# University of BRISTOL

Engineering Mathematics

Mathematical Data Modelling 3 (EMAT30005)

# Forecasting Bristol's Air Pollution Using Bayesian Networks

*Authors*
Alfred Brown
Andrew Corrigan
Jake Farren-Price
Luke Moorcroft
Jianing Wang

*Supervisors*
Filippo Simini
Eddie Wilson

**Abstract**

Bristol has been in breach of the current air quality standards for years with significant research showing that these levels cause serious harm to humans [1]. This project uses statistical analysis in the form of Bayesian Networks to predict the levels of $NO_2$ on an hourly basis, using collected wind speed, temperature and traffic data at three different locations in Bristol. Our findings suggest that traffic flow alone does not dictate the levels of $NO_2$ accurately, and in order for nitrogen dioxide levels to exceed the legal limit, wind speed and temperature must be low and traffic flow high. Our report aims to provide a foundation for further predictive models to be built upon and ultimately enable the city of Bristol to take more action towards limiting its dangerous pollution levels.

# 1    Introduction

Health warnings have been issued over parts of the United Kingdom over the last decade during episodes of high pollution [2]. Bristol's air quality has been in breach of the current nitrogen dioxide regulations set by the European Union for a number of years. It was recently estimated that the current annual health costs due to the level of air pollution are around £16 billion [3].

These levels are largely due to emissions from motor vehicles, especially diesel combustion engines. Road traffic contributes to around 40% of the total nitrogen dioxide levels [4]. Having the ability to predict these levels of harmful air pollutants would be advantageous as it could be used to advise people who are at risk and potentially prevent an episode of high pollution.

In this paper Bayesian Networks are used to explore how the factors wind speed, temperature and traffic at three different locations interact in order to increase or reduce pollution levels in Bristol. A years worth of data is used with the majority used for training the model. The remaining data is used to validate the accuracy of the model. Given the random nature of these factors throughout a city, it is appropriate to use Bayesian Networks as these are able to capture uncertainty in the form of probability distributions. Preliminary research showed there is not a simple linear relationship between the three factors (see Figure 6) which further inspired the use of Bayesian Networks as these do not need to fit any functions to data. Furthermore, Bayesian Networks allow one to easily model complex dependencies between different factors.

# 2    Data Acquisition

The year of 2014 was chosen as the year to take the data from as it is the latest year in which all factors, except traffic, had accessible data sets.

## 2.1    NO$_2$ Pollutant Levels

Hourly nitrogen dioxide pollutant data was collected from three locations across Bristol; the A37, Brislington and St Pauls, sourced from Open Data Bristol [5]. The data sets available all had 851 missing data points. To overcome this issue the mean average of the given time three days before and three days after was inserted instead.

## 2.2    Traffic

The traffic data from Bristol City Council was different to the rest of the data sets, as the data was recorded every five minutes and only the week of 01/04/2019 was available. With 12 five minute intervals in an hour, having filled in the 116 missing values, the grouping of the data was simple. Due to the lack of data, it was necessary to extrapolate this week for a years length. To do this, this weeks data was replicated 52 times, with the off school term time weeks multiplied with a scale factor of 0.97 [6]. Finally to make the traffic data more

applicable to the model, it was necessary to rescale the data as the volume of motor vehicles in Bristol has increased. According to data from the Department of Transport, between 2014 and 2017, the number of motor vehicles on the road has increased on average by 1.019% per year [7]. To implement this each data point was multiplied by a scale factor of $(\frac{1}{1.019})^5 = 0.91$.

The three sites at which data was collected was very close to where the $NO_2$ was recorded, being along the same roads of; Well's Road for the A37, St Philip's Causeway for Brislington and Ashley Road for St Pauls.

## 2.3   Wind Speed

The wind speed data was again taken from the government's Department of Environment, Food and Rural affairs [8], and again this data set had an issue with missing data. This set was 336 hours less than the yearly length of 8760 hours. Having found the position of the absent data points, again they were replaced with the average of the given time 3 days either side. Although wind speed is taken as a city wide variable, these observations were made in St Pauls, suggesting that the model presented further on will predict St Pauls nitrogen dioxide levels more accurately compared to Brislington and the A37.

## 2.4   Temperature Flow

The temperature data was collected from Open Data Bristol [9], and as with all the other data, the 16 absent data points were filled using the average of the neighbouring days. The measurement point for this data is located on the roof of the CREATE centre in Bristol.

## 2.5   Assumptions

Taking this data and using it later in the models, the following assumptions are made:

- Due to the limited availability of traffic data, it is assumed that the week recorded has a representative distribution of traffic across the times of each day for the rest of the year. This is a reasonable assumption to make as the most influential factor in the changing of traffic flow is school term time [6], which affects whole weeks and this is accounted for.

- The year 2019 traffic flow data is assumed valid to be used as 2014 data. This has been accounted for as best as possible by scaling the 2019 data.

- Wind speed and temperature, having been collected in only one location, are considered to be consistent throughout the city at any one time. This assumption is made under the presumption that meteorological conditions are more consistent within smaller areas of land, and Bristol, being roughly 110km$^2$, can be considered small enough for this to hold.

## 2.6 Data Processing and Thresholds

Bayesian Networks are computationally very expensive when working with multiple probability distributions whose joint and conditional distributions together are unknown [10]. What is clearly seen in the data plots of the traffic in all three locations in figures 15,16 and 17 is that none of the traffic data seems to be normally distributed. This was verified by performing a Kolmogorov-Smirnov test on each of the locations. The null hypothesis was that this data came from a standard normal distribution, and this was rejected every time. The same happened with the wind speed data with a normal distribution and even with a Weibull distribution, which is what would be expected [11]. With no distribution applicable to the input factors the conditional probabilities become too difficult to compute. To overcome this the data was converted into binary, being either above or below a given threshold.

The following thresholds were used to split the data into binary values:

- **$NO_2$ Level**: The evaluation criteria for $NO_2$ is defined in two air quality standards: The 'hourly standard', which is defined as the concentration of $NO_2$ averaged over a period of one hour, and the 'annual standard' defined as the concentration of $NO_2$ averaged over a period of a year. The European Union (EU) and the UK Air Quality Strategy has developed legislation with the aim of limiting air pollutants. For the hourly standard, the concentration objective is $200\mu\mathrm{gm}^{-3}$ with 18 exceedances per year [12]. For the annual standard, the concentration objective is $40\mu\mathrm{gm}^{-3}$ [13]. We will be focusing on hourly data so our initial threshold will be $40\mu\mathrm{gm}^{-3}$. Anything above this will be considered high and anything below this will be considered low.

