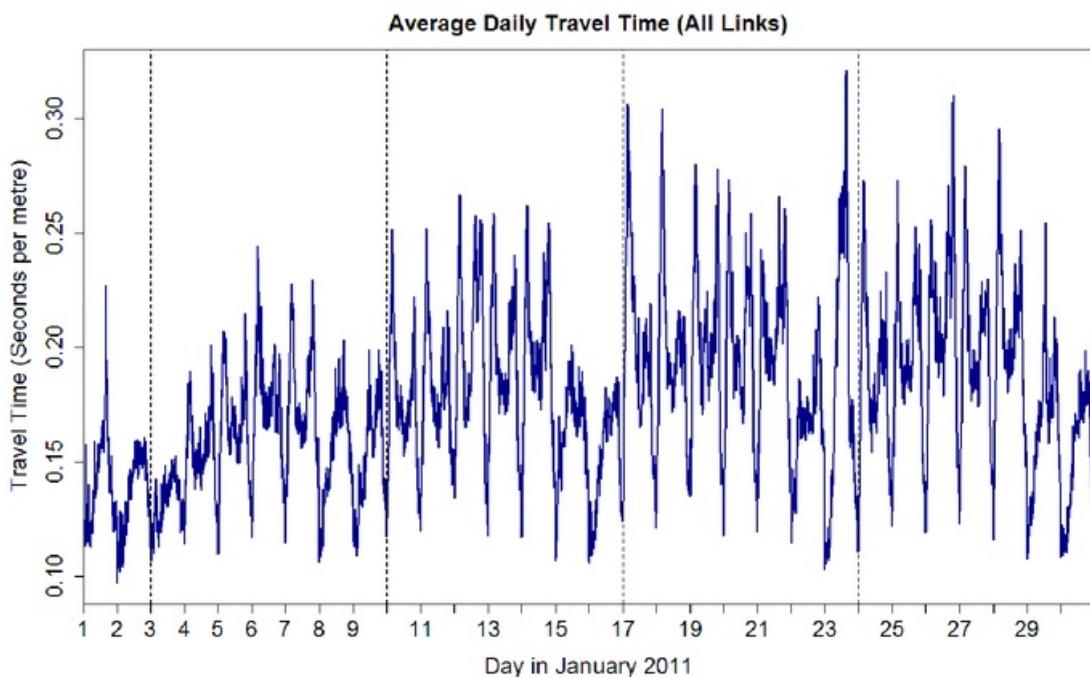


Figure 4: Mean travel times for all links at 5 minute intervals from 1st - 30th January 2011. Dotted lines show 00:00 on Monday.

Figure 5 shows the average travel time for each day of the week. For weekdays, there are two distinct peaks of longer travel times: between 08:00 and 09:00, and again between 18:00 and 19:00. The morning peak is consistently higher than the afternoon peak across all weekdays. This appears most pronounced on Thursdays and Fridays. A third, subtler, peak is observed around 16:00. These afternoon peaks likely show the different working schedules between those working shifts, and employees working normal business hours. The morning peak is likely higher as people working both shifts and normal business hours start work at the same time.

On weekends, the travel times show no clear peaks of longer travel times throughout the day, with Sunday morning between 06:00 and 08:00 showing the fastest travel time of the week at around 0.11 s/m.

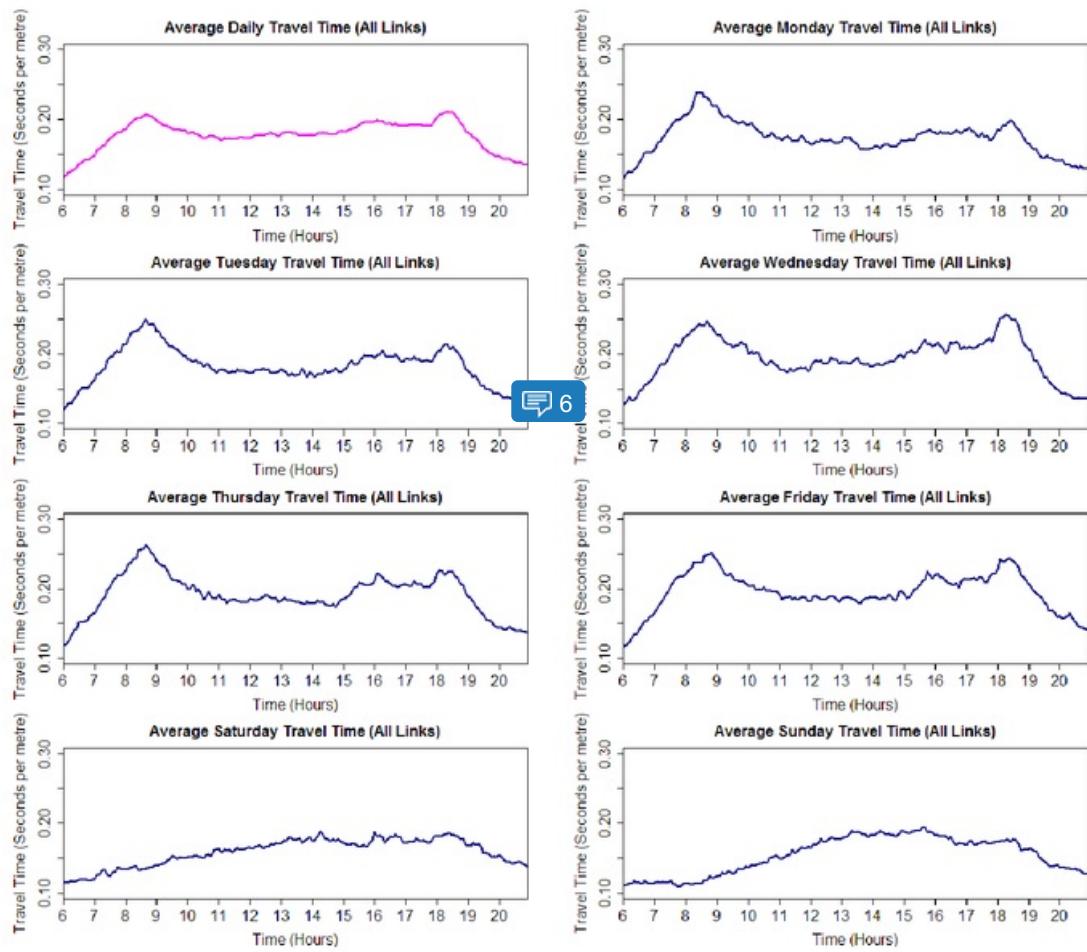


Figure 5: Average daily travel time for each day of the week at 5 minute intervals, averaged across all links. Weekdays and weekends show clearly different trends, with weekdays showing a morning peak around 08:00, and two afternoon peaks at 16:00 and 18:00.

Figure 6 shows the mean daily travel time profiles for all links. There is significant variation in travel time profiles across the links, with some roads showing little variation over the day. Links 434 and 2358, for example, show relatively consistent travel times throughout the day, suggesting that these roads are not particularly vulnerable to congestion at peak times. This is reflected in the low range of travel times of 0.45 and 0.40 s/m respectively. In contrast, link 1882 shows the greatest variation in travel time (0.98 s/m range), with significant increases in travel time seen during the morning and late afternoon rush hours. In the morning, the travel time is seen to increase rapidly to a peak of 0.5 s/m, yet this decreases relatively quickly throughout the morning. The evening peak is even more pronounced with travel time increasing by 0.13 s/m over the space of 20 minutes, yet returning to normal only 45 minutes later, indicating that this road segment may be particularly susceptible to small variations in traffic volume. It could be that the road has few lanes.

Some links show increases in travel times only during the morning or afternoon peak times. It is likely that these are inbound and outbound roads (or their adjoining roads), such as 2472 and 1620 respectively, since it is expected that most commuters would be heading into London during the morning for work, and returning home in the afternoon.

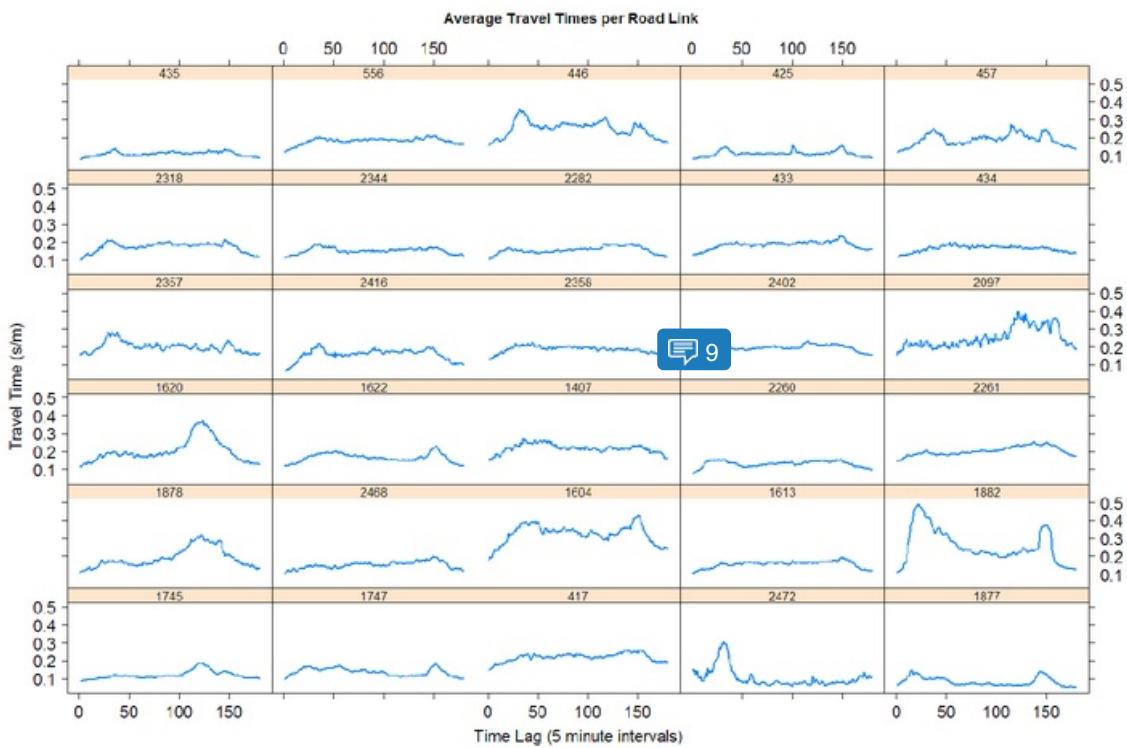


Figure 6: Time series plot showing average travel times (in seconds per meter) per road link between 1st and 30th January 2011. Travel times are plotted by time lag from 08:00 to 21:00, where 1 unit = 5 minutes.

2.2. Spatio-Temporal Characteristics

Spatial analysis of the travel times for each road link at 08:00, 11:00, 16:00, and 18:00 (Figure 7) show the links immediately north of the intersection of the two main roads have the slowest travel times during both peak and off-peak times (links 417 and 1604), with travel times of between 0.3 and 0.4 s/m recorded at all four time periods. The links forming the section of road directly south of this intersection (446 and 433) also show relatively high travel times of 0.25 – 0.3 s/m, although this is only observed at 08:00 and 11:00. The intersecting road only appears to experience travel times of over 0.3 s/m during early afternoon peak at 16:00, with travel times of only 0.15 – 0.2 s/m measured earlier in the day.

The two roundabouts located to the south of this intersection show slow travel times at all time of the day for drivers turning west (links 2097 and 1882). This could indicate an obstruction, such as roadworks, or traffic lights controlling the flow of the traffic.

To the south-west of the study area, link 2472 shows high travel times only at 08:00, and low travel times at all other times of the day. At this point, the road joins 2402, a major route out of central London. It is likely that this morning increase in travel times is due to drivers queuing to join the main road. The selection of links comprising the road through that follows the north bank of the Thame shows faster travel times than the previously mentioned slow travel time links, with the highest travel times of only 0.2 - 0.25 s/m seen at 16:00. People wanting to access this area of London may utilise public transport instead to avoid the congestion charge.

Viewing the average travel times per road link can be misleading, as it does not account for differences in speed limits between the roads, the number of lanes, or any factors contributing to high travel times. Figure 8 shows the differences in travel times during the previous hour, highlighting areas with abnormally high travel times.



Figure 7: Spatial variation in travel time at peak (08:00, 16:00, 18:00) and off-peak (11:00) times. It should be noted that some roads links are situated on top of each other and so some link attributes are obscured using this visual analysis method

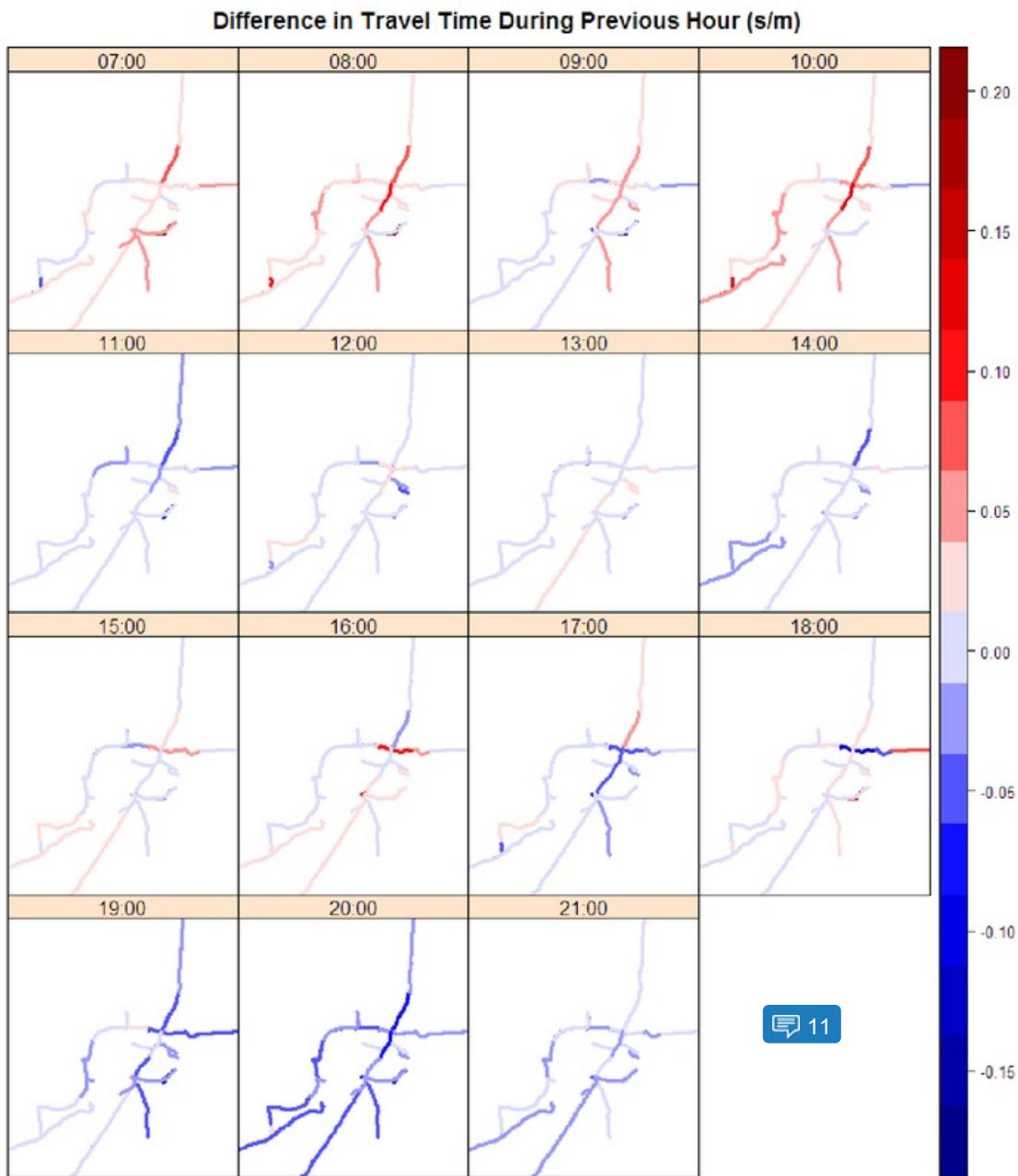


Figure 8: Difference in travel times for all road links during the previous hour (s/m). It should be noted that some roads links are situated on top of each other and so some link attributes are obscured using this visual analysis method.

2.3. Temporal Autocorrelation

To further examine the temporal characteristics of the dataset, the autocorrelation function in R (ACF) was used to determine the relationship between travel times at differing time lags. Correlation is measured on a scale of -1 to 1, where 0 indicates no correlation between the observations.

