

Introduction

In this group project we will be conducting an investigation into the use of Autoregressive Integrated Moving Average (ARIMA), Vector Autoregressive (VAR) and Long Short-Term Memory (LSTM) models for forecasting air pollution in Beijing and seeing how our models would work on a large scale.

We will be exploring this [dataset](#). We found it on Kaggle and it is a dataset that reports the weather and air pollution level at the US embassy in Beijing, China. The data from this dataset takes readings every hour for five years. It was specifically put on Kaggle so that users could explore LSTM architecture for time series forecasting. The reason that this dataset was favourable was its data completeness but also the fact that many people had already used it suggesting that the dataset might offer some interesting insights.

1 What we're doing

We'll be interested in analysing the predictive and computational performance of our chosen models and consider how these models will generalise to a larger dataset.

Initially, we'll look at the ARIMA and VAR models and take these as our baseline. They make intuitive sense in the context of the problem and are mathematically accessible models. The VAR model looks at all the variables from the previous few days of data and uses a linear combination of these variables to try and predict what each of these will be the next day. In a sense, ARIMA is a univariate version of this model. They will give us good baselines in order to get a feel for how a simple model performs and the computational complexity involved.

We'll then look at two deep learning models, a univariate Long short-term memory (LSTM) model and a multivariate LSTM model. These are a type of Recurrent Neural Network (RNN) and allow for prediction using sequential data, accounting for long-term dependencies. Whilst univariate relies solely on pollution levels, multivariate will incorporate additional weather variables, introducing complexity and allowing us to interpret if additional features are necessary for LSTM networks. These are more complex models than what we looked at initially and allows us to assess how the computational and predictive

performance changes as the models complexity increases.

2 Investigating Computational Performance

Our main areas of interest with regards to assessing the computational performance will be the time it takes the models to train, the amount of memory the model uses, and the complexity of the model. These will be key when we come to think about how the work we've carried out can be applied to a much larger amount of data, such as datasets that look at pollution and its related variables over the course of decades.