- **Wind Speed**: The Beaufort Wind Force Scale [14] is used to determine a threshold for wind speed. The pivot wind speed chosen for our data is a force 3 on the Beaufort Wind Scale (Gentle breeze), which corresponds to wind speeds of $4ms^{-1}$. While this is not considered 'strong wind' according to the scale, dust particulates become afloat and disperse at this speed, and thus, we assume that $NO_2$ will be displaced too.

- **Traffic**: The threshold for traffic levels are difficult to decide upon as different streets have varying levels of traffic and different cities will have different averages. It was decided to simply use the average of the data for each of the locations.

- **Temperature**: Like traffic, it was decided to use the average as the threshold for the temperature data.

# 3 The NO₂ Pollutant Model
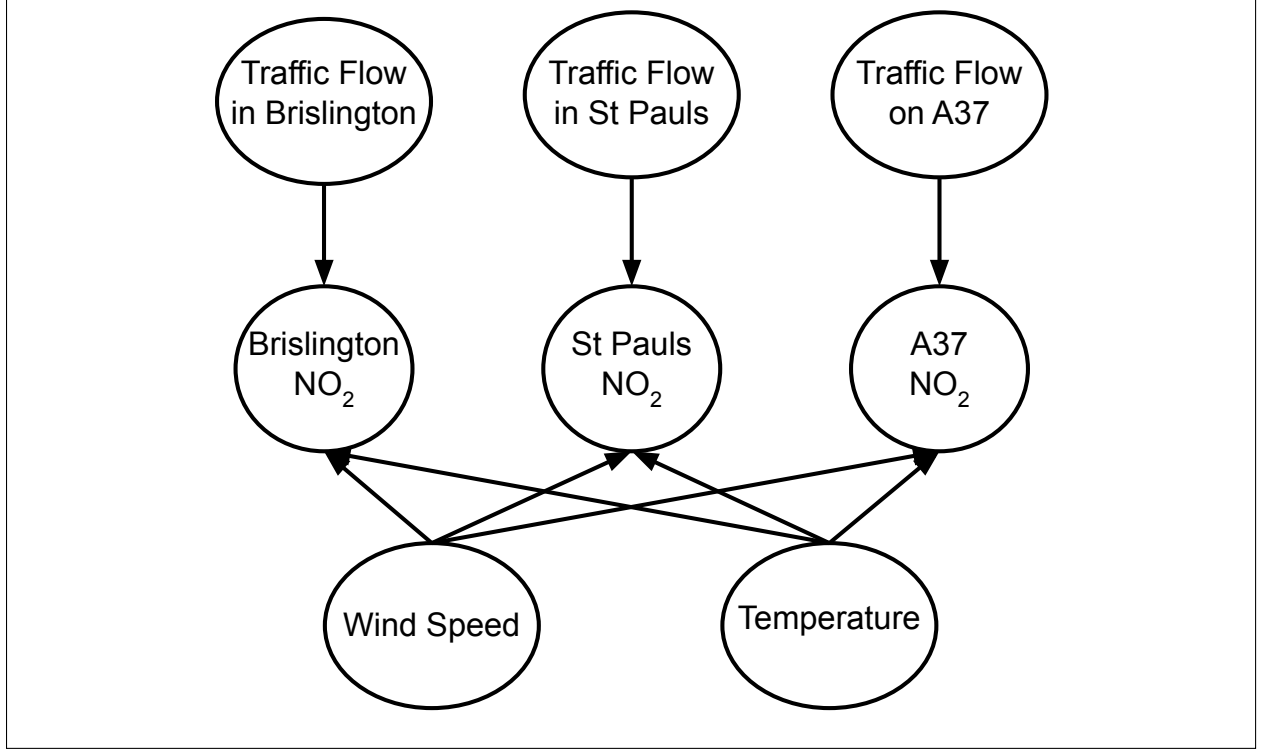
## 3.1 Bayesian Network Graph



Figure 1: The Pollution Bayesian Network for $NO_2$ levels in three different locations across Bristol. The nodes represent random variables and the directed edges represent dependencies.

## 3.2 Determining the Joint Distribution

Given the network presented above, the joint the distribution of all the random variables can be determined. The formula for determining the joint distribution of a Bayesian network is

$$P(x_1, x_2, ..., x_n) = \prod_{i=1}^{n} P(x_i | \text{Parents}(x_i)),$$

where Parents($x_i$) are the random variables (nodes) that $x_i$ depends on, i.e., the nodes that feed into it. Thus, the joint distribution for our model is expressed as

$$P(B_T, S_T, A_T, B_N, S_N, A_N, W, T) = P(W)P(T) \sum_{X \in \{B,S,A\}} P(X_T) \sum_{X \in \{B,S,A\}} P(X_N | X_T, W, T). \quad (1)$$

where $B_T, S_T, A_T$ and $B_N, S_N, A_N$ are the traffic flow levels and $NO_2$ levels in Brislington, St Pauls and along the A37, respectively, and $W$ is the wind speed and $T$ is the temperature across Bristol. For the sake of simplicity, we merge the three locations into one for the

4

following sections. Let $V$ (for vehicles) and $N$ represent the combined traffic flow and $NO_2$ levels, respectively. Then we can rewrite the joint distribution as

$$P(N, T, V, W) = P(N|T, V, W)P(T)P(V)P(W). \tag{2}$$

## 3.3   Marginalisation and Predicted NO$_2$ Distributions

Given the simplified joint distribution we are now in the position to find the probability distribution for $NO_2$. By marginalising out the factors, we can obtain the marginal distribution for $NO_2$ in the form

$$\begin{aligned}
P(NO_2) &= \sum_{T \in \{0,1\}} \sum_{V \in \{0,1\}} \sum_{W \in \{0,1\}} P(N0_2, T, V, W), \\
&= \sum_{T \in \{0,1\}} \sum_{V \in \{0,1\}} \sum_{W \in \{0,1\}} P(N0_2|T, V, W)P(T)P(V)P(W). \tag{3}
\end{aligned}$$

The exact same process shown above is applied to our model for each location. The distributions for $P(B_N), P(S_N)$ and $P(A_N)$ now contain all the information from the traffic flow data for each of the locations and the wind speed and temperature data.