Figure 4 was seen to show both daily and weekly patterns in the mean travel time across all links. This is confirmed by Figure 9, where two separate correlation patterns over a period of a week (1st – 7th January) are seen. Strong positive correlation is observed at multiples of 180 lags, with moderately strong correlation occurring at 180 lag intervals starting at 100 lags. 100 and 180 lags refer to the times 06:00/21:00 and 14:20. Although the strength of these correlations initially decreases with time, the autocorrelation increases again towards the end of the week. This is likely indicating that weekends are more similar to each other than weekdays, since the 1st January was a Saturday. Apart from at time lags of 360 and 1620, all correlation, whether positive or negative, is significant to with 95% confidence levels.

Analysis of autocorrelation for individual links (Figure 10) shows this trend is common across most links.

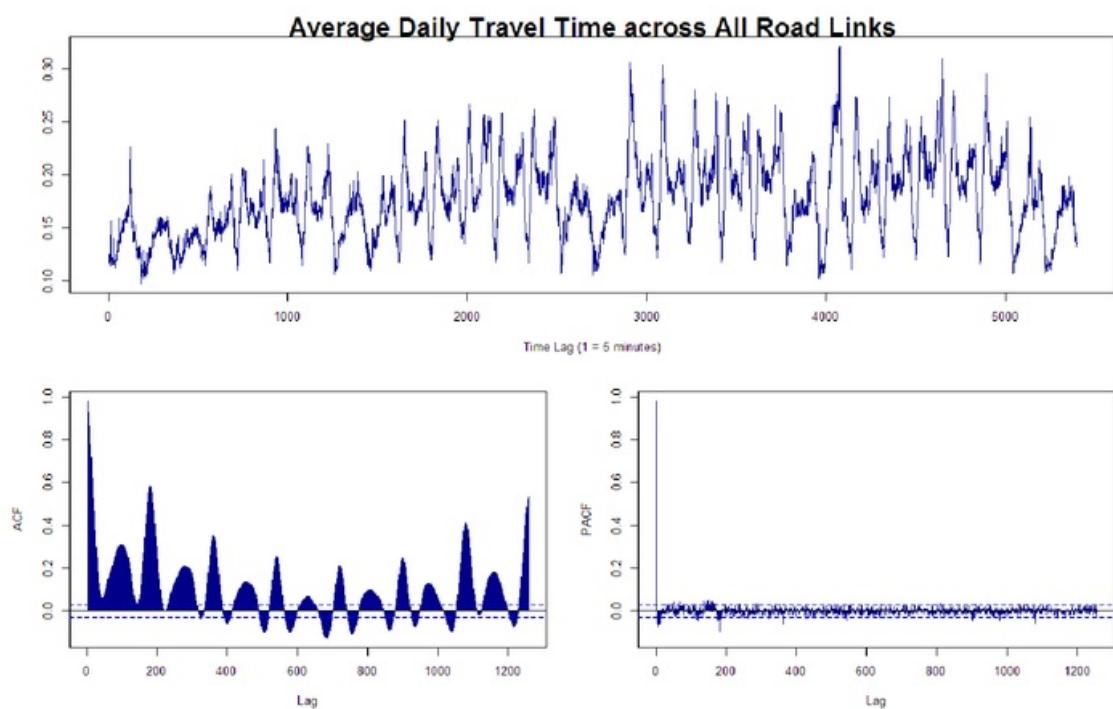


Figure 9: Time series (top), autocorrelation function (left), and partial autocorrelation function (right) for average daily travel time across all links. Dotted blue lines indicated 95% significance.

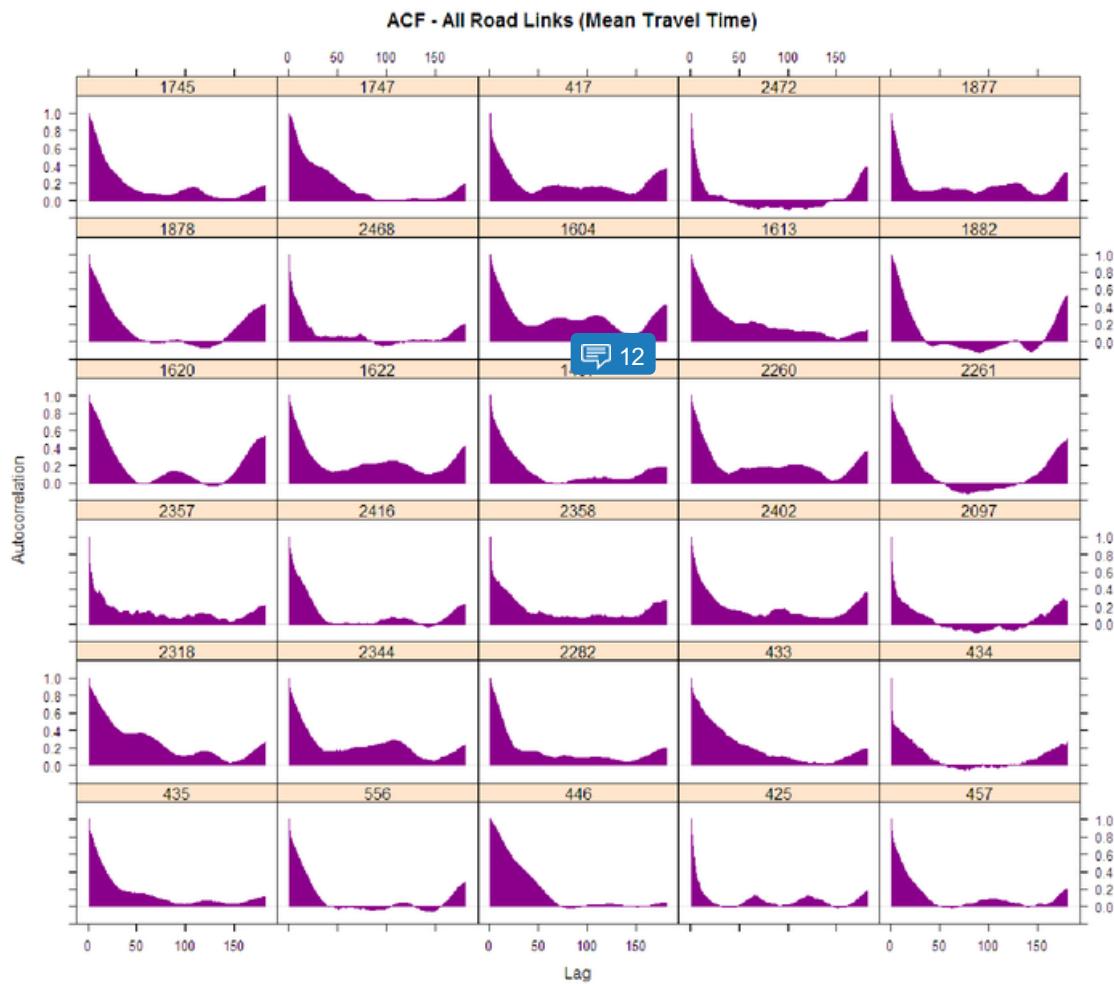


Figure 10: Autocorrelation plots for all road links using mean travel time per time lag. All road links show a cyclical trend in autocorrelation, peaking at the beginning and end of the days. Some roads show an additional increase in autocorrelation during the middle of the day, with either one or two peaks autocorrelation peaks observed.

Scatter plots of mean travel time observations against lagged observations at lags ranging from 1 to 36 (Figure 11) show observations recorded within 5-15 minutes of each other are strongly correlated. This is also shown by calculation of the Pearson correlation coefficient, which is 0.98 for a time lag of 1. Travel times measured up to 30 minutes apart are also well correlated ($r = 0.89$), however the spread of the points around the diagonal is clearly increased, indicating a weakening relationship between the observations.

Partial autocorrelation (PACF) plots show autocorrelation per time lag after accounting for all autocorrelation at lower lags. Despite this, some significant autocorrelation exists (Figure 12), although it is only present during the first 2-5 time lags (Figure 13), and then again at 160 lags (Figure 12). This information could be useful for ARIMA forecasting. Neither the ACF nor PACF plots show the autocorrelation decaying to zero, this indicates that the data is not stationary, and so must be transformed to stationarity before ARIMA analysis.

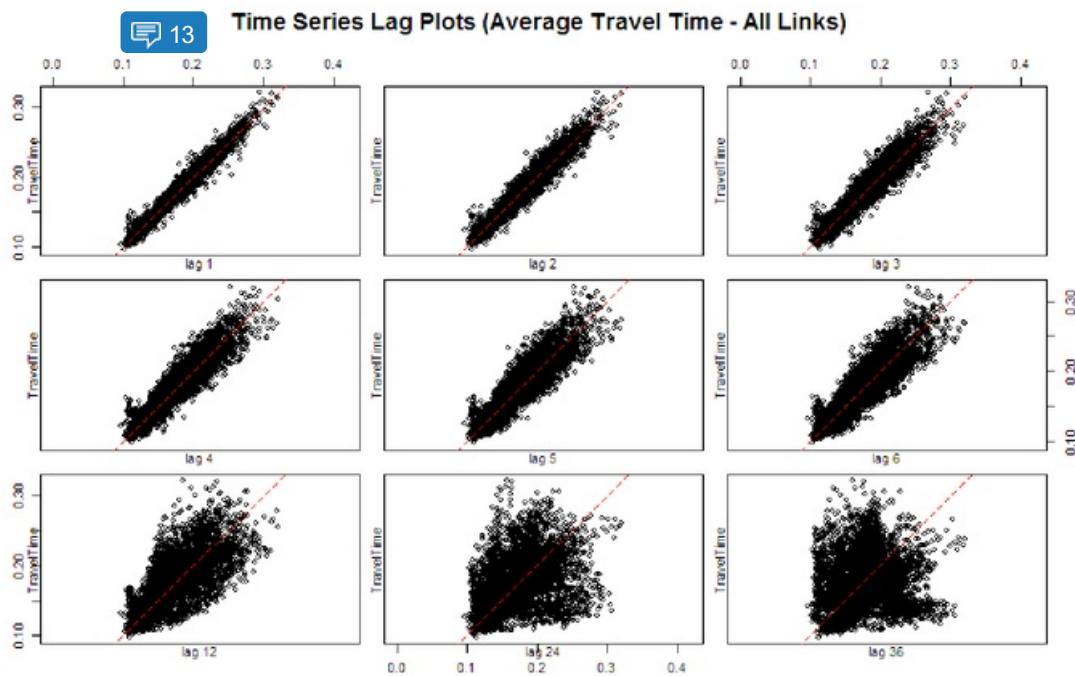


Figure 11: Lag plot for average travel time across all links, showing observations plotted against observations at lags of 1-6, 12, 24, and 36. These lags represent time intervals of 5, 10, 15, 20, 25, 30 minutes, and 1, 2, and 3 hours. Travel times measured up to 30 minutes apart are well correlated ($r = 0.89$), however the spread of the points around the diagonal is clearly increased, indicating a weakening relationship between the observations. This is particularly true for high travel times, as is indicated by the heteroscedasticity (cone-like shape) of the distribution. Observations made 3 hours apart show very little correlation overall, as with a correlation coefficient of 0.11. However, some subtle, weak, correlation may still exist across a few of the links, as it appears that two separate correlation patterns are overprinting each other.

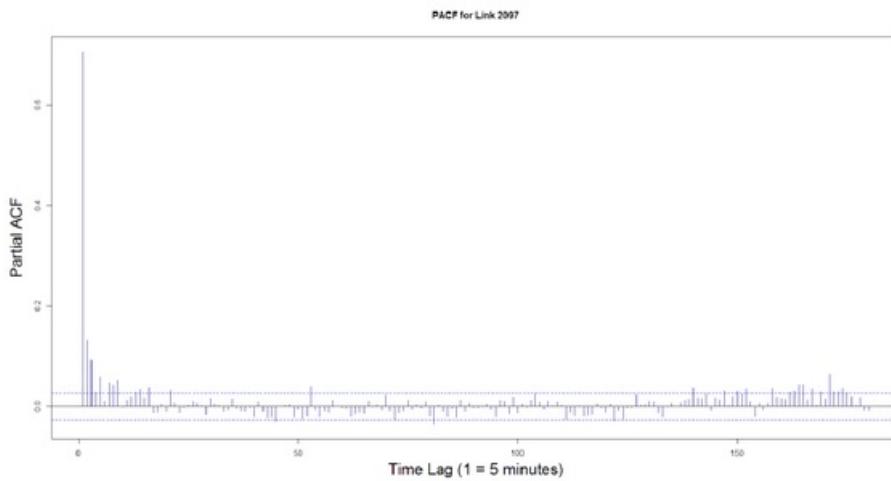


Figure 12: Partial autocorrelation function plot for link 2097, for the 1st January 2011, showing that although autocorrelation initially becomes insignificant at the beginning of the day, it becomes significant again around 160 lags. Blue lines indicated 95% confidence limits.

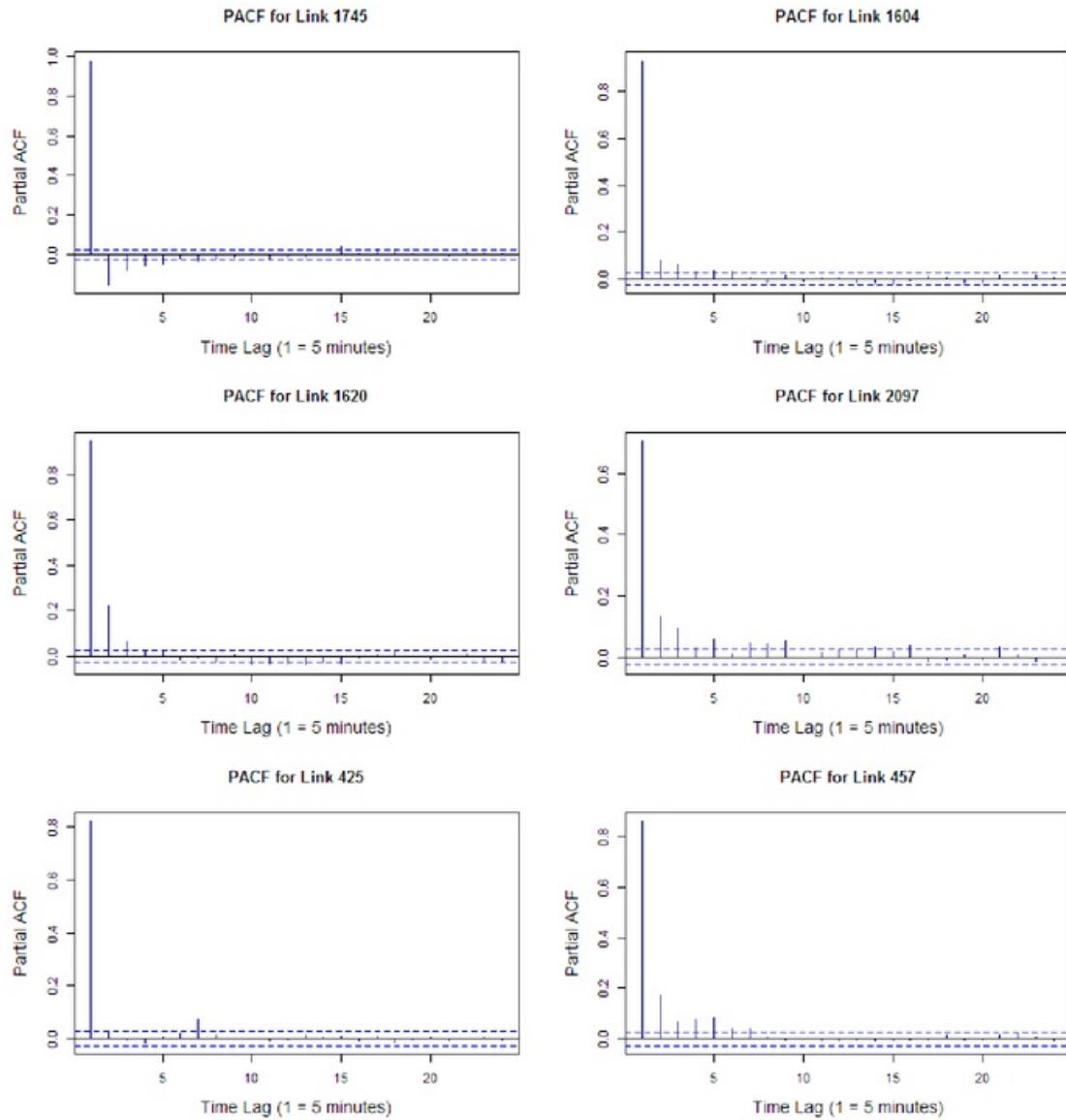


Figure 13: Partial autocorrelation function plots showing the autocorrelation at a particular lag, after accounting for autocorrelation at all lower lags. With the exception of link 425, all other links show some correlation remains within the first 5 time lags. Link 1745 shows the negative autocorrelation, whilst the remaining links show positive correlation.

