

NUS Business School Honors Dissertation

Peng Seng Ang

AY 2019/2020 Semester 2

Abstract

This paper studies how we can use various spatial-temporal time series model to better predict demand across different locations.

1 Introduction

Having an accurate forecast of delivery demand for food service providers would help them more effectively and efficiently assign orders to drivers to improve the overall delivery time. Currently, most Autoregressive (AR) or Autoregressive Integrated Moving Average (ARIMA) models only consider temporal features when predicting demand. However, we believe including spatial features between the data points might improve forecast accuracy. This paper would focus on and explore models that include both spatial and temporal features to improve forecast accuracy.

2 Literature Review

test

3 Data

The data source used was an operational dataset from a food delivery service provider from Shanghai that includes delivery information for a 2-month period from 10 August 2015 to 30 September 2015 (excluding Saturdays) in 2015. The provider only provides delivery service for 90 minutes during lunchtime and the

dataset has split the data into 15-minute time periods, and as such, each day would only consists of demand data for 6 time periods. Hence, our dataset has 839 locations with demand data for 204 time periods in total.

To include other exogenous variables, data from <https://www.worldweatheronline.com/shanghai-weather-history/shanghai/cn.aspx> was used to include weather and rainfall data as well as encoding of the weekadys for all the respective days.

3.1 Exploratory Analysis

We would first do some exploratory analysis and check if there are any obvious relationships between the variables.

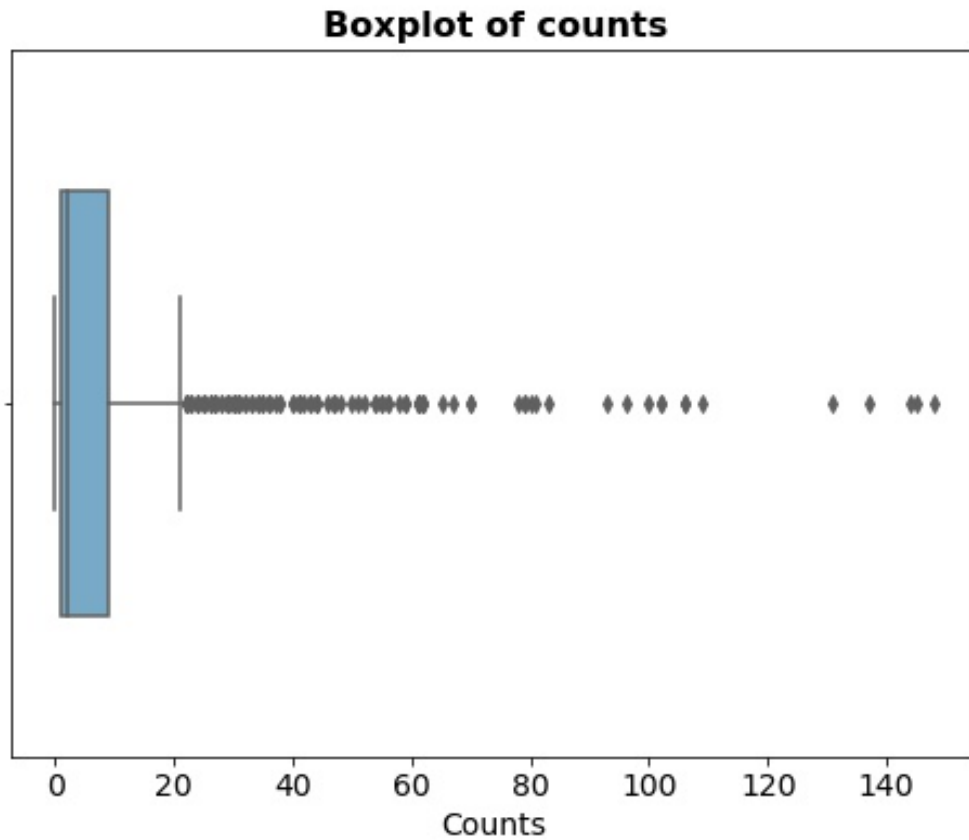


Figure 1: Boxplot of counts

We can see from the boxplot that most of the locations have extremely low number of non-zero orders, with about 335 locations having just a maximum of one non-zero order throughout the 204 time periods.