In practise, the marginalisation was computed using a brute force method [15]. This simply meant iterating over all possible combinations of the random variables within the conditional and marginal distributions. (Using this technique, it becomes very apparent how computationally expensive Bayesian Networks are and why it was necessary to threshold the data into binary values.)

## 3.4   Other Inferences

The following gives a simplified example on how we used our Bayesian Network model to perform inferences other than just finding the probability distributions for nitrogen dioxide levels at the three locations. Our model lets us explore questions like: Given that the wind is strong and traffic high, how likely is it that the nitrogen dioxide levels are also high? We can answer this using the brute force method to compute

$$P(NO_2{=}1|V{=}1, W{=}1) = \frac{\sum_{T \in \{0,1\}} P(NO_2{=}1|T, V{=}1, W{=}1)P(T)P(V{=}1)P(W{=}1)}{\sum_{T \in \{0,1\}} \sum_{NO_2 \in \{0,1\}} P(NO_2|T, V{=}1, W{=}1)P(T)P(V{=}1)P(W{=}1)}.$$

The model uses training data to find the marginal and conditional probabilities of all the random variables, so that when we trial it with test data, it computes the probability of $NO_2$ being high *given* whatever the binary values for temperature, wind speed and traffic for each of the locations are for every hour. In the following, we refer to the $NO_2$ pollutant model we just presented as the 'Simple Model' in order to avoid confusion with a more complex model.
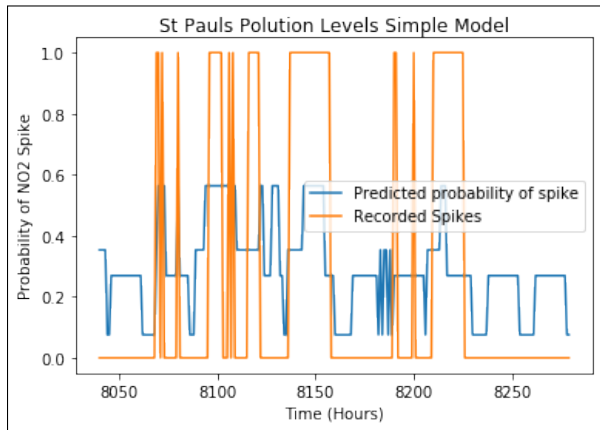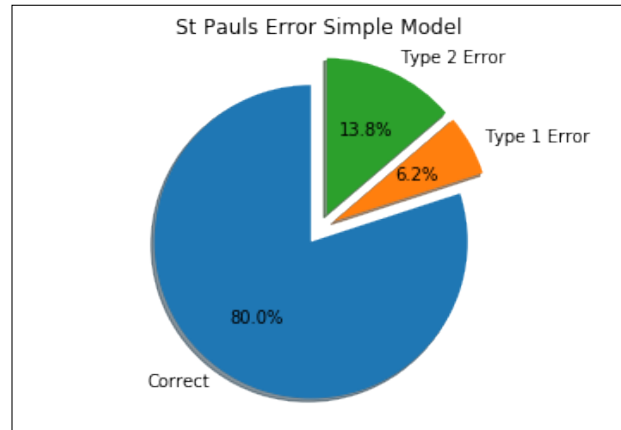
# 4 Results

## 4.1 Training and Testing

In order to test the accuracy of our developed model the data has been split into training and test data subsets. Using the training data to develop the model, the robustness of the model can then be tested using the unseen test data. The test data consists of 2.7% of the whole data set. This gives us a test period of ten days. This seems reasonable as that is typically the furthest reliable weather prediction for when the model is used in practise.

## 4.2 Simple Model

The Bayesian Network outputs a probability given certain evidence. There are multiple ways for this probability to be interpreted and used for predictions. The first way the output of the model has been assessed is by comparing it directly with the the $NO_2$ data it is attempting to replicate. Figure 2a shows the probability output (in blue) graphed against the recorded $NO_2$ spikes (in Orange). This shows that the model predicts the majority of the spikes recorded in the model. To quantify this further the probability output of the Bayesian network is passed through a threshold, whereby over a certain probability high nitrogen dioxide is predicted and under, low $NO_2$ is predicted. This means the output of the network can be directly compared to the recorded $NO_2$ spikes and a percentage accuracy can be computed. This is graphed in Figure 2b. This also shows the type of error associated with the model. Similar results for the two other locations can be found in the appendix (see figures 9 10).



(a) The predicted levels of $NO_2$ vs the recorded spikes in $NO_2$ pollution in St Pauls.

(b) A pie chart of the prediction error for St Pauls.

With an accuracy of 80% the model seems accurate however due to the nature of the nitrogen dioxide spikes in this time frame, predicting 0 spikes would be an accurate model. A better measure is the type of error. In our case type 1 error would be predicting an $NO_2$ spike where one did not occur. Conversely a type 2 error would not be predicting an $NO_2$ spike where one had occurred. Whilst all error is undesirable type 2 error is much worse in our case. Missing a large $NO_2$ spike could be problematic. This means that a type 2 error of 13.8% requires improvement.