2.4. Spatial Autocorrelation

Space-time autocorrelation analysis for all links, that also considers the travel times of adjoining links, shows a similar pattern to that seen in Figure 14. However, the autocorrelation of observations is not as strong at only ~ 0.25 at lag 2, compared to ~ 0.9 for the non-space-time autocorrelation. This suggests that although there is some correlation between travel times of connecting links, there correlation is not as strong as might be expected. PACF shows some minor, statistically significant, negative correlation at lags 6-7, perhaps indicating that when travel times are high on one road, they are lower on the roads connecting to them (Figure 15).

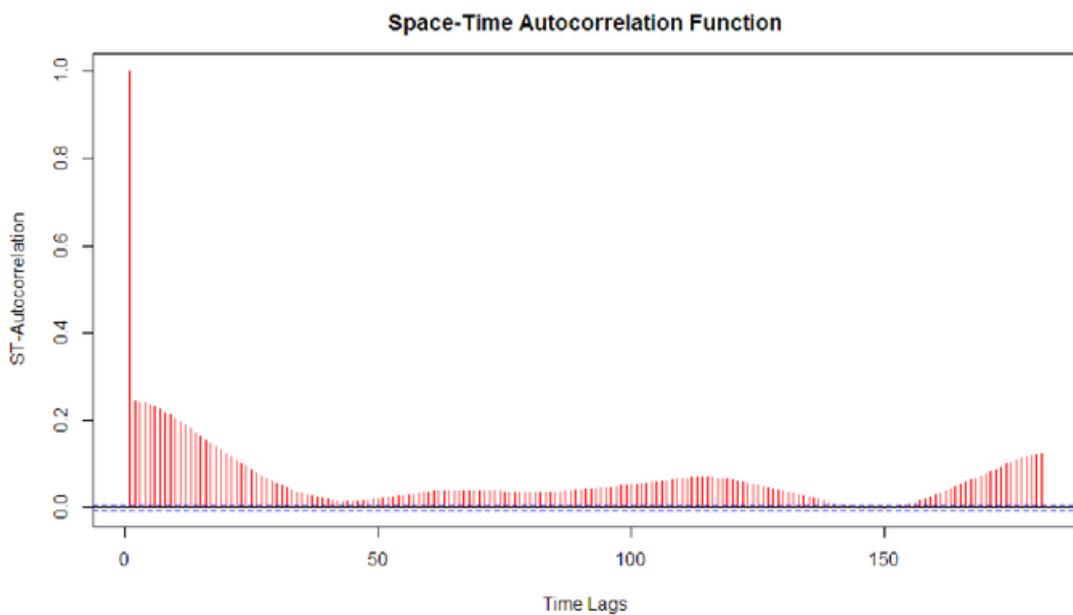


Figure 14: Space-time autocorrelation plot for all road links. Blue lines indicated 95% confidence limits.

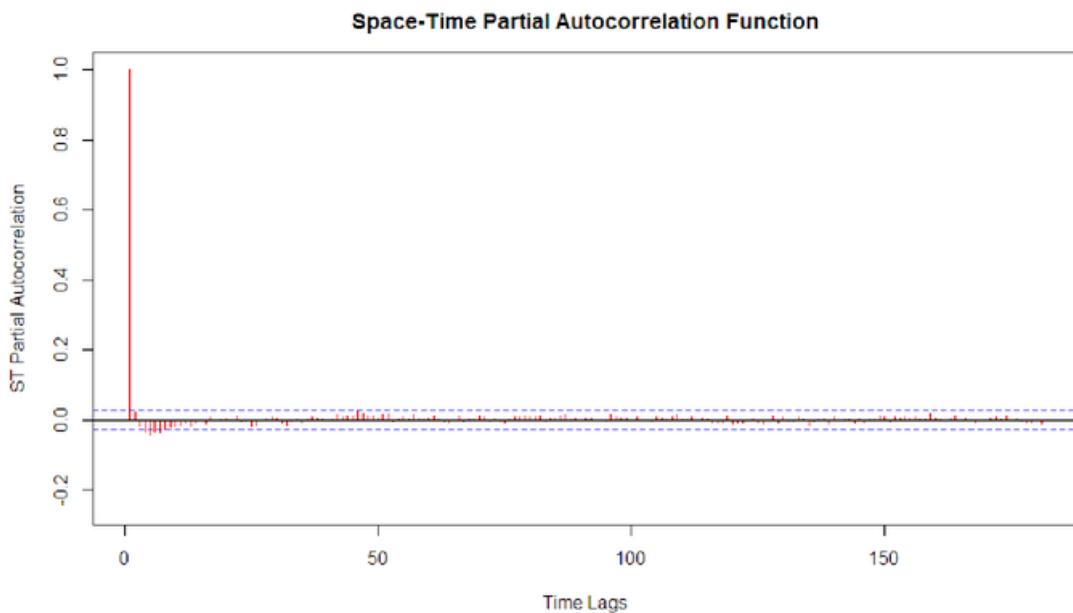


Figure 15: Partial space-time autocorrelation plot for all road links. Blue lines indicated 95% confidence limits.

A global Moran's I test was carried out to determine whether travel times across London are spatially clustered, however using the data from the selected links showed very poor results, with Moran's I and P-values of 0.08 and 0.13 respectively. The local Moran test was equally as poor, and so the Moran's I tests were repeated using the entire dataset of 256 road links. This resulted in a global Moran's I of 0.42 for both the test under randomisation, and the Monte-Carlo simulation. Both results produced statistically significant to within 99% confidence limits, indicating that there is a negligible probability that the travel times across London are not somewhat autocorrelated.

The local Moran's I test for the whole dataset gives values ranging from -3.34 to 5.90, with varying levels of significance. Links showing strong positive autocorrelation with their neighbours (Figure 16) are the north-eastern and north-western parts of the North Circular (A406), and it's connecting roads (A1, A12, A13, and A40), as well as some select roads in Central London. The most strongly negatively autocorrelated links are found in South London. The P-values of these correlation coefficients showing statistically significant autocorrelation to within 95% confidence limits (Figure 17) are consistent with the roads showing strong positive autocorrelation. Those with low statistical significance are the roads showing strong negative autocorrelation.

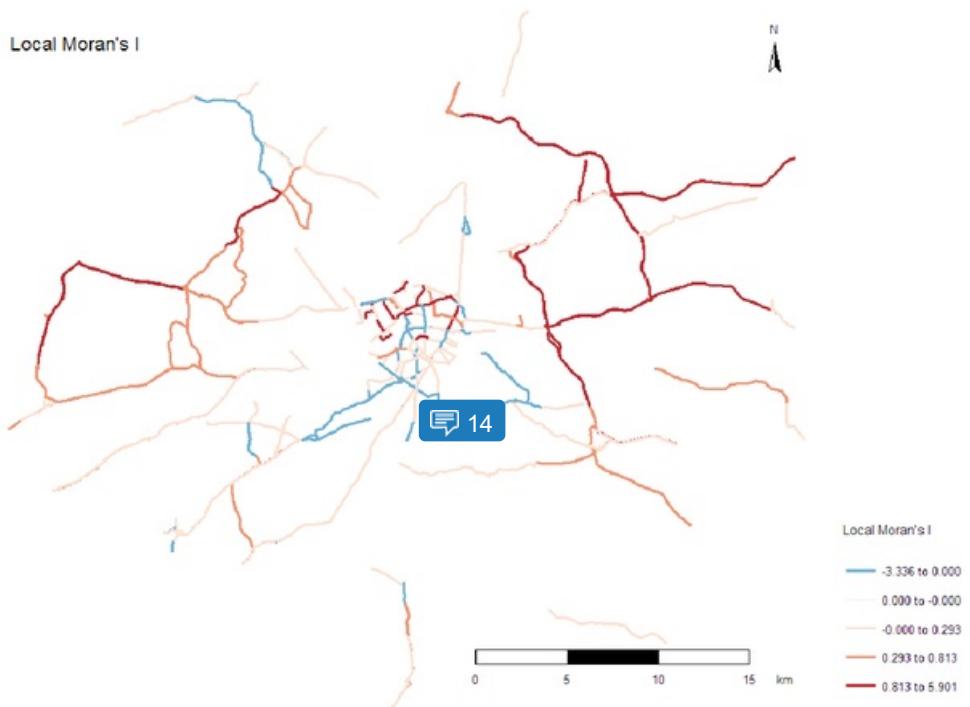


Figure 16: Local Moran's I values for all 256 road links with travel time information in London. Usually the Moran statistic should not exceed 1.0, however, as mentioned previously the travel time data is strongly skewed and contains many extreme outliers. This may have affected the Moran statistic calculation and resulted in calculated values outside of the range -1 to 1

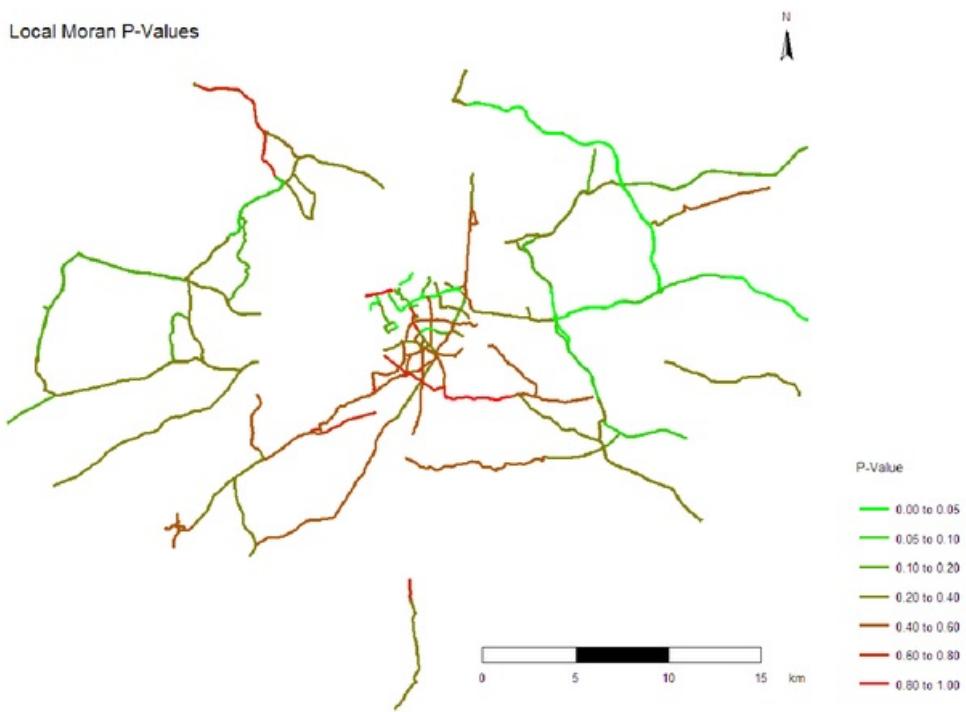


Figure 17: P -values for Local Moran's I test carried out using all 256 road links with travel time information in London.

3. Space-Time Forecasting

3.1. Autoregressive Integrated Moving Average (ARIMA) – Christopher Baxter

3.1.1. Method

An autoregressive integrated moving average (ARIMA) model was used in the experiment to predict urban travel times in London. ARIMA is a popular model used to forecast space-time series. Cheng & Wang (2011) note that the model represents each observation as at time t as a weighted linear combination of the previous observations (Cheng & Wang, 2011).

The ARIMA model has several steps. The Box-Jenkins three-step estimation procedure is widely used (Cheng & Wang, 2011). This includes: model identification using the autocorrelation (ACF) and partial autocorrelation (PACF) functions, model estimation and diagnostic checking (Cheng & Wang, 2011).

3.1.2. Experimental Set-Up

The dataset was divided into a training subset and a test subset to validate the results of the model. The training subset consisted of travel time measurements for the first 23 days of the dataset. The model was then used to predict urban travel times for the remaining 7 days of the dataset. The observed measurements for the test subset were used to validate the predictions of the model and assess quantitatively its accuracy.

Five links were selected to predict urban travel times. These were link 2097, 2358, 446, 1613 and 1747. These links were chosen because, when using SVR to model them, they formed a selection of models producing a range low and high RMSE and R² values between them.

Several parameters were used to train the ARIMA model:

- Autoregressive order (denoted as p)
- Integration order of the model (denoted as d)
- Moving average order (denoted as q)

To meet the assumptions of the ARIMA model, data for each road link was transformed to stationary using differencing where travel times showed strong cyclical patterns.

Parameters were chosen for each road link based upon the results of the autocorrelation and partial autocorrelation analysis, as shown in the results section:

- **Link 2097:** ACF plot after differencing (Figure 20) showed alternating positive and negative correlations, decaying to zero, therefore an autoregressive model was chosen. A autoregressive order of 5 was chosen because the PACF plot (Figure 21) indicated up to lag 5 as statistically significant.
- **Link 2358:** ACF plot after differencing (Figure 20) showed alternating positive and negative correlations, decaying to zero, therefore an autoregressive model was chosen. PACF plot (Figure 21) indicated a AR order of 8 should be chosen.
- **Link 446:** ACF plot (Figure 19) showed alternating positive and negative correlations, decaying to zero. Therefore, an autoregressive model was chosen. An AR order of 2 was chosen after examination of the PACF plot (Figure 21).
- **Link 1613:** ACF plot (Figure 19) displayed alternating positive and negative correlations, decaying to zero. As a result, an autoregressive model was chosen. The PACF plot (Figure 21) used to select an AR order of 6 for the ARIMA model for this link.
- **Link 1747:** ACF plot (Figure 19) showed alternating positive and negative correlations, decaying to zero, so an autoregressive model was chosen for this link. Additionally, after differencing (Figure 20) a cyclical pattern of positive and negative patterns is evident at approximately 10 lags. Therefore, a seasonal component with an autoregressive term ($t-10$) is also included in the model for this link.

3.1.3. Results

Time series plots were used to visualise the temporal patterns present in travel times at each road link. None of the selected links displayed an upward or downward trend. Link 2097 and 2358 showed strong cyclical patterns as can be seen in Figure 18. Link 1747 also showed cyclical trends. In order to build an ARIMA model these links were differenced in the pre-processing stage to satisfy ARIMA's model that the dataset is stationary.

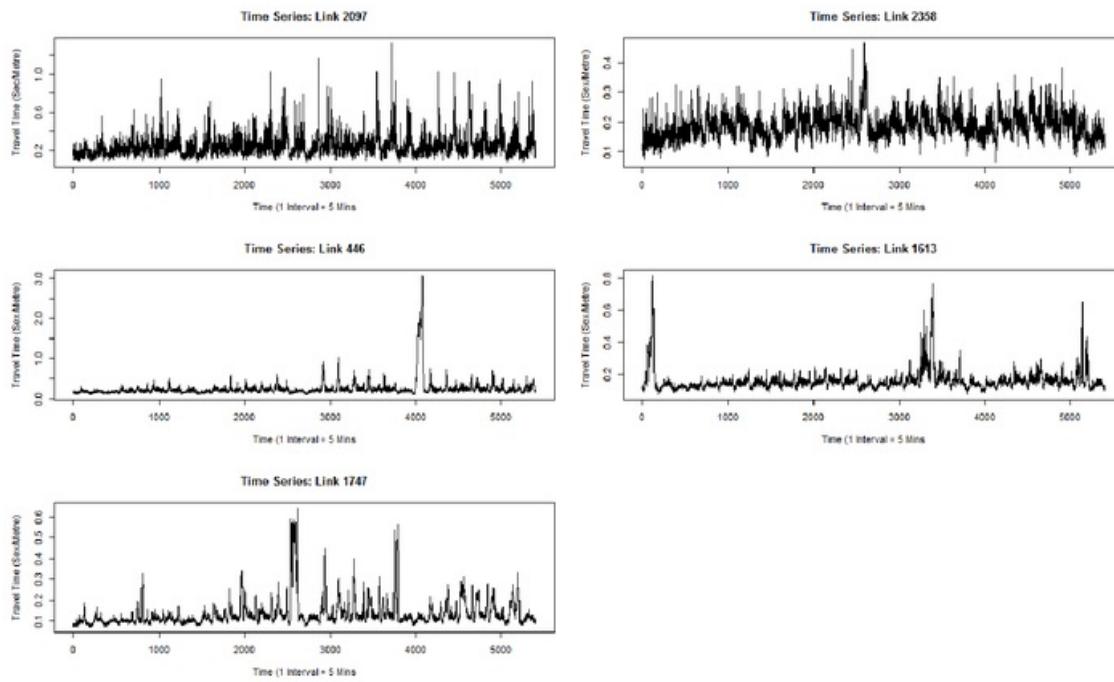


Figure 18: Time series plots for each road link selected for analysis.