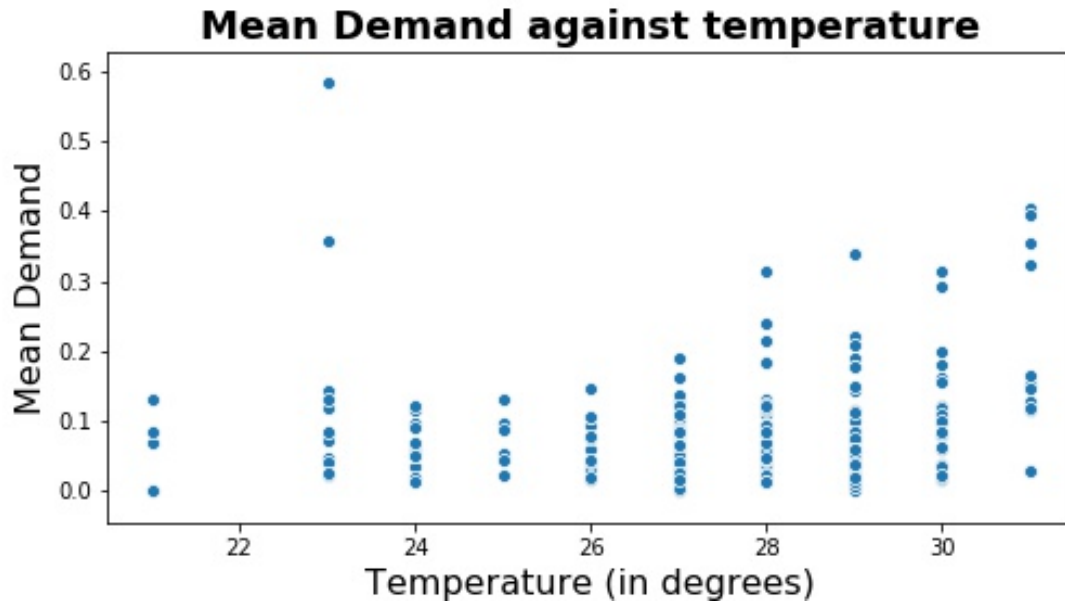


Figure 2: Scatter plot of mean counts against temperature

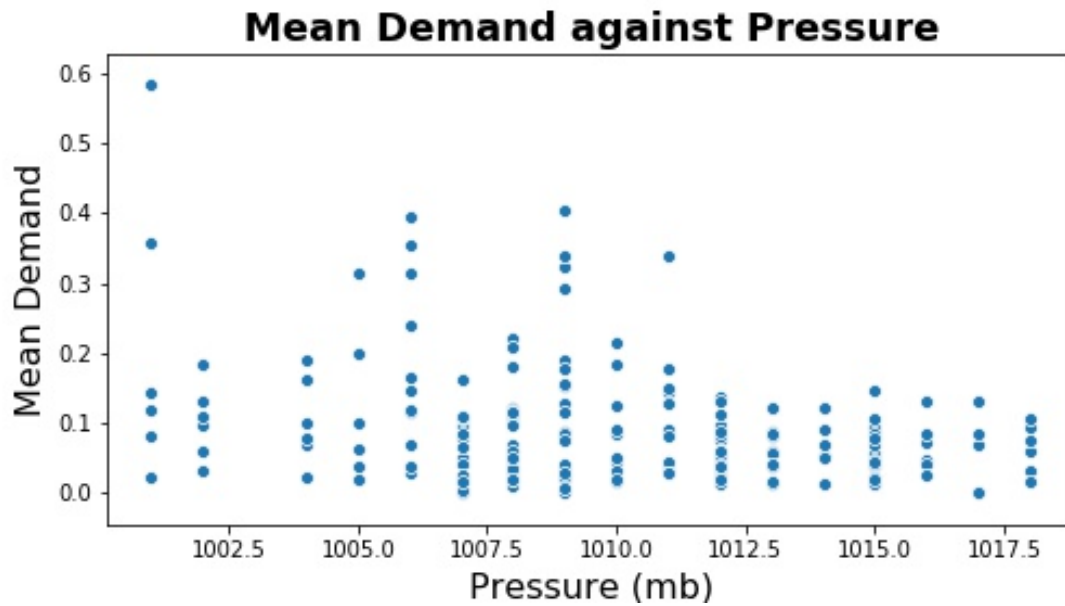


Figure 3: Scatter plot of mean counts against pressure

The scatter plot in Figure 2 visually display a slight positive relationship between temperature and mean demand across all locations whereas Figure 3 visually display a slight negative relationship between pressure and mean demand across all locations.

The distribution plots for the rest of the exogenous variables can be found in the appendix.

4 Baseline Model

In this section, we would build a simple baseline model. Following which, we would try other different spatial temporal time series models and compare the results to the baseline model.

4.1 Metric Used

The main metric that would be used for comparison would be Mean Squared Forecast Error (MSFE), which is calculated by:

$$MSFE = \frac{1}{n} \sum_{t=1}^n \|\hat{y}_t - y_t\|_2^2$$

where n is the number of data points, \hat{y}_t is the predicted demand at time t and y_t is the actual demand at time t .

4.2 Train-Test Split

From Figure 1 in Section 3.1, the data is very sparse as there are many locations that have no demand counts for the majority of the time period. Hence, to get a better idea of how our models would work, only locations with at least 50 non-zero counts across the time period would be used initially, leaving us with 42 locations that meet this criteria. The dataset was then split into training and test set by considering the first 27 days as the training set and the next 7 days as the test set. Our training set would then have 162 demand data for each location and test set would have 42 demand data for each location.