## 4.3 Complex Model

Bayesian Networks can be easily extended to include more links and more nodes with the only cost being processing power. This means that exploring new links within the same data set is an easily implemented extension. Thus the traffic data from one location has been connected to the $NO_2$ data of the other outputs. This has linked all of the locations together so that they are dependent on one another. This gives the following network shown in Figure 3.
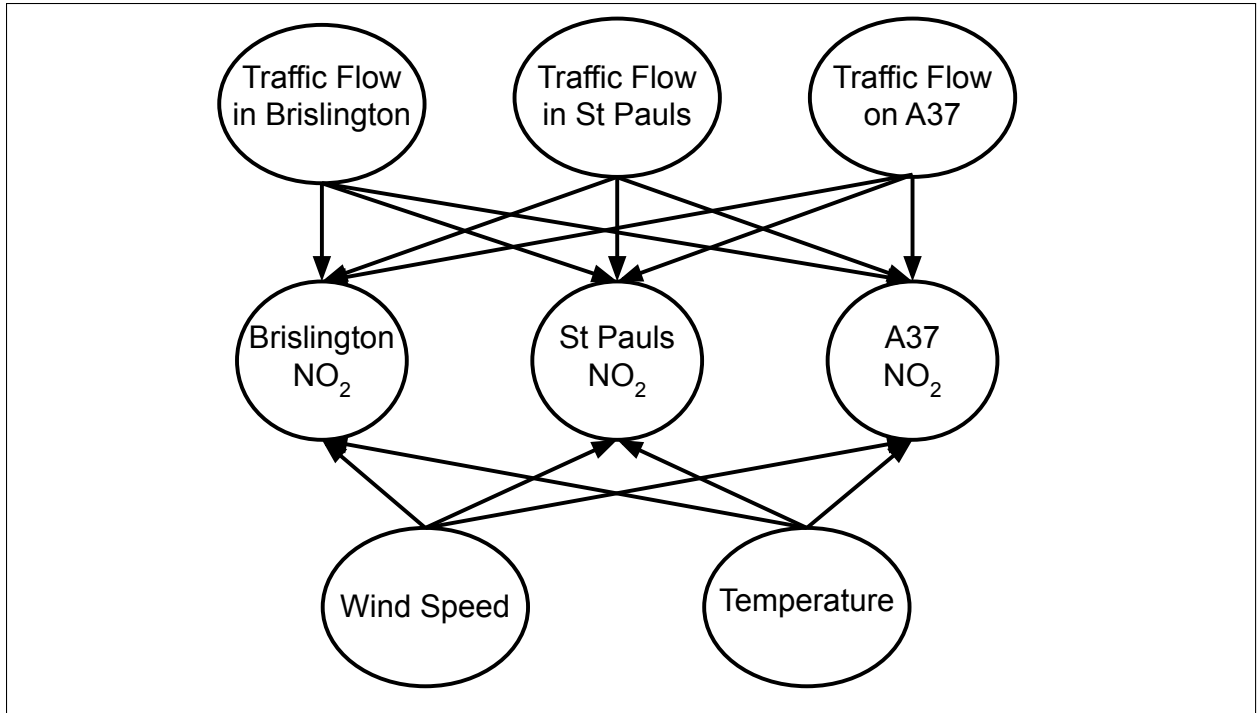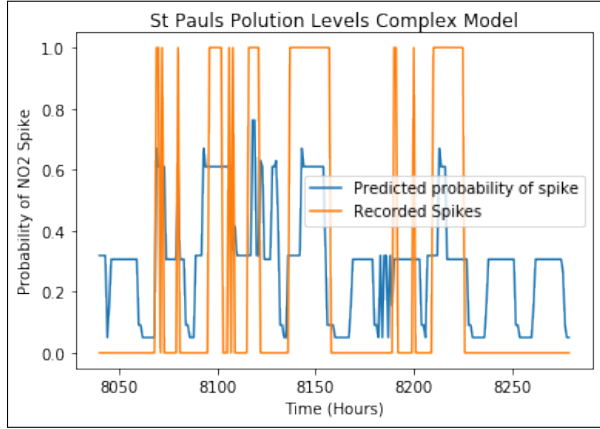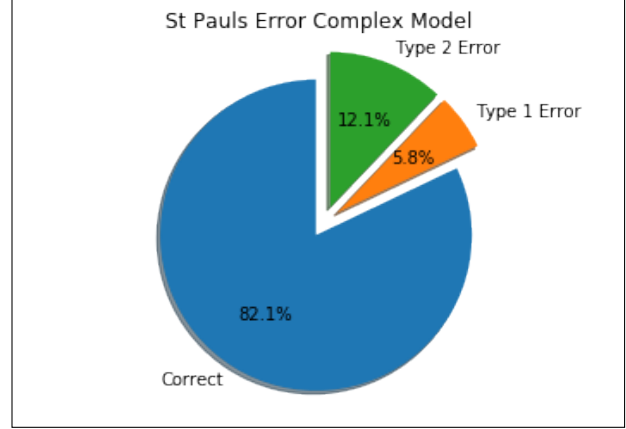


Figure 3: A more complex Bayesian Network. This time the $NO_2$ levels depend on all of the locations.

From the new network further predictions can be made for the same training and test sets, this means the models can be easily compared and improvements quantified. The same graphs generated for the simple model have also been generated for the complex. The probability output of the model against the recorded spikes can be seen in Figure 4a and the error in Figure 4b.

(a) The predicted levels of pollution of the more complex model vs the recorded spikes in pollution for St Pauls.
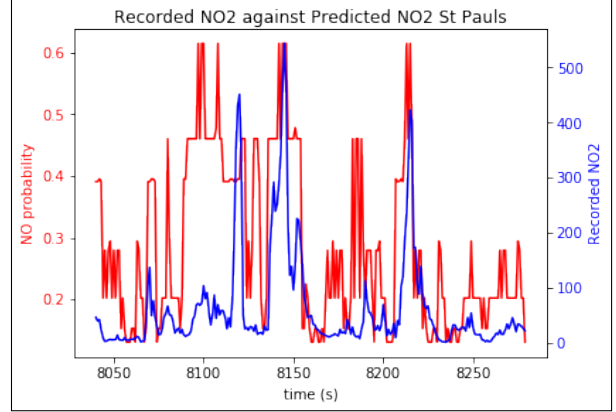


(b) A pie chart of the prediction error for St Pauls.

Figure 4: St Pauls predicted $NO_2$ levels using the complex model and its errors.

Both models predict spikes at similar times, however due to the increased complexity of the network the variation in the prediction has also increased. This has resulted in a reduction of error by 2.1% for the St Paul's Location with the type 2 error reducing to 12.1%. Whilst not offering a huge improvement for predicting binary $NO_2$ spikes, the increase in variation of the model is more useful for predicting real world $NO_2$ data.

In the following graphs, the probability output of the simple and complex models have both been graphed against the recorded $NO_2$. Figures 5a and 5b show that it is clear that the more complicated model gives a more usable output. In the first model it predicts the spikes, but cannot differentiate between the much larger spikes and smaller spikes. This is particularly useful as the model is intended to predict the spikes that can cause health damages so that they can be averted. Therefore quantifying predictions is essential.

(a) The recorded $NO_2$ data (blue) against the probability output of the simple model (red).

(b) The graph shows recorded $NO_2$ data (blue) against the probability output of the complex model (red).