Figures 19 and Figure 20 show the results of the differencing. Both the autocorrelation function (ACF) and partial autocorrelation function (PACF) were used in selecting the appropriate model for each road link. In all three links, using the differencing technique was only able to remove some of the cyclical patterns, primarily towards the end of the 30-day dataset.

The partial autocorrelation function was also used to identify the autoregressive order of the model for each road link (p).

All the links showed statistically significant partial autocorrelations. Link 2097 and link 2358 showed statistically significant correlations to lag 4 and lag 8. This shows that travel times at a certain temporal point were influenced by travel patterns in the last 20 – 40 minutes previously. This also informed the autoregressive order chosen for each link when defining the parameters for the ARIMA model, as discussed in the experimental setup.

The results of the predictions of all ARIMA models is shown in Figure 22 - 24Figure 24. All prediction models are not able to predict urban travel times for the last 7 days of the dataset (predicted values are shown by the red line; actual values are shown by the black line). All models predict increasing travel times initially, followed by a constant travel time. In most links, this was close to the average shown by the black line over the 7-day period. However, for link 1747 the model predicted values significantly higher than the observed values.

CEGEG076: Spatio-Temporal Analysis and Data Mining
Christopher Baxter and Lucille Ablett

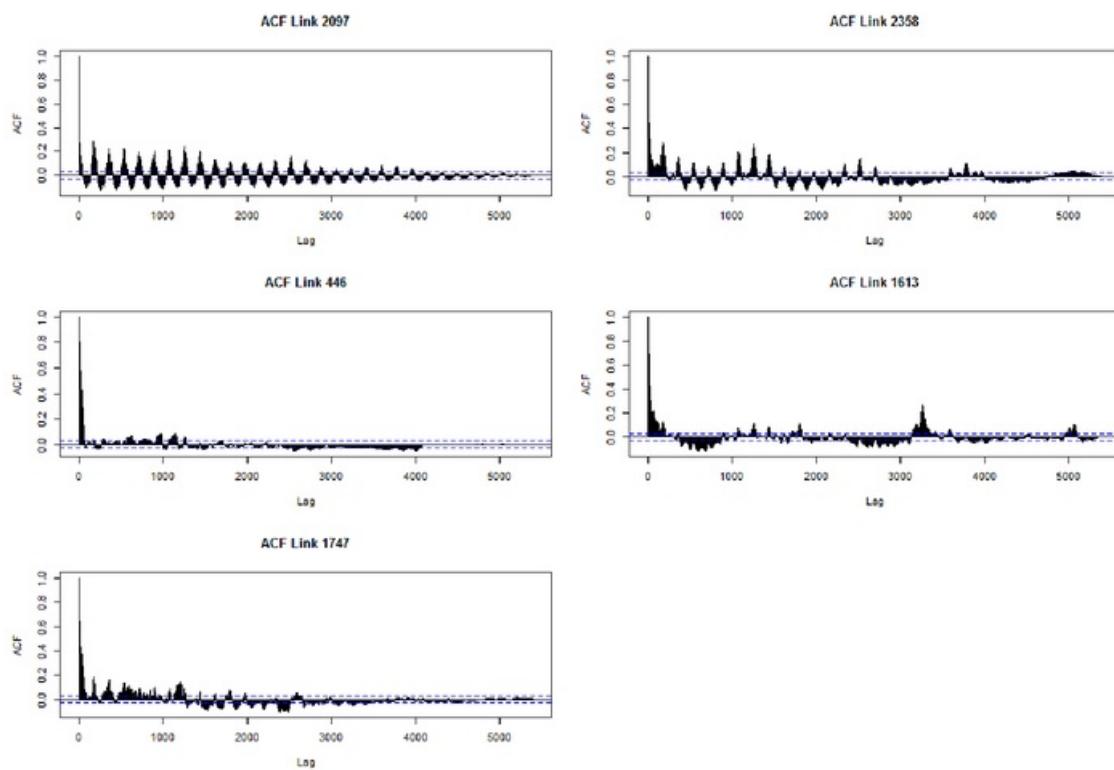


Figure 19: Autocorrelation function plots for chosen road links.

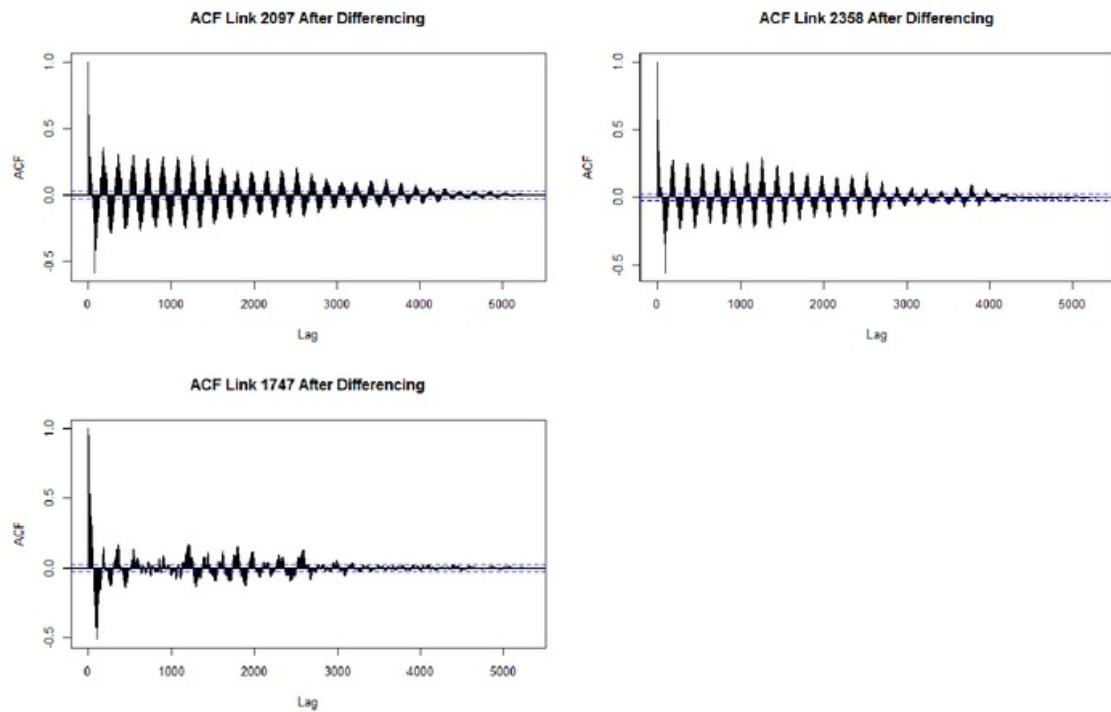


Figure 20: Autocorrelation function plots for road links 2097, 2358 & 1747.

CEGEG076: Spatio-Temporal Analysis and Data Mining
Christopher Baxter and Lucille Ablett

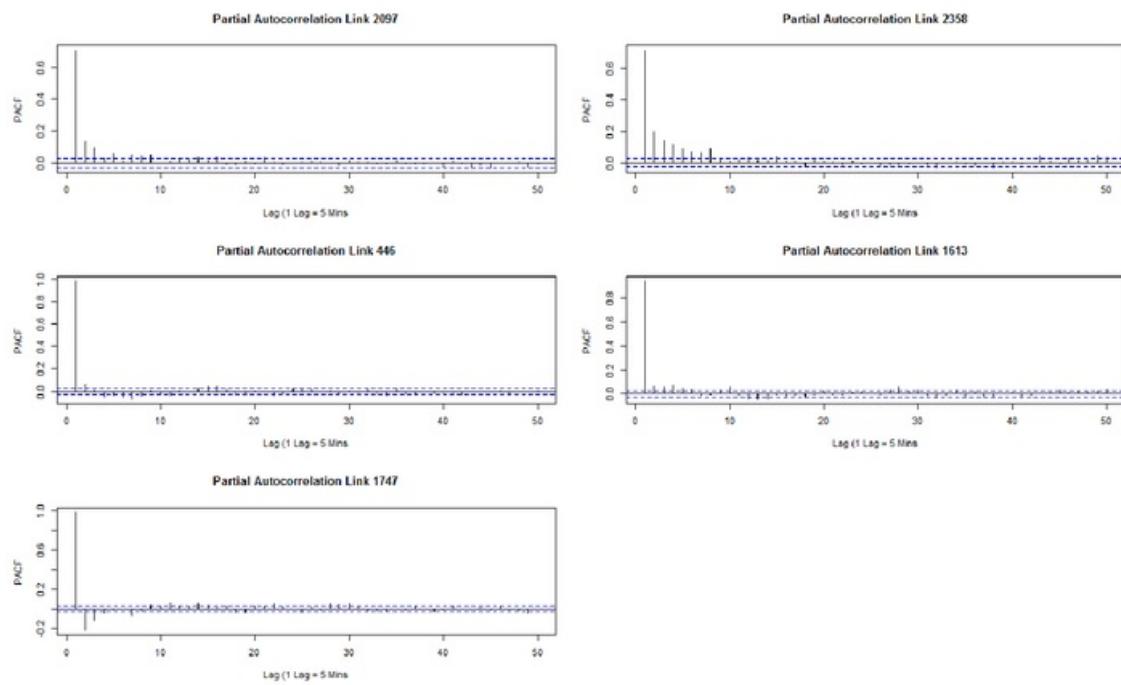


Figure 21: Partial autocorrelation function (PACF) plot for selected road links.

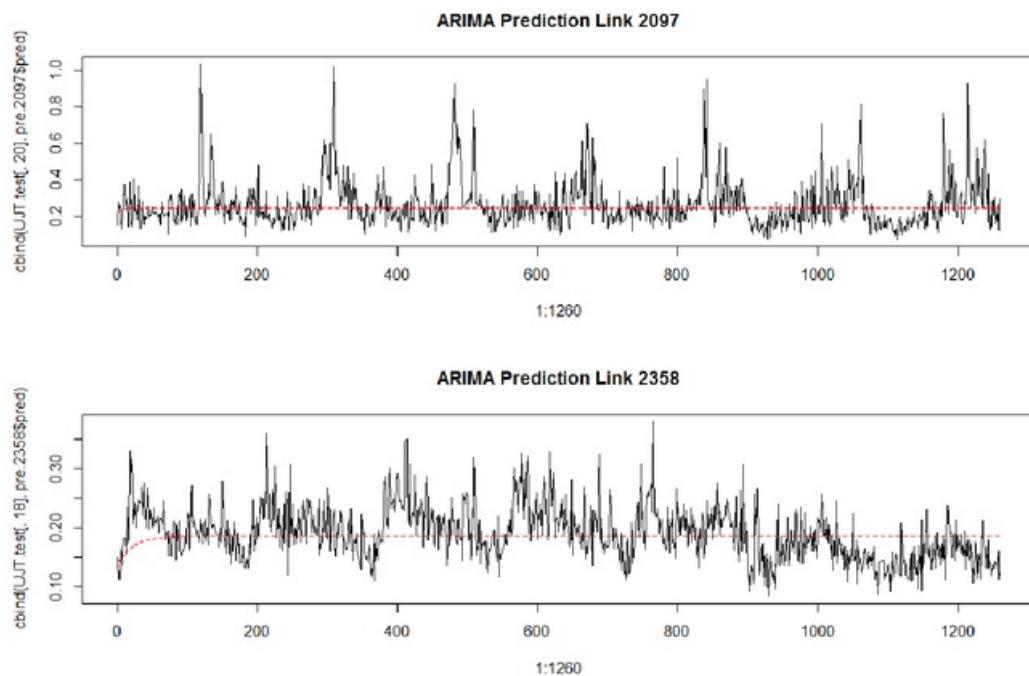


Figure 22: ARIMA predicted and observed travel times: links 2097 & 2358.

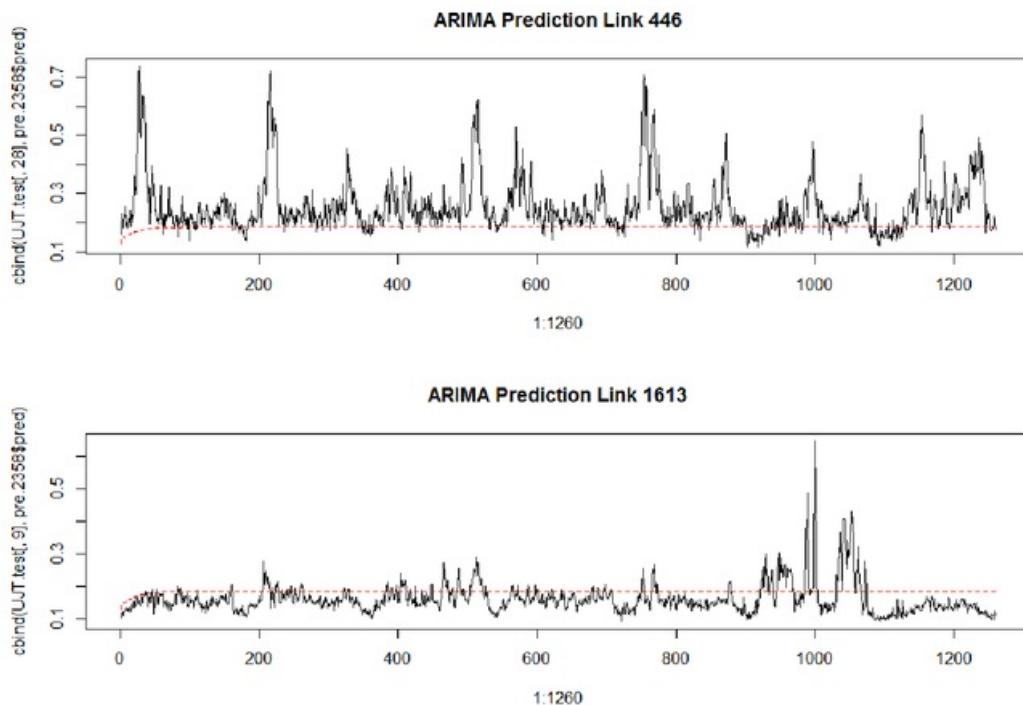


Figure 23: ARIMA predicted and observed travel times: links 446 & 1613.

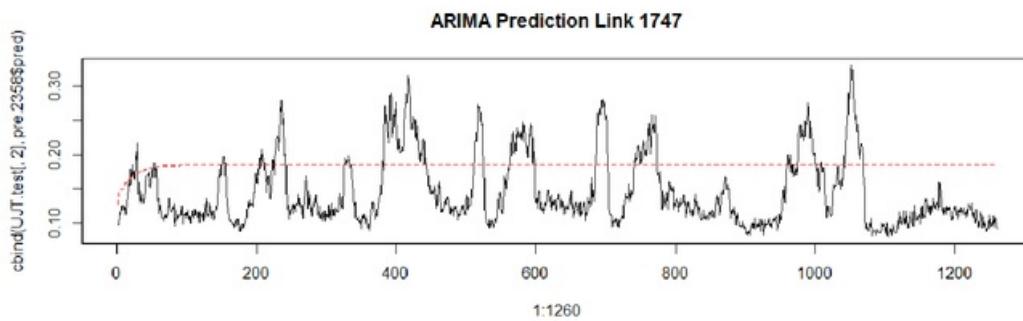


Figure 24: ARIMA predicted and observed travel times: link 1747.

The normalised root mean square error (NRMSE) was used to quantitatively assess the performance of each model. This value represents the square root of mean squared error and is a widely-used measure of model performance. Table 2: NRSME values for each road link shows the NRMSE values for each road link.

Model performance, as measured by NRSME, differed significantly between links. Links 2097 and 2358 performed considerably worse than the prediction models for links 446, 1613 and 1747, both scoring a NRSME value above 0.6.

Link	NRSME Value
2097	0.7029223
2358	0.6632987
446	0.1680887
1613	0.2830049
1747	0.2005903

Table 2: NRSME values for each road link.

3.2. Support Vector Regression (SVR) – Lucille Ablett

3.2.1. Method and Experimental Set-Up

As mentioned previously, SVR is a machine learning method that uses kernels to map data into a high dimensional feature space, allowing the use of linear algorithms to model non-linear datasets. For this investigation, SVR was utilised to predict the last 7 days' worth of data, based on the first 23 days of January.