4.3 ARIMA models

Autoregressive Integrated Moving Average (ARIMA) models are one of the most commonly used models for time series (Z. Asha Farhath (2016)). ARIMA models are made up of 3 processes, mainly the Autoregressive

(AR) process, the Integrated (I) process and the Moving Average (MA) process (Jamal Fattah (2018)). The AR process assumes that each observation can be expressed as a linear combination of its past values. An AR(x) process would mean using x lagged values. The MA process assumes that each observation can be expressed as a linear combination of its current error term as well as its past error terms. The Integrated Process states that the time series can undergo differencing to ensure that the series is stationary. A MA(x) process would mean using x number of past observations. Hence, an ARIMA model is usually represented by ARIMA(p,d,q), where p represents the number of autoregressive terms, d represents the number of differences needed for stationarity, and q represents the number of lagged forecast errors.

4.4 Baseline ARIMA Result

As a baseline model, each of the locations was assessed individually and a suitable ARIMA model was built for each location. Auto-arima function from Python was used to implement this. The out-of-sample MSFE for this baseline model on the 42 locations is **58.80**. A sample forecast plot is shown below:

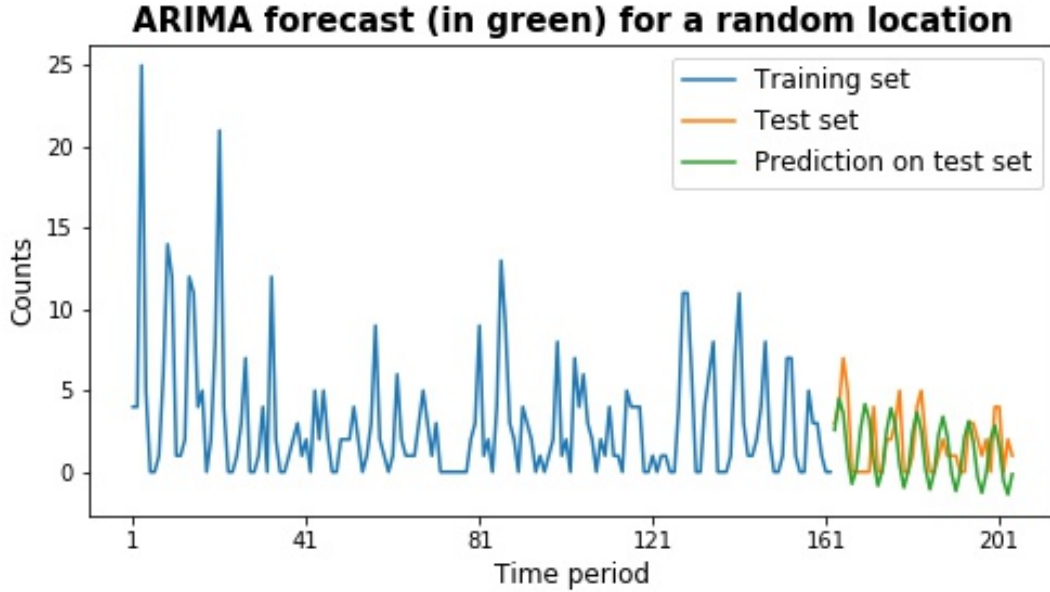


Figure 4: ARIMA forecast on a random location

5 VAR Model

Vector Autoregressive (VAR) models are the most commonly used model for multivariate time series, particularly in economics and financial time series as shown in Bjørnland (2000). VAR models are very similar to multivariate linear regression models and methods used to perform inferencing on linear regression models can also be applied to VAR models. VAR(p) represents a VAR model of order p if the time series can be written as:

$$x_t = \sum_{i=1}^p \phi_i x_{t-i} + \alpha_t$$

where x_t is the value of time t , ϕ_0 is a constant vector and ϕ_i are coefficient matrices for $i > 0$ and α_t are independent and identically distributed random vectors. To validate if the multi-variate time series is stationary, the Johansen's test for cointegrating time series would be performed.

Assumption 1. *The first assumption of a VAR model is....*

Assumption 2. *It is generally true...*

5.1 VARX Model

VAR models can also be extended to include exogenous variables.

5.2 Model Checking

To validate and verify if our fitted model is adequate, model checking would be performed by performing the following residual analysis:

test

5.3 Results

BigVAR library in R was used to implement the VAR models. The results from the VAR model without exogenous variables gives an out-of-sample MSFE of 46.641 on the 42 locations.

6 GLM Model

Assumption 3. *Generalised Linear Models (GLM) are.....*

Assumption 4. *An assumption is that the data follows a poisson process or a non-homogenous poisson process.*

6.1 Model Checking

To validate our model,

6.2 Insights and Implementation

Any findings from the results? If there are any benefits or issues in implementing the proposed model...

7 Conclusion

Conclude your efforts and main findings.

8 Appendix

Append extra plots, graphs, analysis, etc.

References

- Bjørnland, H. C. (2000). Var models in macroeconomic research.
- Jamal Fattah, Latifa Ezzine, Z. A. H. E. M. A. L. (2018). Forecasting of demand using arima model.
- Z. Asha Farhath, B. Arputhamary, L. A. (2016). A survey on arima forecasting using time series model.