Figure 5: Recorded $NO_2$ levels against probability output of the Bayesian Network for the simple and complex model.

# 5 Discussion of Results

## 5.1 Sensitivity of Thresholds

Comparing the graphs in Figure 7, it is apparent that the decision of where to put the threshold is vital. Both type 2 and type 1 errors increase with just a small adjustment of wind speed. For example, with the threshold for the $NO_2$ level at $40\mu\text{gm}^{-3}$ (see Section 2.6) the accuracy of the model is 78.3%. This accuracy can be improved to 81.7% by increasing the $NO_2$ threshold to $60\mu\text{gm}^{-3}$. This seems reasonable as the model would only be required to predict the more extreme episodes of high $NO_2$ levels. Adjusting the wind speed threshold also has a major effect. Using the wind speed threshold of $10ms^{-1}$ gives us an accuracy of 63.7%. Decreasing this threshold to $4ms^{-1}$, increases the accuracy to 81.7%.

## 5.2 Error Types

Limiting one type of error was more important to us than the other. Type 2 error, predicts low levels of $NO_2$ when in fact the actual levels are high. We would rather predict high nitrogen dioxide levels when they are not, rather than predicting them to be low when they are high. Looking at our error figures (see figures 13 and 14 in the appendix), we can see that generally the amount of type 1 error and type 2 error is very similar. Taking this project forward it would be a useful development to alter the model to account for a smaller type 2 error.

## 5.3  Limitations

- The model is only as accurate as the meteorological and traffic data. Weather forecasts are often inaccurate predicting too far into the future (a few days usually). If they are incorrect, it is likely that our model will be as well. With lots of training data, however, this issue can be reduced.

- The binary nature of our model limits the range of $N0_2$ levels it can output. Having more discrete levels would be useful for people to know precisely how hazardous the nitrogen dioxide levels are.

- The model only includes three factors for each location. To get a better picture of the pollution levels across Bristol, more locations and factors could be considered.

# 6  Conclusion

## 6.1  Findings

The first major finding is that the three factors had to align for the pollutant to be high. Wind speed had to be low, temperature had to be low and traffic flow had to be high. Before the project began, we expected high levels of traffic would be a key factor. We also thought that high wind speeds would decrease the level of pollutants as it would simply disperse the particles. However, the finding that a high temperature decreases the level of $NO_2$ was surprising.

Another interesting discovery is how the complex model increased the accuracy of prediction. This meant that all three locations were affected by traffic from different areas, roughly a 1% increase in accuracy. Taking this project forward it would be interesting to implement a more dynamic network of traffic into the model to see if that increased the accuracy of predictions.

Finally, one of the most important findings is the model's sensitivity to thresholds. With an increased number of discrete levels, the sensitivity of the thresholds could decrease, resulting in a more stable model. This would also give more variation in the output, thereby producing a more realistic model.

## 6.2  Further Work

The most accessible improvement to the model would be the addition of more variables. This could be in the form of additional data sets such as wind direction, humidity and time of year. Due to the nature of Bayesian networks, adding variables is easily done, as the network automatically assesses the dependence of the output on the input data.

One could also introduce an increased number of discrete levels. For example, the $NO_2$ data could be split into background, above the EU recommended and dangerously high levels. This could then be used to quantify the risk the model predicts and result in a more useful

model. The contributing factors to each level of $NO_2$ could be identified so that the levels can ultimately be reduced and dangerous episodes avoided.

The assumption of conditional independence between the factors temperature, wind speed and traffic flow can be relaxed. By introducing dependencies between them one could generate more realistic results. For example, due to increased population and building density in the city centre, temperatures may be higher than the surrounding suburbs, resulting in an influx of wind towards the centre.

There is much that can be done to improve modelling pollution levels in Bristol. Bayesian networks have shown to be a useful and easy tool for modelling purposes of this kind, but have also proven to be computationally expensive and limited to the number of variables and links used in their design. While it is of paramount importance to be able to predict pollution, we feel it is just as important to make these predictions accessible to the people of Bristol. We created a website that allows the user to predict the $NO_2$ levels given temperature, wind speed and traffic flow levels (see image 18). We hope that it will inspire future development of a website that uses meteorological and traffic forecasts to predict pollution levels, so people can be informed and prepare themselves appropriately to reduce health risks.

# References

[1] Icopal Noxite. Nitrogen oxide (nox) pollution. `http://www.icopal-noxite.co.uk/nox-problem/nox-pollution.aspx`. Assessed: 08/05/2019.

[2] BBC News. Uk air pollution: How bad is it? `https://www.bbc.co.uk/news/uk-26851399`. Accessed: 08/05/2019.

[3] Public Health England. Health matters: Air pollution. `https://www.gov.uk/government/publications/health-matters-air-pollution/health-matters-air-pollution`. Accessed: 08/05/2019.

[4] ICOPAL. Noxite. `http://www.icopal-noxite.co.uk/nox-problem/nox-pollution.aspx`. Assessed: 2015.

[5] Open Data Bristol. Real time air quality data. `https://opendata.bristol.gov.uk/explore/dataset/nox_wide/table/?disjunctive.location&sort=date_time&refine.date_time=2014`. Assessed: 10/03/2019.

[6] Department for Transport. Road traffic estimates: Great britain 2017. `https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/741953/road-traffic-estimates-in-great-britain-2017.pdf`. Assessed: 07/05/2019.

[7] Department of Transport. City of bristol, traffic profile for 2000-2017. `https://www.dft.gov.uk/traffic-counts/area/regions/South+West/local-authorities/Bristol%2C+City+of`. Assessed: 08/05/2019.

[8] Food Department of Environment and Rural Affairs. Uk air: Air information resource. `https://uk-air.defra.gov.uk/data/`. Assessed: 10/03/2019.

[9] Open Data Bristol. Meteorological data - create centre. `https://opendata.bristol.gov.uk/explore/dataset/meteorological-data-create/information/`. Assessed: 15/03/2019.

[10] K. Vala. Probabilistic graphical models: Bayesian networks. `https://towardsdatascience.com/probabilistic-graphical-models-bayesian-networks`. Accessed: 08/05/2019.

[11] D. Delaunay. Wind speed distribution. `https://www.sciencedirect.com/topics/engineering/wind-speed-distribution`. Assessed: 08/05/2019.