Based on this requirement, the data was divided into training and testing datasets. The number of data points in each training and testing set depended on the number of embedding dimensions, since the higher the number of previous observations, the fewer number of values remain for prediction. The splitting point of the embedded time series was calculated using $split = 5400 - 1260 - m$, where 5400 is the total number of observations for one road link, 1260 is the total number of observations for one road link in 7 days, and m is the embedding dimension.

Initially, the effect of different embedding dimensions was investigated. The accuracy of the model was expected to increase with the increase in dimensionality, however an increase in dimensionality leads to an increase in processing time. The aim was to find an optimum number of lagged observations to include in the model [15] to achieve a good balance between training time and accuracy. This was established through modelling of three links using *kernlab*, and testing varying embedding dimensions until an optimal number determined. *Kernlab* was used as it allows all model parameters to be controlled, unlike *Caret*, where epsilon cannot be optimised. The model was run using sigma = 0.1, C = 10, epsilon = 0.2, and 2-fold cross-validation. The links chosen for this test were 446, 1882, and 2097, as these showed the most variation in travel times (Figure 6), and so were assumed to be more difficult to model, and therefore varying embedding dimensionality was hoped to have a more apparent effect than for other, more easily predictable, time series.

The optimal number of embedding dimensions was then used to create travel time forecast models for all links. *Caret* was selected for this task as it allows models to be trained using a grid of parameters for both sigma and C. Due to the large number of links, this allowed the optimal model to be determined more quickly and efficiently than using *kernlab*. *Kernlab* was then utilised to create models for links 1745, 1747, 417, 2472, 2097, 434, and 446, with the aim of further optimising the models through the inclusion of epsilon in the tuning process. These models were selected due to the varying range RMSE and R² values between them. 5-fold cross-fold validation was used for the creation of all models to avoid overfitting.

Finally, the travel times of neighbouring links were used to create space-time forecasting models for the links optimised using *kernlab*. This was accomplished through use of the *st_embed* function to

embed the previous observations from the link and its neighbours into the required dimensions. Due to processing time limitations, only the previous 3 observations were included in the model.

A radial basis function kernel was used for all forecasting models, and the optimal models were considered to be those with the lowest root mean squared error (RMSE).

It should be noted that the number of support vectors cannot be controlled using either or the methods used in this investigation.

17

3.2.2. Results

Embedding Dimensionality Optimisation

Tests to determine the optimum embedding dimensionality (Figure 26) show 7 dimensions is the optimal number. For links 446 and 1882, models with 7 embedding dimensions gave the lowest processing time, and the lowest number of support vectors. Processing time remained similar between 2 and 7 dimensions, after which it increased. Training error decreased in a linear fashion with increase in dimensionality, and so the optimum model was concluded to be that with the maximum number of dimensions that did not impact on processing time. Models for link 2097 had larger errors and so more support vectors were required, thus increasing processing time for models with 6 dimensions and higher.

Individual Road Link Forecasting and Space-Time Modelling

The sigma and C values resulting in the optimal models for all links are detailed in Table 3, with corresponding training error plots in Appendix B. Sigma values range from 0.0005 to 0.01, whilst C values range between 1 and 250. The links optimised further using *kernlab* are highlighted in yellow. However, it was found that optimising epsilon only produced minor improvements where the training error was already low, and so the original models created using *caret* were retained.

Inclusion of travel times from neighbouring links reduced the number of support vectors required for three of the seven links tested (Table 5). The greatest reduction was seen for link 1747 (1.1%), and this also resulted in a minor reduction in RMSE error, and an increase in R^2 . The second largest change was seen for link 446 (0.8%), however this resulted in a significant increase in model error of 0.008 s/m, and an increase in R^2 .

The errors between the observed and predicted values for all links, for both non-spatial and spatial models can be seen in Figure 27 to 40. SVR successfully predicts the general trends in the time, however struggles with excessive variation (noise), such is observed for links 417 and 2097 (Figures 31 and 32, and Figures 35 and 36 respectively). Although SVR successfully predicts unusually large increases or decreases in travel time throughout the time series, particularly the rush hour peaks, it also consistently fails to predict the extent of them, even for the models with the lowest RMSE and R^2 values, such as link 446 (Figures 39 and 40).

ACF plots of the residuals for non-spatially and spatially dependent models show some significant autocorrelation is still apparent for all link models, particularly between lags 2-5. For some links this is positive autocorrelation, and others it is negative. For a few links, the inclusion of spatial data in the forecasting of travel times leads to increased autocorrelation (434, Figure 38), sometimes significant (1747, Figure 25). However, for the majority of links autocorrelation either remains the same, or is reduced (for example 446).

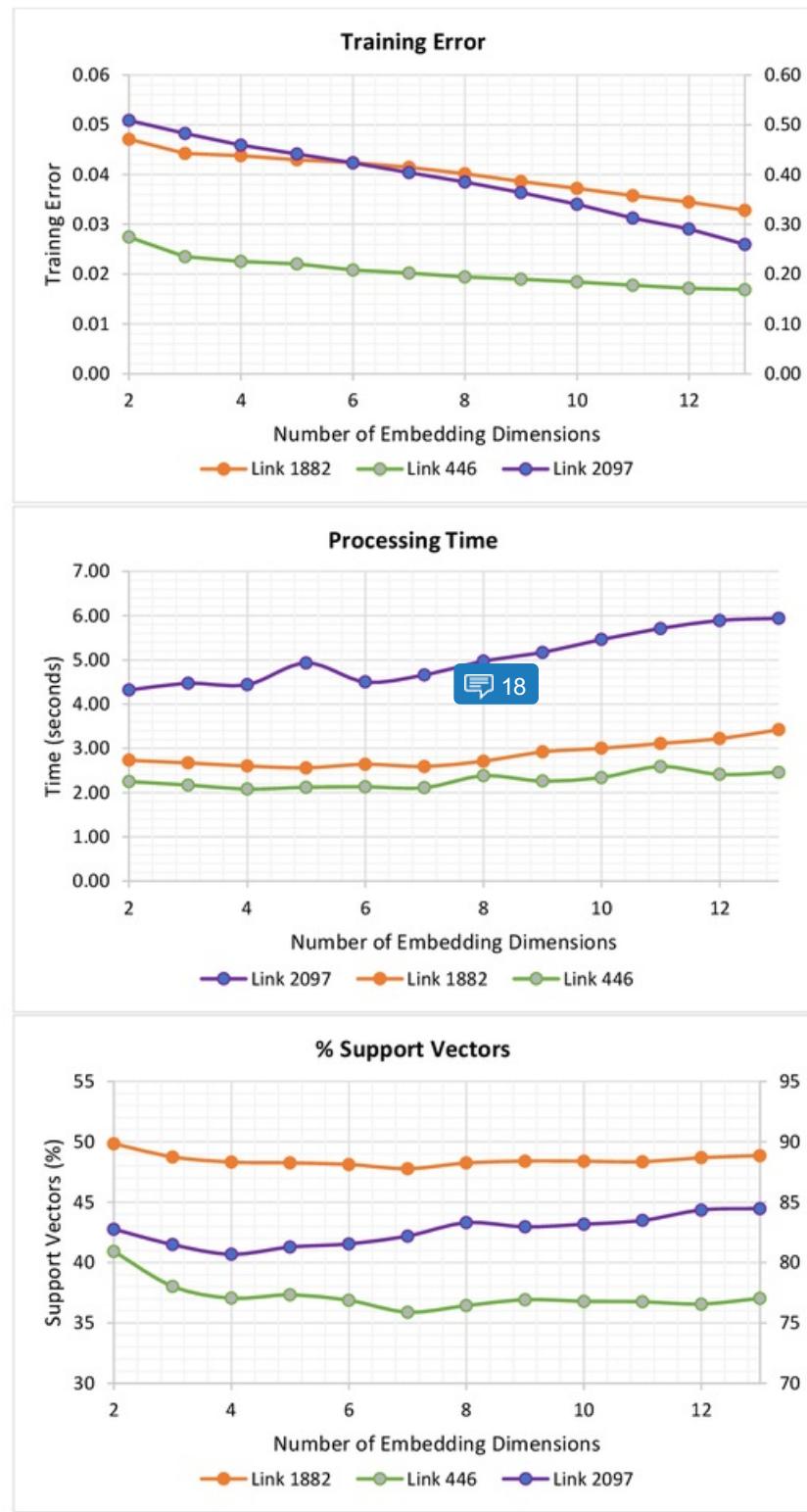


Figure 26: Training error (top), processing time (middle), and percentage of support vectors (bottom) required to produce a forecasting model for links 1882, 446, and 2097, using $\sigma = 0.1$, $C = 10$, $\epsilon = 0.2$, and 2-fold cross-validation. Note that, where applicable, link 2097 is plotted on the secondary axis.

CEGEG076: Spatio-Temporal Analysis and Data Mining
Christopher Baxter and Lucille Ablett

<i>Count</i>	<i>Link ID</i>	<i>Sigma</i>	<i>C</i>	<i>RMSE</i>	<i>R</i> ²	<i>Figure</i>
1	1745	0.01	10	0.0138	0.9635	27
2	1747	0.02	220	0.0134	0.9669	29
3	417	0.0005	50	0.0264	0.7019	31
4	2472	0.01	1	0.0581	0.6537	33
5	1877	0.01	10	0.0220	0.7445	
6	1878	0.01	10	0.0361	0.8794	
7	2468	0.01	10	0.0269	0.6835	
8	1604	0.01	10	0.0401	0.8767	
9	1613	0.01	10	0.0190	0.9237	
10	1882	0.01	200	0.0342	0.9557	
11	1620	0.01	10	0.0346	0.9115	
12	1622	0.01	10	0.0216	0.8504	
13	1407	0.001	100	0.0380	0.7850	
14	2260	0.0005	250	0.0155	0.8234	
15	2261	0.001	175	0.0171	0.8113	
16	2357	0.01	1	0.0500	0.6985	
17	2416	0.00025	10	0.0393	0.7676	
18	2358	0.001	100	0.0293	0.5547	
19	2402	0.001	70	0.0193	0.7677	
20	2097	0.0005	175	0.0805	0.5239	35
21	2318	0.01	10	0.0192	0.8709	
22	2344	0.01	10	0.0217	0.8235	
23	2282	0.01	10	0.0159	0.7841	
24	433	0.0015	50	0.0264	0.7019	
25	434	0.0005	75	0.0278	0.4224	37
26	435	0.002	50	0.0241	0.8307	
27	556	0.002	228	0.0201	0.8543	
28	446	0.004	30	0.0385	0.9747	39
29	425	0.001	220	0.0198	0.7763	
30	457	0.001	100	0.0390	0.7855	

Table 3: *Sigma* and *C* values resulting in the optimal forecasting model for all links. Highlighted links are those further optimised using kernlab, and modelled using spatial observations.

Link ID	Sigma	C	RMSE	R ²	+/- RMSE	+/- R ²	Figure
1745	0.01	30	0.0136	0.9646	-0.0002	0.0011	28
1747	0.01	20	0.0131	0.9675	-0.0002	0.0006	30
417	0.01	5	0.0264	0.7014	0.0000	-0.0005	32
2472	0.01	4	0.0588	0.6450	0.0008	-0.0087	34
2097	0.01	75	0.0831	0.4871	0.0026	-0.0368	36
434	0.01	8	0.0282	0.4072	0.0003	-0.0152	38
446	0.01	200	0.0467	0.9646	0.0081	-0.0100	40

Table 4: Sigma and C values resulting in the optimal forecasting models, taking into account travel times for adjacent links. Epsilon = 0.1 for all models. All models were created using 5-fold cross validation, and included travel times for the previous 3 time lags. Root mean squared errors (RMSE) and correlation coefficients (R²) for each model, and changes in these compared to the non-spatial models, show some improvement when taking into account spatial trends in travel times.

Link ID	Support Vectors				
	Non-Spatial	%	Spatial	%	+/- Change
1745	1377	33.3	1381	33.4	0.1
1747	1663	40.2	1618	39.1	-1.1
417	3371	81.5	3395	82.1	0.6
2472	2429	58.8	2441	59.0	0.3
2097	3397	82.2	3393	82.1	-0.1
434	3623	87.6	3612	87.4	-0.3
446	1591	38.5	1556	37.6	-0.8

Table 5: The number and percentage of observations required as support vectors to produce the optimal non-spatial and spatial forecasting models. Inclusion of spatial trends in travel time resulted in the reduction of necessary support vectors in three of the seven cases, with the greatest change seen for link 446 at 0.8%.

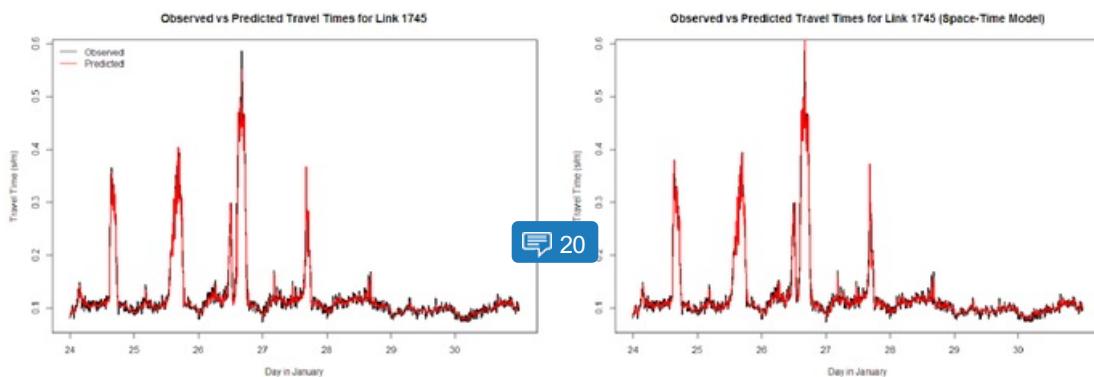


Figure 27 (left): Observed vs Predicted travel times for link 1745, for the 24th - 30th January (non-spatial).
 Figure 28 (right): Observed vs Predicted travel times for link 1745, for the 24th - 30th January (space-time model). The spatial model more accurately predicts the noise in the dataset, as well as the unusually large increase in travel time during the afternoon of the 26th January.

CEGEG076: Spatio-Temporal Analysis and Data Mining
Christopher Baxter and Lucille Ablett

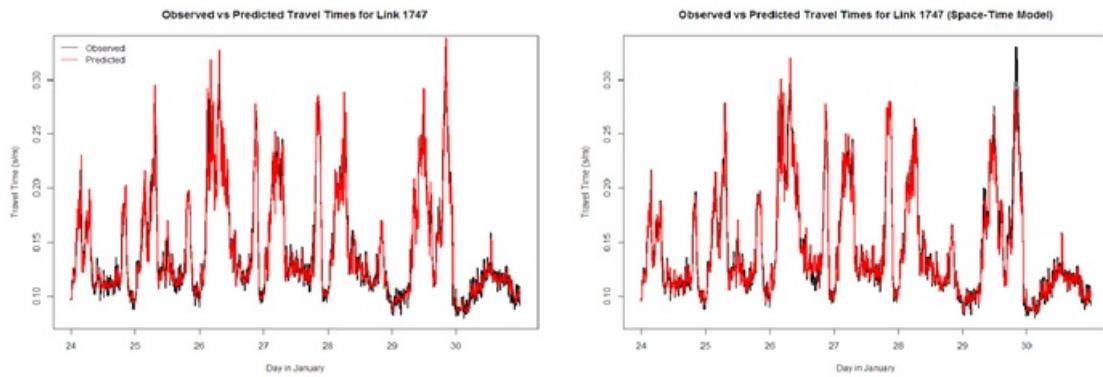


Figure 29 (left): Observed vs Predicted travel times for link 1747, for the 24th - 30th January (non-spatial).

Figure 30 (right): Observed vs Predicted travel times for link 1747, for the 24th - 30th January (space-time model). The spatial model better predicts the noise in the dataset, however it is less successful at predicting the unusually large travel time observed on the afternoon of 29th January.

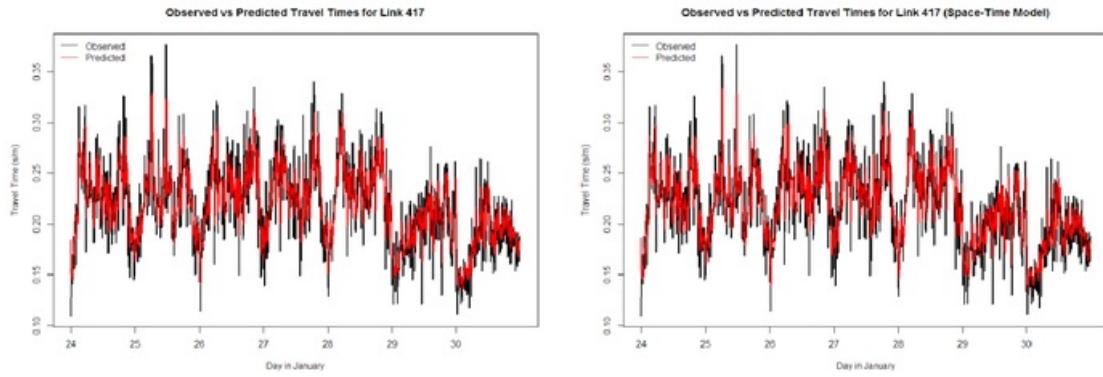


Figure 31 (left): Observed vs Predicted travel times for link 417, for the 24th - 30th January (non-spatial).

Figure 32 (right): Observed vs Predicted travel times for link 417, for the 24th - 30th January (space-time model). The spatial model shows no significant improvement in either noise prediction, nor prediction of unusually high or low travel times.

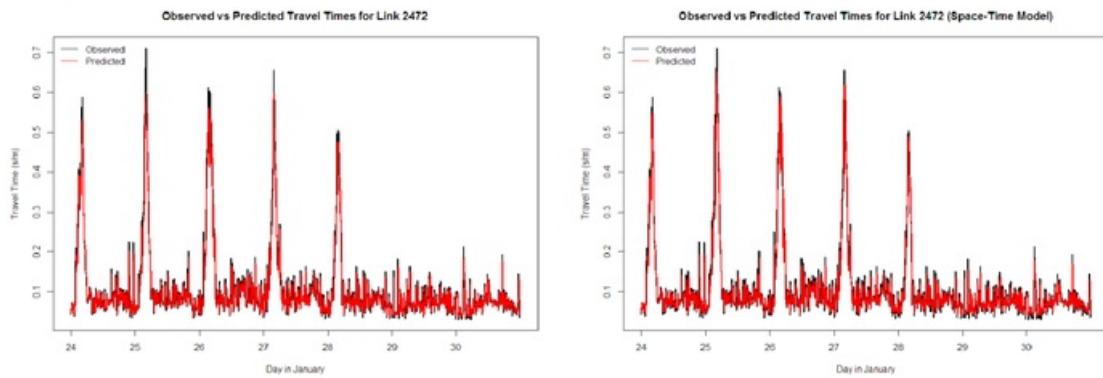


Figure 33 (left): Observed vs Predicted travel times for link 2472, for the 24th - 30th January (non-spatial).

Figure 34 (right): Observed vs Predicted travel times for link 2472, for the 24th - 30th January (space-time model). The spatial model shows minor improvements in the prediction of the peak travel time periods, however does not visibly improve in noise prediction.

CEGEG076: Spatio-Temporal Analysis and Data Mining
Christopher Baxter and Lucille Ablett

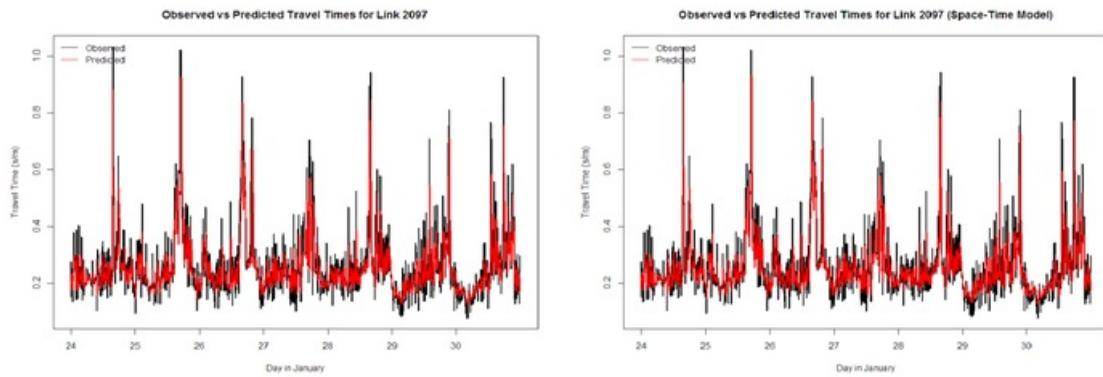


Figure 35 (left): Observed vs Predicted travel times for link 2097, for the 24th - 30th January (non-spatial).

Figure 36 (right): Observed vs Predicted travel times for link 2097, for the 24th - 30th January (space-time model).

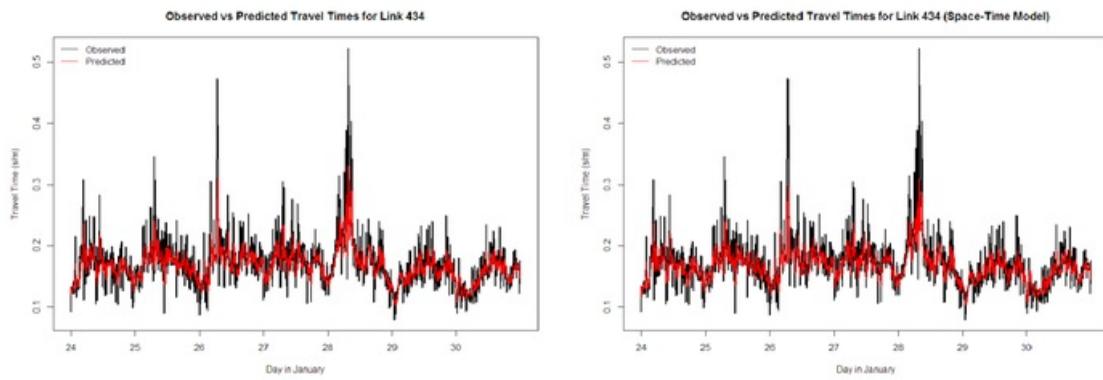


Figure 37 (left): Observed vs Predicted travel times for link 434, for the 24th - 30th January (non-spatial).

Figure 38 (right): Observed vs Predicted travel times for link 434, for the 24th - 30th January (space-time model). The spatial model does not visibly improve the prediction of travel times for this link.

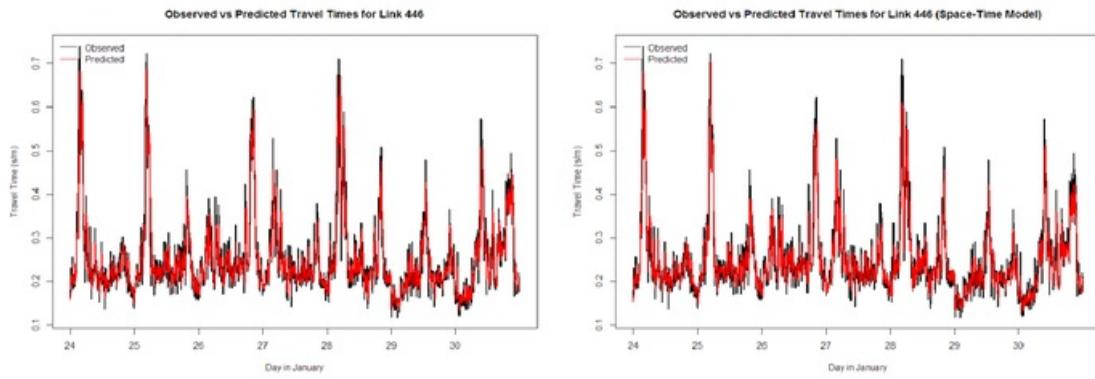


Figure 39 (left): Observed vs Predicted travel times for link 446, for the 24th - 30th January (non-spatial).

Figure 40 (right): Observed vs Predicted travel times for link 446, for the 24th - 30th January (space-time model). The spatial model does not visibly improve the prediction of travel times for this link, if anything the space-time model is a worse estimation than the non-spatial model.

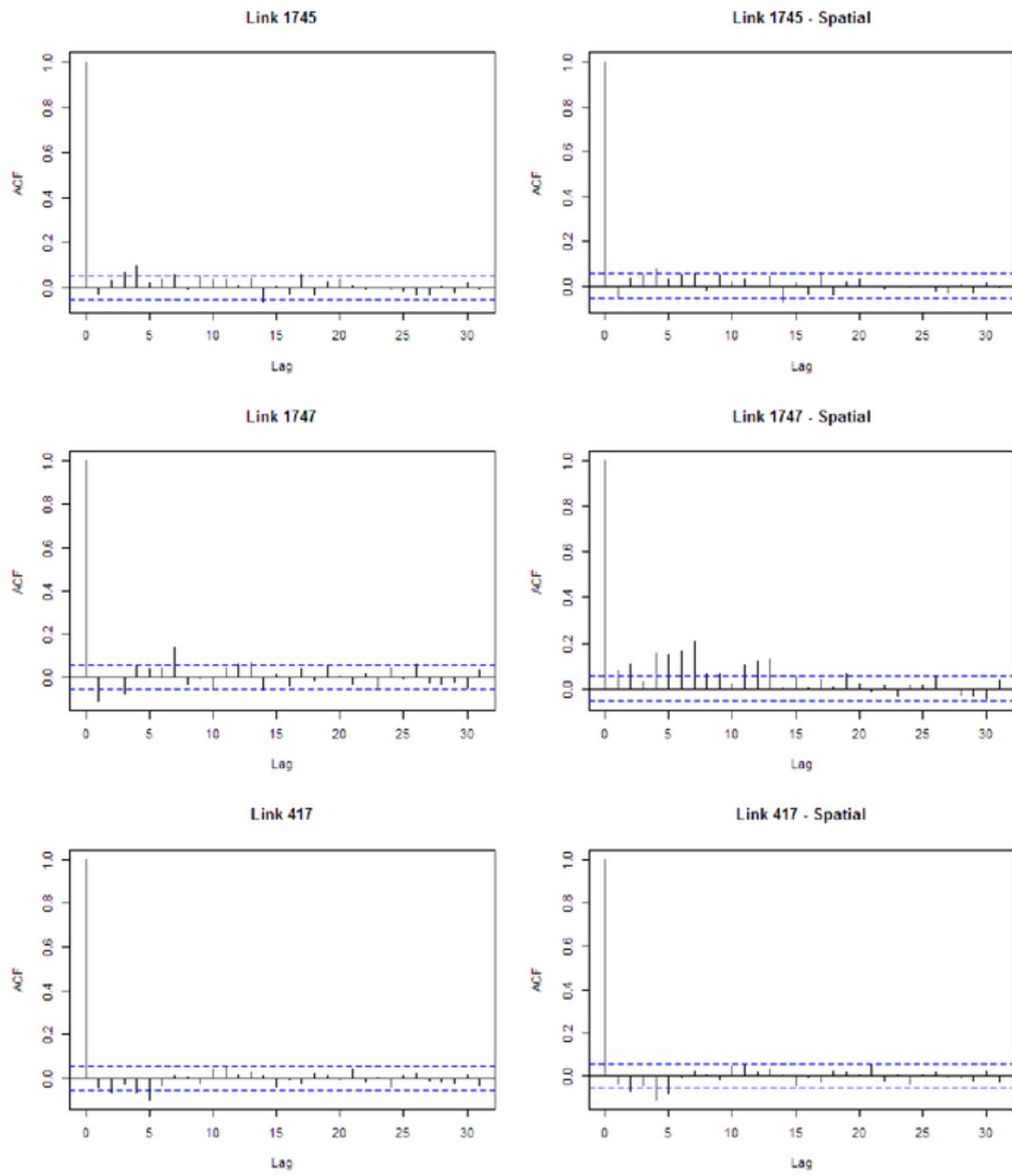


Figure 41: ACF plots of the residuals for non-spatial (left) and spatial (right) model travel time predictions for links 1745 (top), 1747 (middle), and 417 (bottom). Blue lines indicated 95% confidence limits. Some significant autocorrelation is still apparent for all link models. Link 1745 shows autocorrelation remaining at lags 4 and 14, whilst for link 1747 autocorrelation has increased significantly between lags 2 and 13. Minor negative autocorrelation has increased at lags 4 and 5 for link 417.

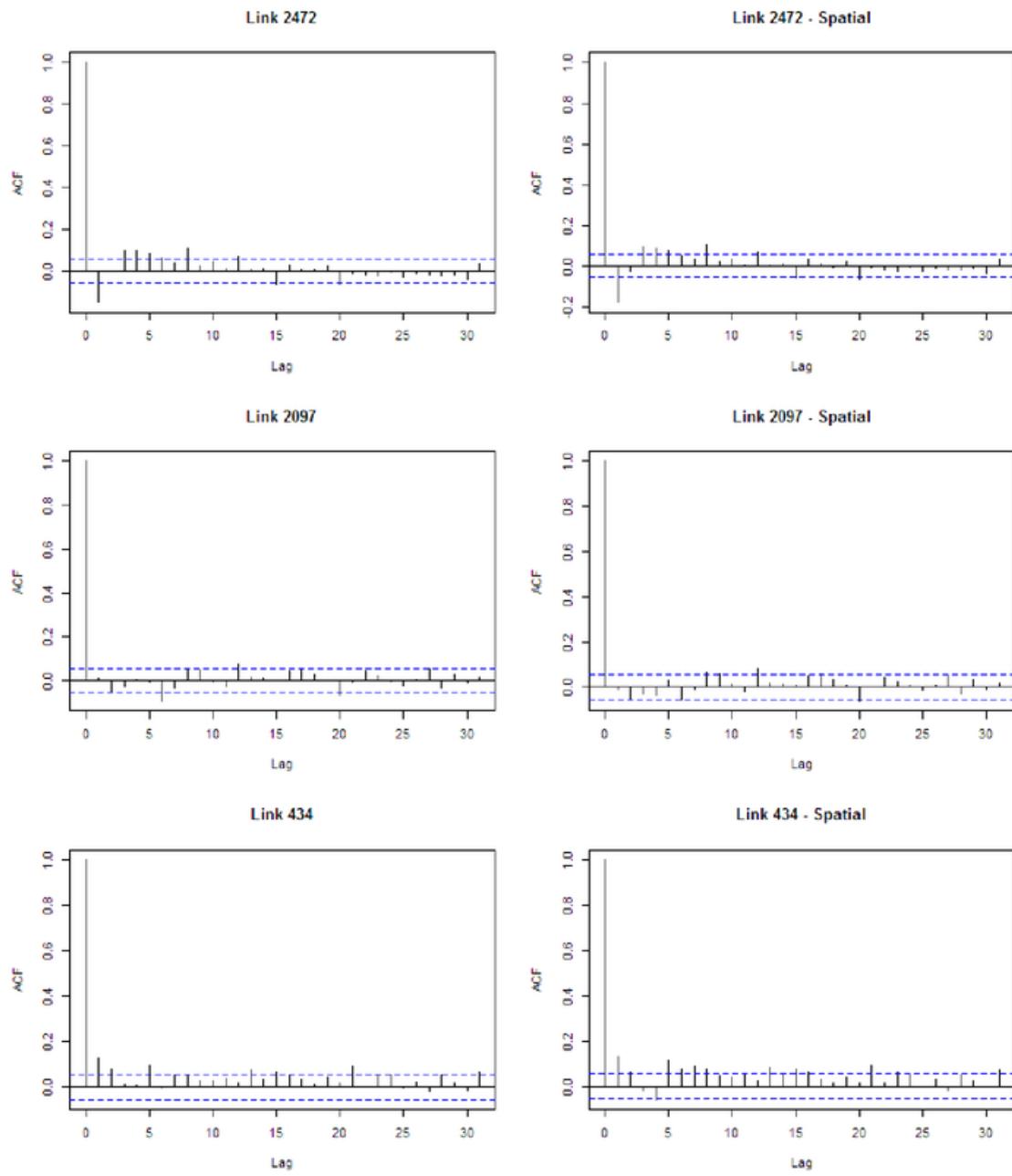


Figure 42: ACF plots of the residuals for non-spatial (left) and spatial (right) model travel time predictions for links 2472 (top), 2097 (middle), and 434 (bottom). Blue lines indicated 95% confidence limits. Some significant autocorrelation is still apparent for all link models, particularly for link 2472 at lags 2-8 for both models. The inclusion of spatial data for link 434 increases the autocorrelation to within significance levels for a significant proportion of time lags.