[12] Air Quality Expert Group. Nitrogen dioxide in the united kingdom. `https://uk-air.defra.gov.uk/assets/documents/reports/aqeg/nd-summary.pdf`. Accessed: 08/05/2019.

[13] Bristol's 2018 air quality report. `https://www.bristol.gov.uk/documents/20182/32675/Bristol+City+Council+2018+Air+Quality+Annual+Status+Report+ASR/3d5c287b-f379-e484-7924-2aa02fc8bb0a`. Accessed: 05/04/2019.

[14] Met Office. Beaufort wind force scale. `https://www.metoffice.gov.uk/weather/guides/coast-and-sea/beaufort-scale`. Accessed: 05/04/2019.

[15] R. Santos-Rodriguez. Inference. dynamic bayesian networks. `https://www.ole.bris.ac.uk/bbcswebdav/pid-3400135-dt-content-rid-11030055_2/courses/EMAT31530_2018/HMM_with_notes%281%29.pdf`. Accessed: 07/05/2019.
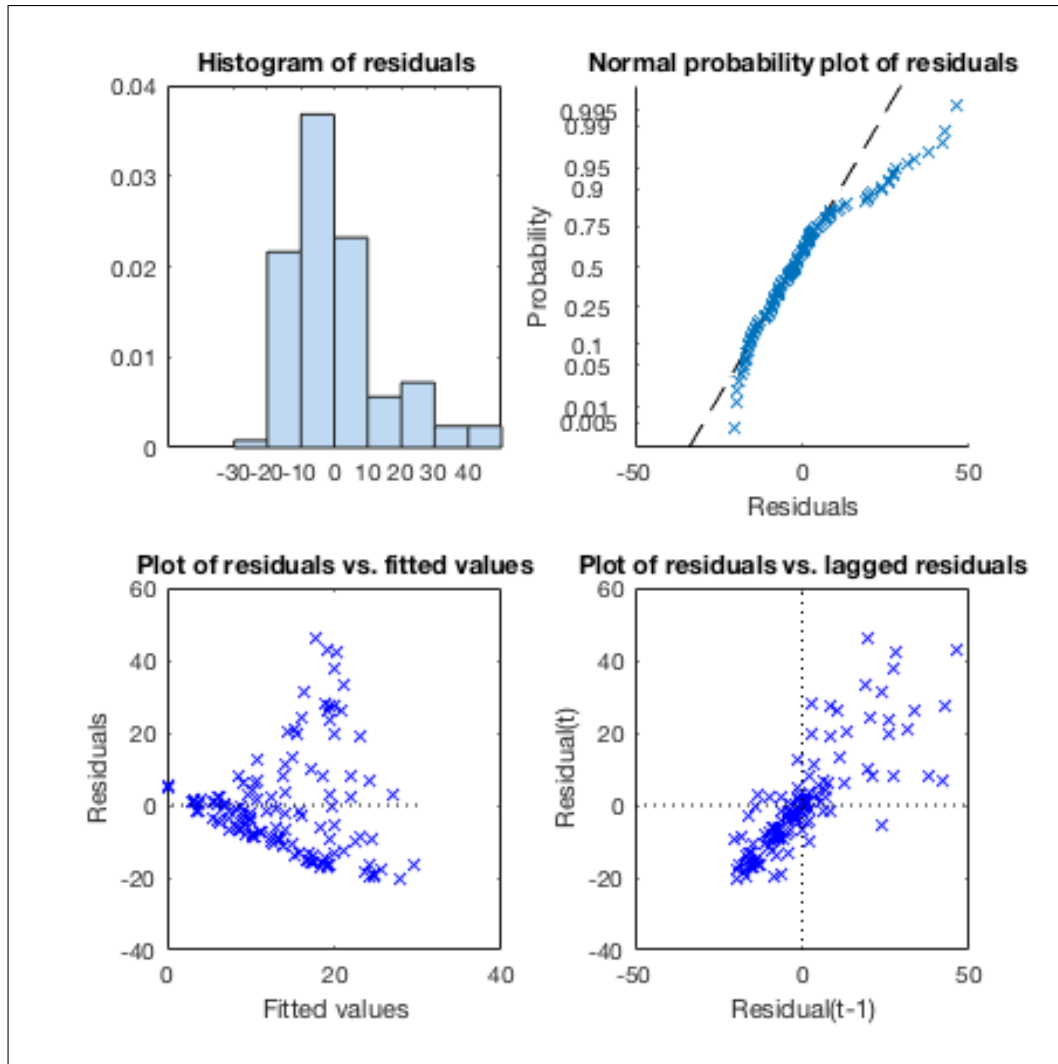
# Appendix

## Residual Plot



Figure 6: These are the residual plots for St Pauls, which are also representative of the other two locations. The figure shows four different plots to test whether there exists a simple linear relationship between the factors and the nitrogen dioxide levels. The histogram and the normal probability plot of residuals suggest the errors do not follow a normal distribution as there is no Gaussian form to the histogram and the trend of the normal probability plot is not linear. The plot of residuals against the fitted values suggests that the assumption of homoscedasticity is violated as the error variance is not constant, and the plot of the residuals against lagged residuals shows that the assumption that the errors are not correlated does not hold. All these graphs provide evidence to suggest that a simple linear model would not be appropriate.

# Thresholds



Figure 7: Percentage Errors with varying wind speed thresholds for the St Pauls location. All errors are recorded over a period of ten days.



Figure 8: Percentage error, with $NO_2$ thresholds of $40\mu gm^{-3}$ and $60\mu gm^{-3}$ for the St Pauls location for a period of 10 days
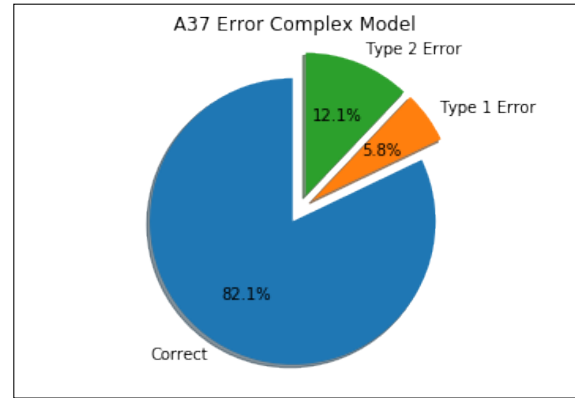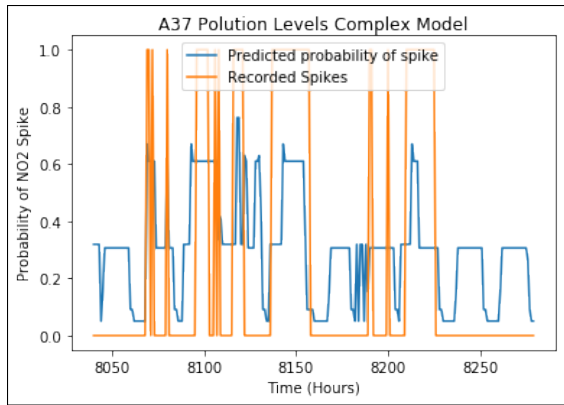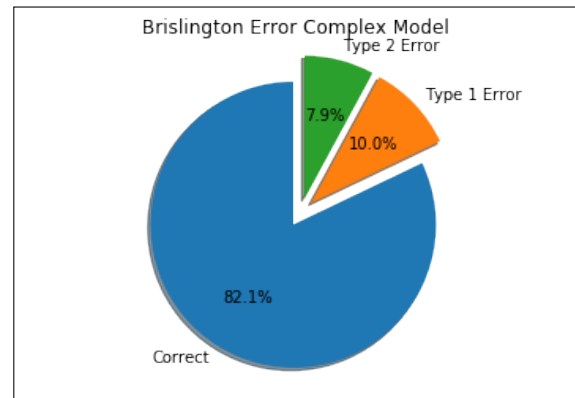
# Results



Figure 9: Recorded $NO_2$ spikes against the probability output for the simple model at the A37 location. Also plotted are the error types for the A37 location. Both are recorded over a period of 10 days.
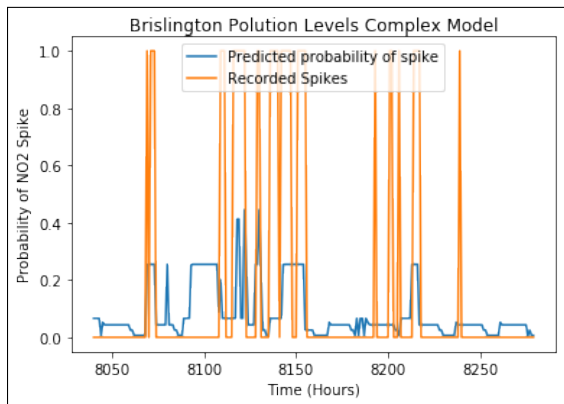


Figure 10: Recorded $NO_2$ spikes against the probability output for the simple model for the Brilington location. Also plotted are the error types for the Brislington location. Both are recorded over a period of 10 days.
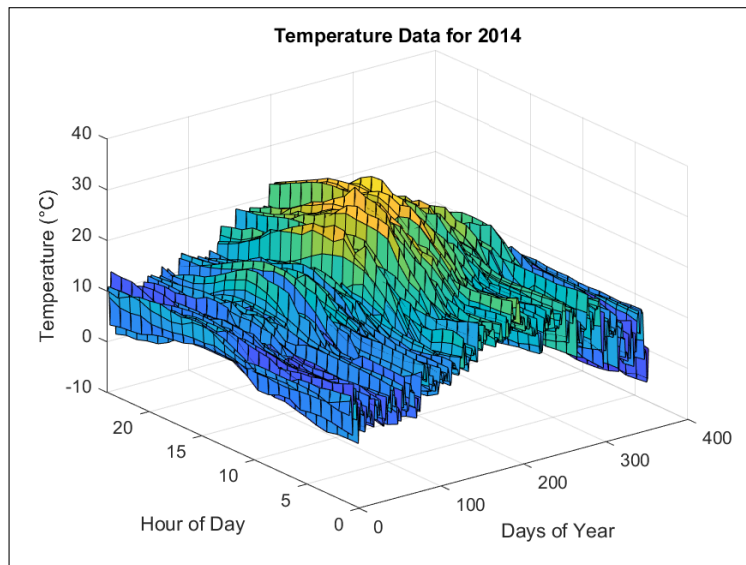
15

Figure 11: Recorded $NO_2$ spikes against the probability output for the complex model for the A37 location. Also plotted are the error types for the A37 location. Both are recorded over a period of 10 days.



Figure 12: Recorded $NO_2$ spikes against the probability output for the complex model for the Brilington location. Also plotted are the error types for the Brislington location. Both are recorded over a period of 10 days.

16

# The Data



Figure 13: This is a graph plotting the raw data of the temperature for the year of 2014.



Figure 14: This graph plots the raw data of the wind speeds recorded in Bristol in 2014.
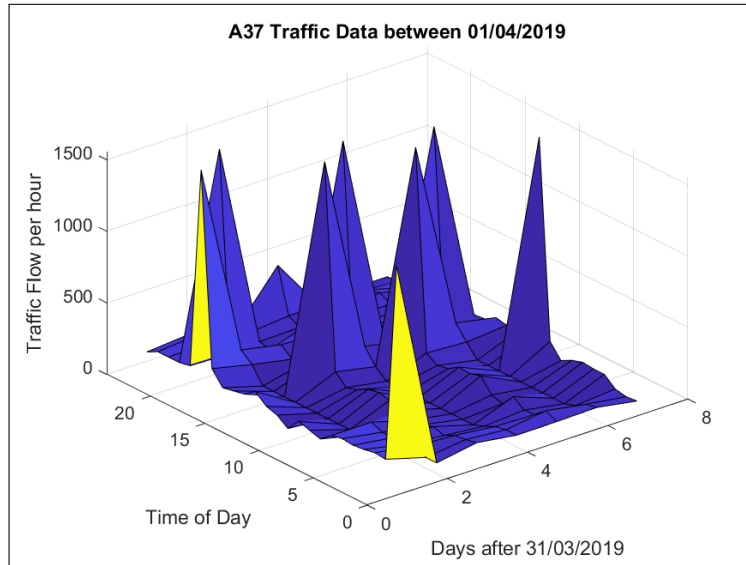
Figure 15: This is a graph plotting the raw data of the traffic flow passing through Well's road, coming off the A37 between the dates of 01/04/2019 and 07/04/2019.
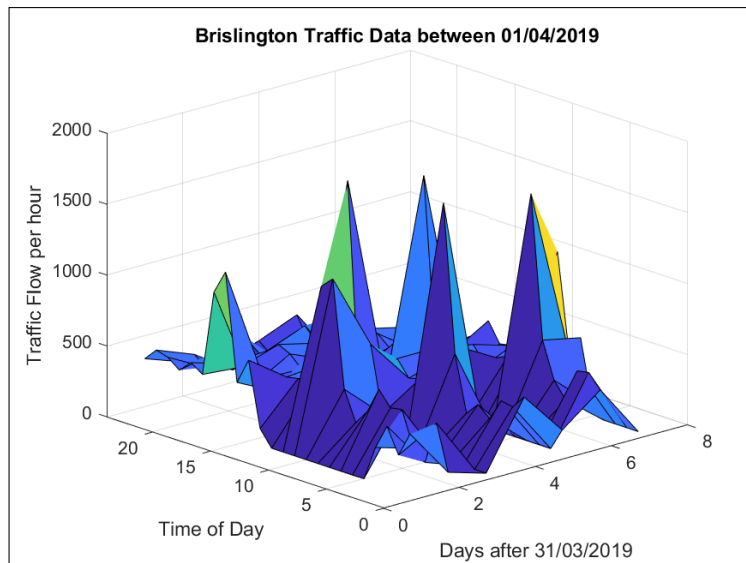