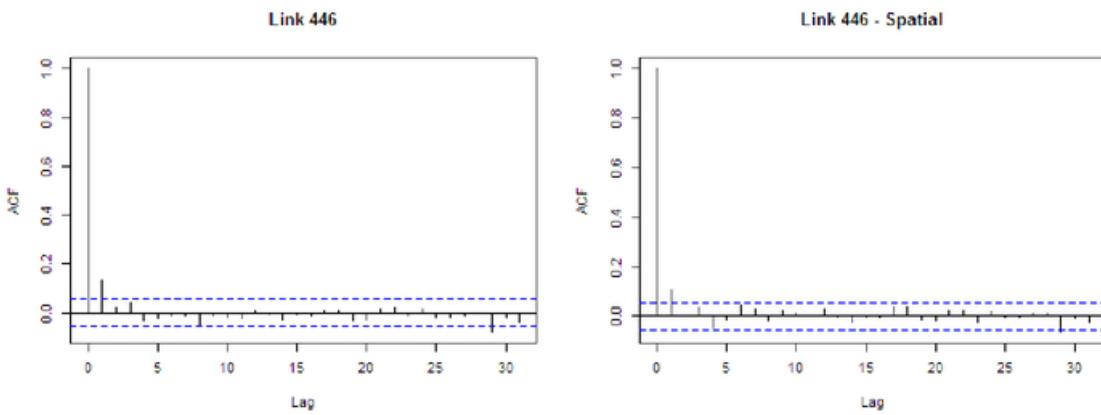


Figure 43: ACF plots of the residuals for non-spatial (left) and spatial (right) model travel time predictions for link 446. Blue lines indicated 95% confidence limits. Some significant autocorrelation is still apparent in both models, however there is less autocorrelation seen when using the spatial model.

4. Discussion

ARIMA found to be the most computationally efficient method, with models taking under a minute to process. In comparison, SVR models often took 10 minutes or more to calculate. This difference in time is offset by the more efficient SVR modelling process, since the algorithms self-optimise over a set of parameters and no initial trend analysis is required. Little input is therefore required from the user, except for interpretation of training error graphs, and the subsequent tuning of parameters, that takes no more than a minute. ARIMA, however, is relatively time-intensive 21 to the pre-analysis required to determine model parameters before the model can be created.

The primary challenge of ARIMA implementation was choosing the correct parameters for each model, as human interpretation is required, which could have been a factor in the quality of the travel time models produced. In the pre-processing stage, the data for each link was not made entirely stationary before the models were created. Differencing was unable to remove the cyclical trends in the dataset, and only reduced the autocorrelation to statistically insignificant levels after 3,000 to 4,000 lags (10 to 13 days). This is likely to have contributed significantly to the poor fit of the resulting models, as ARIMA relies upon a stationary dataset.

Cheng & Wang (2011) note that differencing can be ineffective in removing spatial nonlinear and non-stationary trends in space-time series. It was also found that SVR was unable to account for all autocorrelation, with some remaining present in the residuals. Cheng & Wang (2011) utilised ANN to initially remove global spatio-temporal trends from average annual temperatures across China, before using a Space-Time ARIMA (STARIMA) to create a predictive model. Improvements could therefore likely be made to this experiment through combining SVR with ARIMA.

23

The ability to interpret the results for both methods was found to be good, with predictions easily visualised against observed values using line graphs. Although *caret* does not contain diagnostic tools for analysis of residuals, these were easily calculated and analysed through R's ACF and PACF tools. ARIMA has the ability to produce diagnostic plots that, if used in this analysis, could have provided insight into which model parameters required further tuning.

For the links analysed, SVR produces good predictive models. It can predict noise to some extent, however it struggled with excessive noise and unusually large or small travel times, such as the morning and afternoon rush hours, with model performance varying considerably across the modelled links; more errors are observed for links with more highly variable travel times, such as 434, 2097, and 417. None of these links appear to have any obvious similarities in location, however without having visited the location in question it is difficult to determine why these links are more variable than others. Traffic on these road segments could be controlled by traffic lights with variable cycle times, which would likely cause irregularities in the travel time observations throughout the day. The consistent underreporting of travel times could be due to the faster travel times observed during the first few weeks of January. Performing the same analysis for a different month of the year could prove more successful.

Given that ARIMA failed to successfully forecast any of the travel times, it is difficult to directly compare the performance of the methods across the study area, so it is not known how well it would have predicted noisier datasets.

Finally, travel times are not normalised, and so do not account for variation in speed limits, number of lanes, or other fixed factors affecting travel. Normalisation before predictive modelling could aid in the pinpointing of real areas of congestion, rather than highlighting sections of roads that are always going to result in slower travel times due to external factors.

5. Conclusions and Further Work

In conclusion, distinct spatial and temporal patterns in travel times exist across the links analysed in this report, however the travel times observed are not necessarily directly influenced by the travel times on adjoining roads. SVR was found to be successful at creating predictive models of these travel times, both spatially and non-spatially, and allowed optimal parameters to be determined efficiently via computational methods. ARIMA proved more difficult to utilise, as it depended on the correct interpretation of initial analyses.

Processing time and resources were major limiting factors in this study. If this project were to be extended, one large spatial model should be created for all links using both SVR and STARIMA, a spatio-temporal version of ARIMA, and the residuals plotted spatially and temporally and compared, in a plot similar to that seen in Figure 8. STARIMA models of the same links as modelled using ARIMA would also give insight into the effect of including space-time data on the models.

Residual analysis should also be carried out for the ARIMA method to aid in the parameterisation of the models.

References

- Cheng, T. & Wang, J., 2011. A Hybrid Framework for Space-Time Modeling of Environmental Data. *Geographical Analysis*, Volume 43, pp. 188-210.
- Lau, K. & Wu, Q., 2008. Local prediction of non-linear time series using support vector regression. *Pattern Recognition*, 41(5), p. 1539–1547.
- Peng, Z., Wu, F. & Jiang, Z.-Y., 2009. Prediction of Railway Passenger Traffic Volume by means of LS-SVM. *LNCS*, Volume 3973, pp. 8-14.
- Vanajakshi, L. & Rilett, L. R., 2004. *A Comparison Of The Performance Of Artificial Neural Networks And Support Vector Machines For The Prediction Of Traffic Speed*. Parma, IEEE Intelligent Vehicles Symposium.
- Vapnik, V., 1999. *The Nature of Statistical Learning Theory*. 2nd ed. New York: Springer.
- Wu, C.-H., Ho, J.-M. & Lee, D., 2004. Travel-Time Prediction With Support Vector Regression. *IEEE TRANSACTIONS ON INTELLIGENT TRANSPORTATION SYSTEMS*, 5(4), pp. 276-281.
- Zhang, W., Zhang, H. Z. H., Chen, G. & Wei, Y., 2012. A web partition algorithm based on support vector machine. *Przegląd Elektrotechniczny (Electrical Review)*, Volume 88, pp. 31-33.

Appendix A – Statistical Summary of Travel Times per Road Link

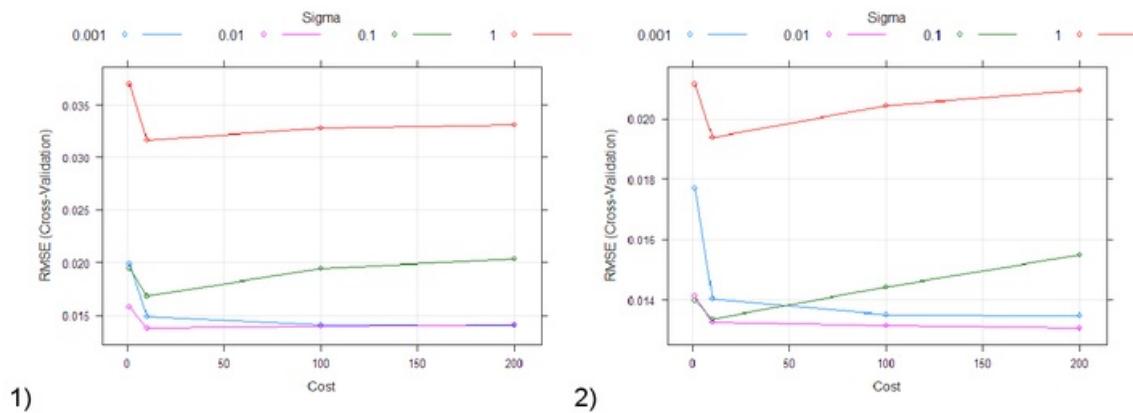
LCAP_ID	Mean	Standard Deviation	Median	Minimum	Maximum	Range	Skew	Kurtosis
1745	0.1195	0.0704	0.1035	0.0645	0.8597	0.7951	5.1472	32.666
1747	0.1368	0.0685	0.1158	0.0703	0.6421	0.5718	3.5667	16.004
417	0.2194	0.0469	0.2174	0.0861	0.6553	0.5692	0.9507	5.840
2472	0.1041	0.0986	0.0745	0.0300	1.0860	1.0560	3.3848	14.983
1877	0.0832	0.0571	0.0685	0.0300	0.8615	0.8315	5.6001	56.120
1878	0.1901	0.1027	0.1586	0.0723	1.2198	1.1474	2.8430	11.791
2468	0.1502	0.0453	0.1415	0.0697	0.8152	0.7456	2.9457	21.889
1604	0.3242	0.1178	0.3017	0.0950	0.9268	0.8317	1.1024	1.830
1613	0.1541	0.0644	0.1405	0.0663	0.8187	0.7525	4.4793	28.006
1882	0.2555	0.1624	0.1838	0.0785	1.0582	0.9796	1.5711	2.324
1620	0.2029	0.1221	0.1654	0.0662	1.0107	0.9445	2.7504	8.581
1622	0.1662	0.0587	0.1508	0.0675	0.6366	0.5691	1.8510	4.984
1407	0.2175	0.0890	0.1928	0.0953	1.1104	1.0150	2.7580	11.954
2260	0.1309	0.0397	0.1254	0.0394	0.4336	0.3942	1.5234	5.668
2261	0.2025	0.0379	0.2002	0.0943	0.4740	0.3797	0.7029	2.400
2357	0.1996	0.0870	0.1773	0.0690	1.0347	0.9657	2.7031	13.598
2416	0.1573	0.0763	0.1477	0.0321	1.0401	1.0081	3.2417	23.672
2358	0.1860	0.0440	0.1818	0.0628	0.4658	0.4029	0.7972	1.606
2402	0.1923	0.0391	0.1868	0.0910	0.4980	0.4070	1.4627	5.106
2097	0.2528	0.1189	0.2270	0.0739	1.3266	1.2528	2.3969	9.236
2318	0.1731	0.0651	0.1586	0.0793	0.6319	0.5526	2.2995	7.372
2344	0.1525	0.0485	0.1433	0.0596	0.6155	0.5558	2.7999	14.403
2282	0.1558	0.0515	0.1501	0.0654	0.8565	0.7912	5.9587	64.808
433	0.1856	0.0557	0.1788	0.0748	0.7218	0.6470	4.0474	27.290
434	0.1640	0.0385	0.1603	0.0680	0.5219	0.4539	1.4200	6.113
435	0.1099	0.0489	0.0990	0.0545	0.8252	0.7706	5.9126	53.379
556	0.1818	0.0408	0.1799	0.0860	0.5914	0.5053	1.7221	9.833
446	0.2498	0.2218	0.2066	0.0935	3.0692	2.9757	7.4290	66.020
425	0.1116	0.0482	0.1005	0.0551	1.5202	1.4651	11.3006	247.401
457	0.1888	0.0807	0.1700	0.0763	1.1417	1.0654	4.0326	27.665

Appendix B – Support Vector Regression Models for All Links (Non-Spatial)

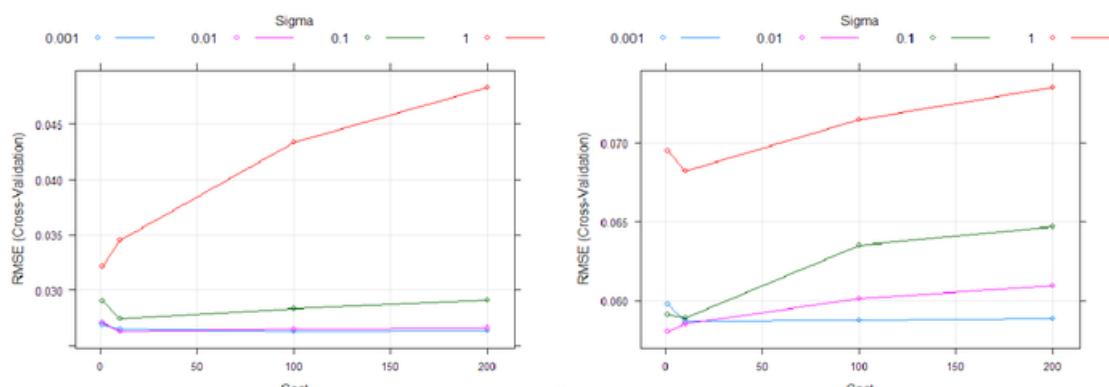
These plots show the training error for all initial SVR models created for all links. Figure numbers and corresponding link IDs are as follows:

Figure	Link ID	Figure	Link ID	Figure	Link ID
1	1745	11	1620	21	2318
2	1747	12	1622	22	2344
3	417	13	1407	23	2282
4	2472	14	2260	24	433
5	1877	15	2261	25	434
6	1878	16	2357	26	435
7	2468	17	2416	27	556
8	1604	18	2358	28	446
9	1613	19	2402	29	425
10	1882	20	2097	30	457

Please note that these do not necessarily show the optimal mode, since further optimisation was carried out after this stage.

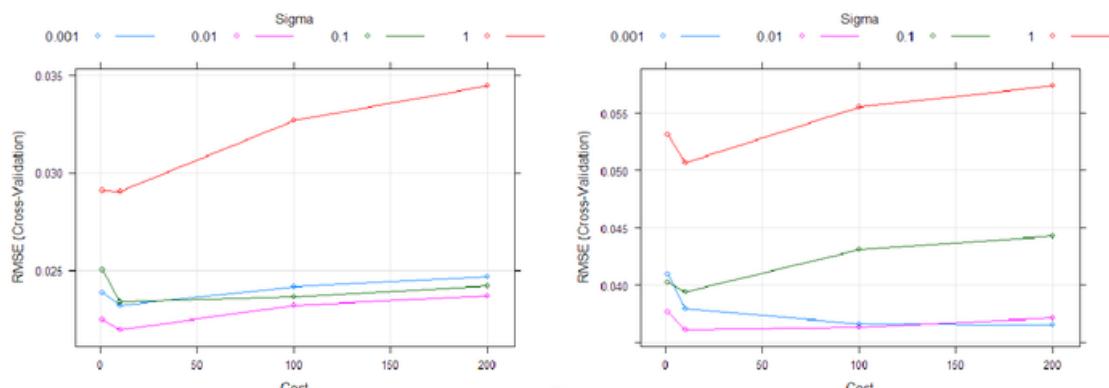


CEGEG076: Spatio-Temporal Analysis and Data Mining
Christopher Baxter and Lucille Ablett



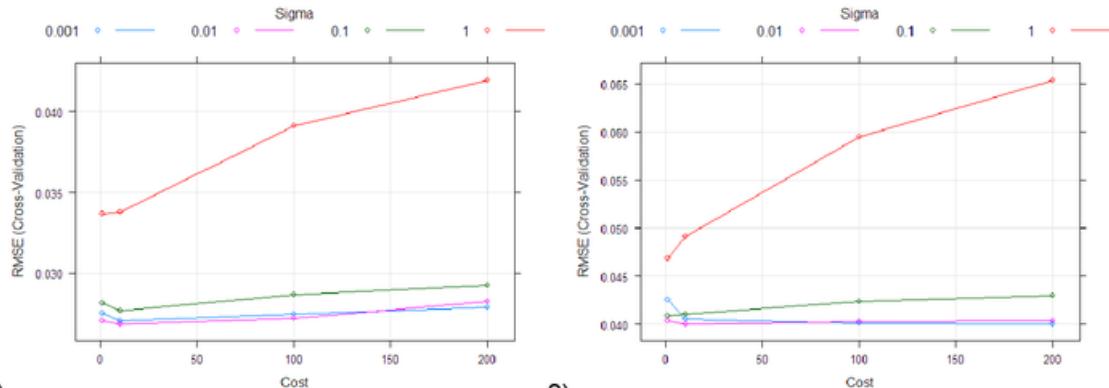
3)

4)



5)