Figure 16: This is a graph plotting the raw data of the traffic flow passing through St Philip's Causeway, Brislington, between the dates of 01/04/2019 and 07/04/2019.
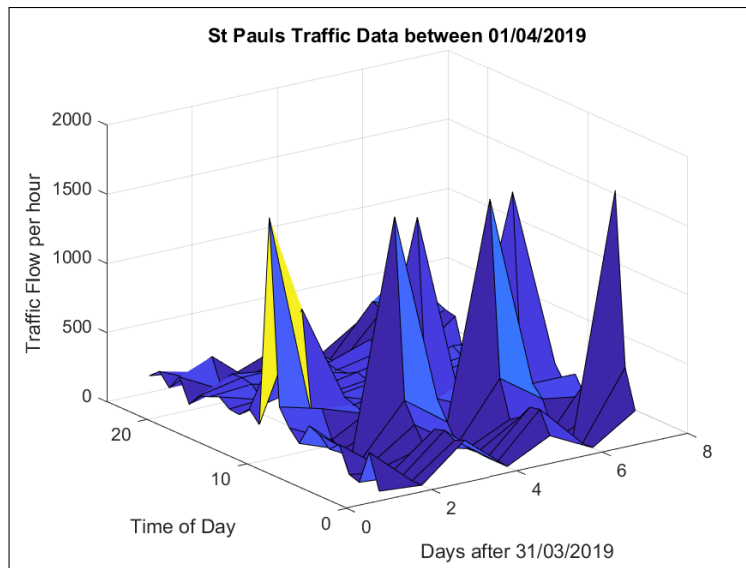
18

Figure 17: This is a graph plotting the raw data of the traffic flow passing through Ashley Road, St Paul's, between the dates of 01/04/2019 and 07/04/2019.
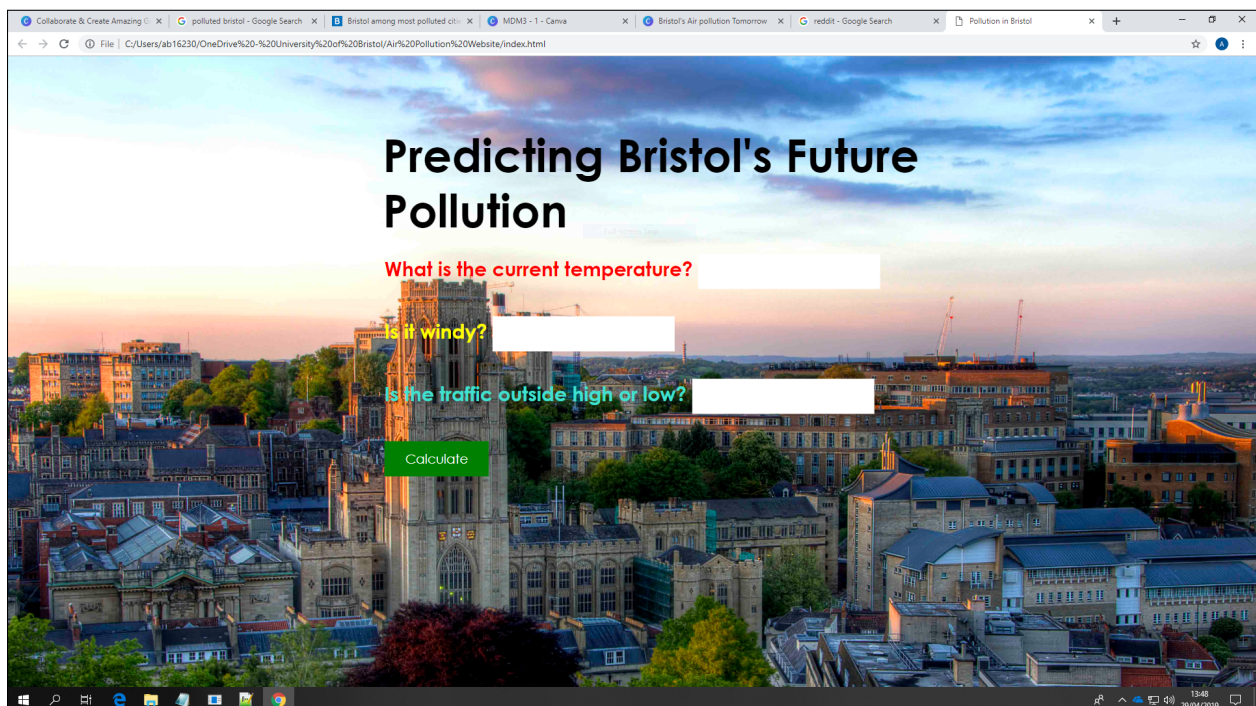


Figure 18: Our website. We hope this will inspire others to create something similar. Ideally, the website could use weather and traffic forecast information provided by other websites to then predict pollution levels for the next coming days for various locations across Bristol.

19