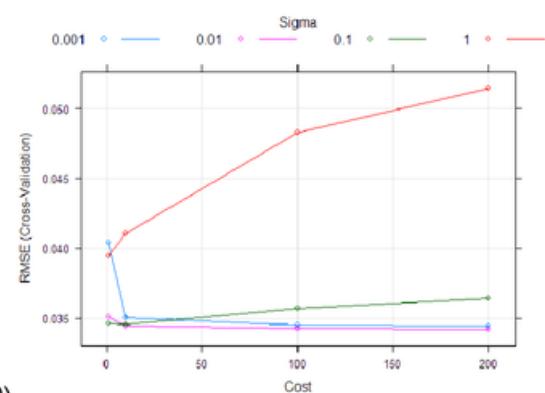
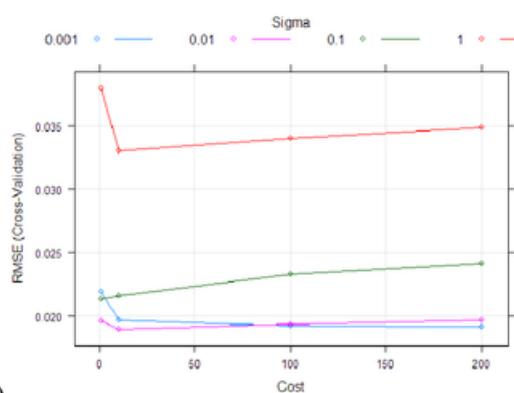
6)



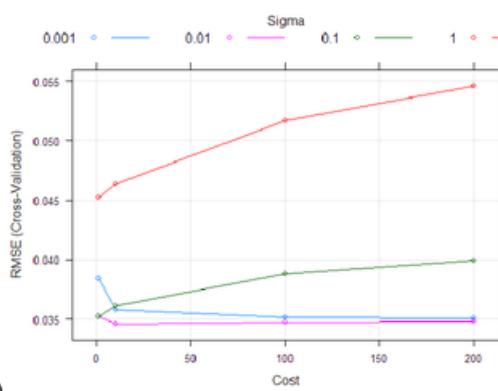
7)

8)

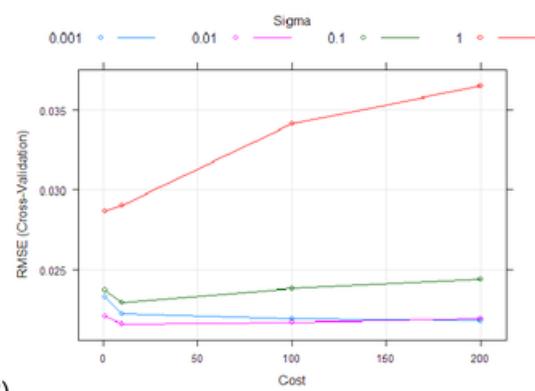
CEGEG076: Spatio-Temporal Analysis and Data Mining
Christopher Baxter and Lucille Ablett



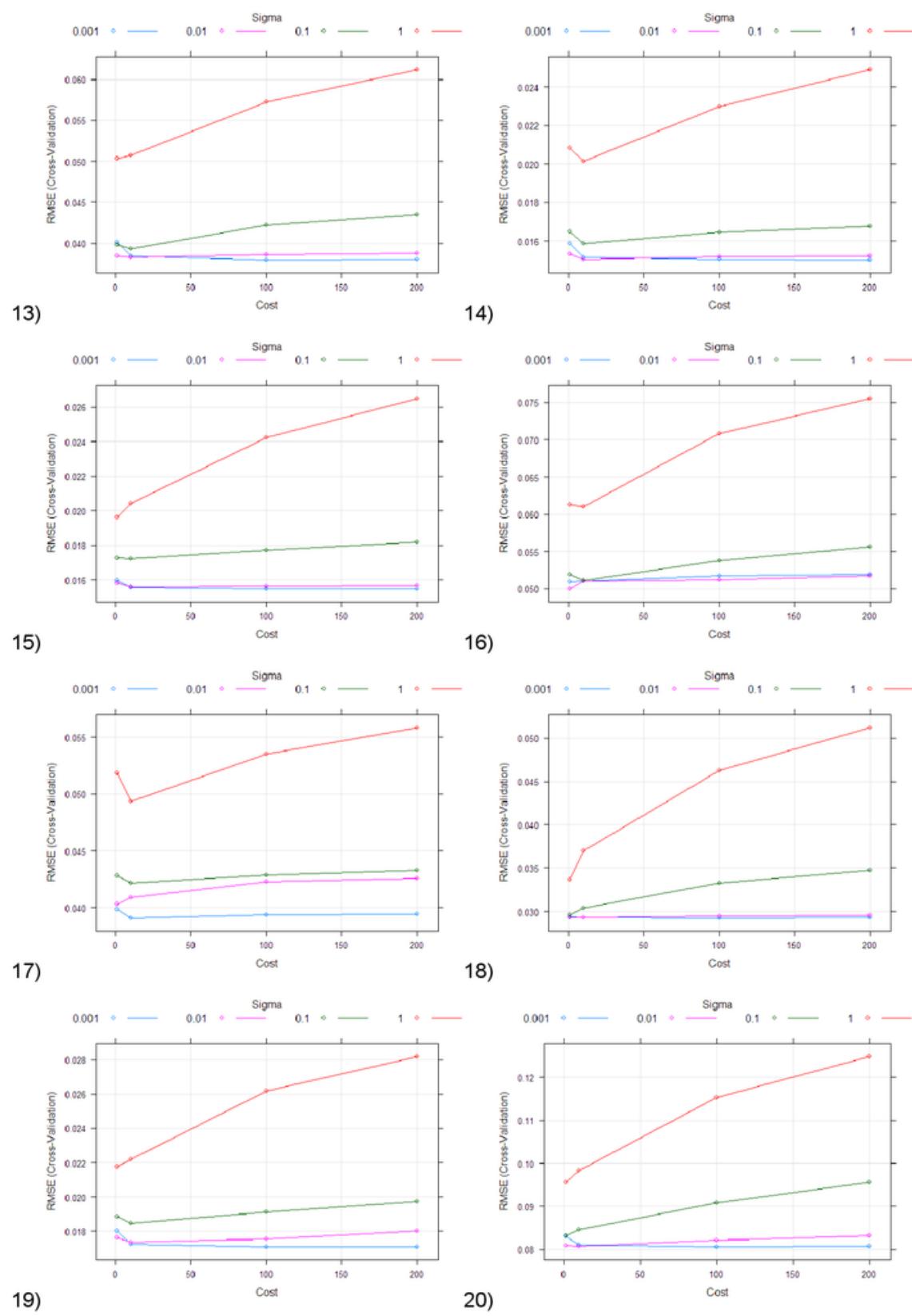
11)



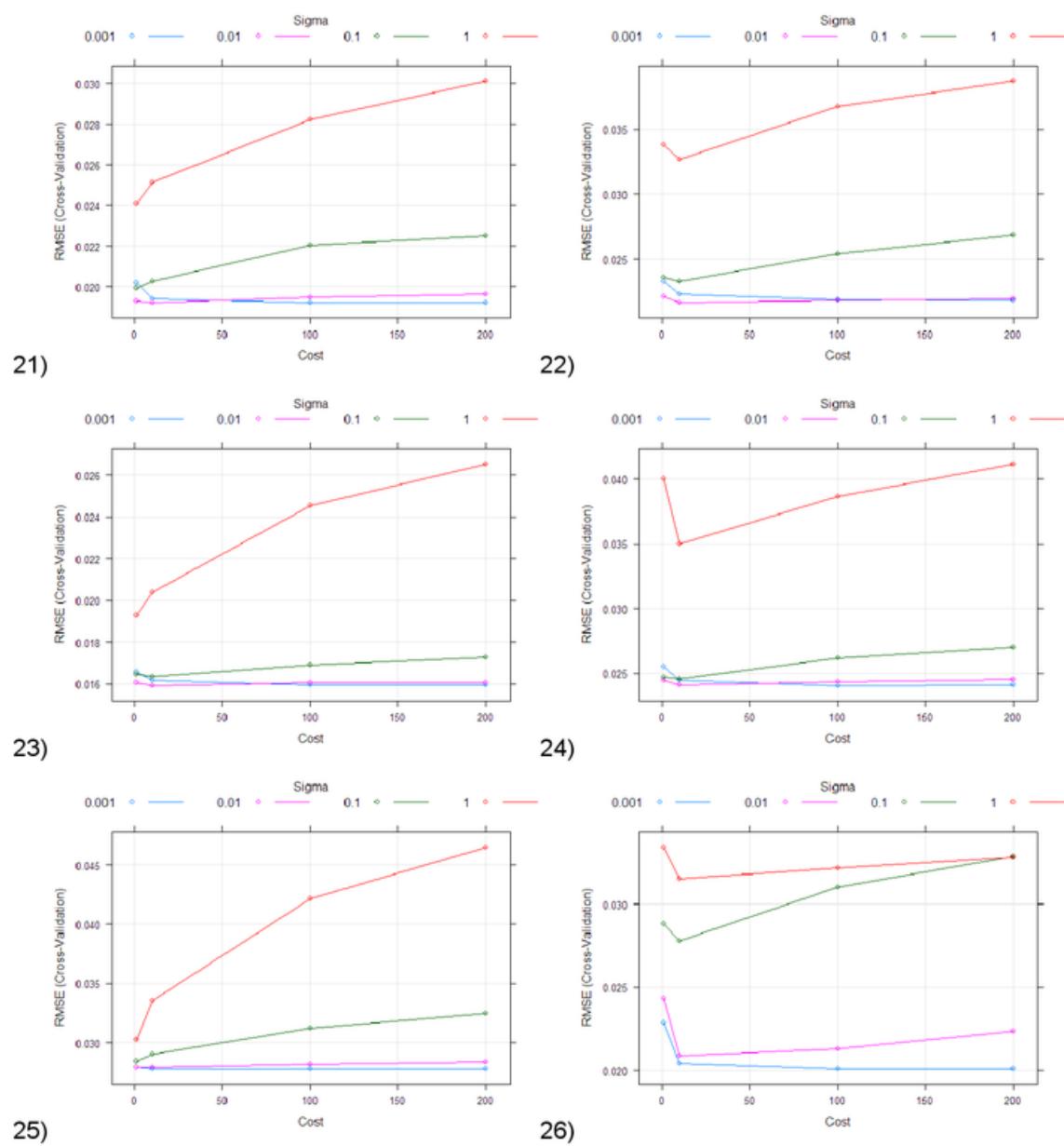
12)



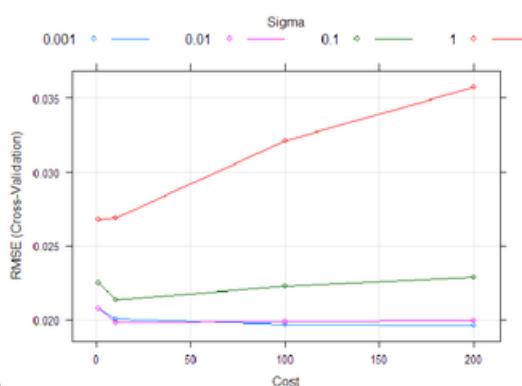
CEGEG076: Spatio-Temporal Analysis and Data Mining
Christopher Baxter and Lucille Ablett



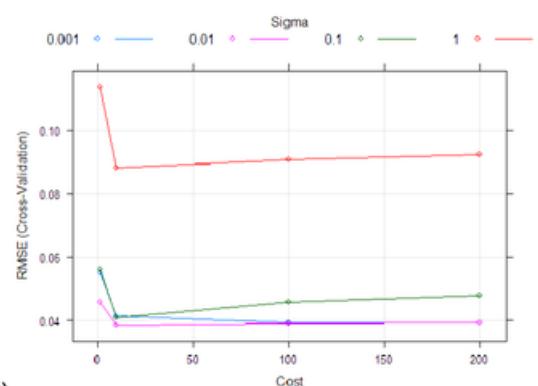
CEGEG076: Spatio-Temporal Analysis and Data Mining
Christopher Baxter and Lucille Ablett



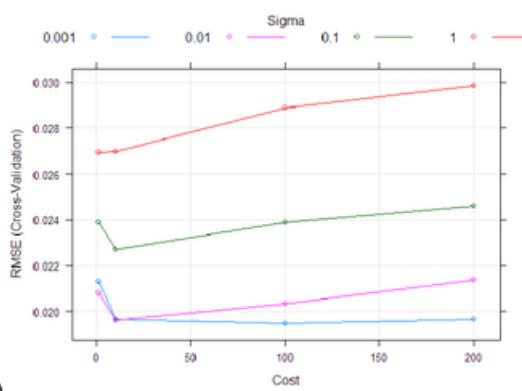
CEGEG076: Spatio-Temporal Analysis and Data Mining
Christopher Baxter and Lucille Ablett



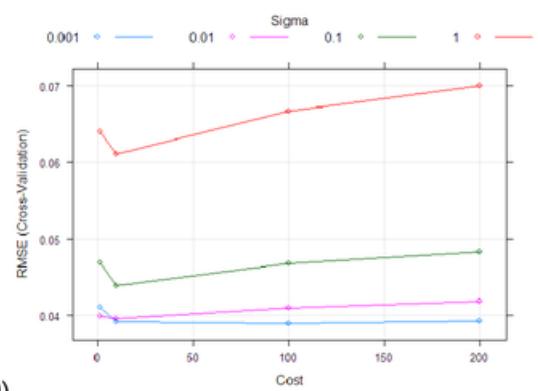
27)



28)



29)



30)

Forecasting Urban Travel Times in London

GRADEMARK REPORT

FINAL GRADE

GENERAL COMMENTS

78 /100

Instructor

Overall this is an excellent report. The introduction to the methods is good and shows understanding. Some more context on traffic congestion in London would improve it further.

The data description and exploratory data analysis are both excellent. You have used a wide range of indices to analyse the data and your interpretations are thoughtful.

The section on SVR is very thorough and demonstrates a lot of work. Your strategy for model training has allowed you to analyse the factors that affect the results. The results themselves are well presented. You have misinterpreted the main computational burden of kernel methods, which is the number of data samples not the number of parameters (embedding dimension). See the inline comments for more on that.

The discussion and conclusions links well to the context and literature and offer some good suggestions for model improvements.

Your use of base R plots and other libraries such as lattice is very good, particularly in your exploratory data analysis section.

Overall, a very good report, well done.



Comment 1

A good overview of the chosen methods that shows understanding. It could be improved by introducing the context of urban traffic congestion to further motivate the study.



Comment 2

This is a very good figure that is well explained in the text.



Comment 3

This title is incorrect.



Comment 4

Nice, clear figures with legend, scale bar etc.



Comment 5

Some interesting interpretations.



Comment 6

Again, a clear, well presented figure.



Comment 7

Good interpretation. Perhaps you could look at the road on Google maps? Although the layout may have changed since 2011.



Comment 8

You can find out.



Comment 9

Good use of lattice plot. The data display is clear.



Comment 10

As you are talking about the context of the roads, it would be useful to have a basemap in one of your figures, e.g. Fig 2.

PAGE 10

PAGE 11



Comment 11

This is a nice way to show the rise and fall in congestion throughout the day.

PAGE 12

PAGE 13



Comment 12

Again, good use of lattice plots. It would be good if you could find a way to put the confidence intervals on here, although it's not necessary for a purely visual analysis.

PAGE 14



Comment 13

Nice illustration of the decreasing dependence between observations as time separation increases.

PAGE 15

PAGE 16

PAGE 17



Comment 14

Well done for producing this for the whole network. This is a good figure.

PAGE 18

PAGE 19

PAGE 20

PAGE 21

PAGE 22

PAGE 23

PAGE 24



Comment 15

This shouldn't be the case for SVR. Number of observations increases the processing time.



Comment 16

Why 2-fold?

PAGE 25



Comment 17

Number of support vectors is determined by the optimisation process. If you have a larger C then the number of support vectors will generally increase, but it depends on the other parameters too, e.g. if the tube size is very small then you are likely to have more support vectors as it needs to 'wrap round' the data points.

PAGE 26



Comment 18

This is a negligible increase in processing time because the important factor is the size of the kernel, which is n^2 , where n is the number of data points. This is one of the main advantages of kernel methods; you can have potentially thousands of variables without a significant increase in processing time. The slight increase may be due to the optimisation algorithm taking slightly longer to find a solution.

PAGE 27

PAGE 28



Comment 19

A very comprehensive set of results.



Comment 20

These figures are a bit too small to be easily interpreted. If using figures of this size, it might be better to combine lines (e.g. for observed) and points (e.g. for predictions).

PAGE 29

PAGE 30

PAGE 31

PAGE 32

PAGE 33



Comment 21

Proof-read



Comment 22

A good summary of the computational vs human burden of the training process.



Comment 23

A good suggestion.

PAGE 34



Comment 24

Some good interpretations.



Comment 25

They are normalised by length, but you are right that the performance is not relative to the nature of the road. Good observation.



Comment 26

Check for typos.

PAGE 35

PAGE 36

PAGE 37

PAGE 38

PAGE 39

PAGE 40

PAGE 41

PAGE 42
