

Data Science Internship at Data Glacier

Project: Retail Forecasting

Week 7: Deliverables

Group name:

Group members: Gordon Poon, Keith Tang, Joseph Xu

Email: gordontxpoon@gmail.com, zx1054@nyu.edu,
keithtang0901@gmail.com

Country: United Kingdom, China

College: UCL, NYU, Durham University

Specialisation: Data Science

Table of Contents:

1. Project description.....	3
2. Business understanding.....	3

1. Project description

Dataset was provided by a large beverage company in Australia. They sell their products through various super-markets and also engage into heavy promotions throughout the year. Their demand is also influenced by various factors like holiday, seasonality. They needed a forecast of each of the products at item level every week in weekly buckets.

2. Business Understanding

Determine Business Objectives:

1. Forecast the item level of 6 products at each week

Assess Situation (Assumptions):

1. Relationship exists between sales and holidays
2. Relationship exists between sales and promotion
3. Relationship exists between sales and Covid
4. Relationship exists between sales and Google Mobility
5. Relationship exists between sales and discount
6. Relationship exists between sales of one product and another

Determine Data Science Goals:

1. Build 4-5 multivariate forecasting model
2. Demonstrate best in class forecast accuracy
3. Write a code in such a way in order to run the model in least time

4. Demonstrate explainability in the form of contribution of each variables

Project Plan:

1. Week 7: Business understanding
2. Week 8: Data Understanding
3. Week 9: Data Cleaning and preparation
4. Week 10: EDA
5. Week 11: EDA presentation and proposed modelling technique
6. Week 12: Model selection and model building

3. Data Understanding

Describe data:

Total number of observations	1218
Total number of features	12
Base format of the file	.xlsx
Size of the data	74 kB

Explore data and check data quality:

1. Data types:

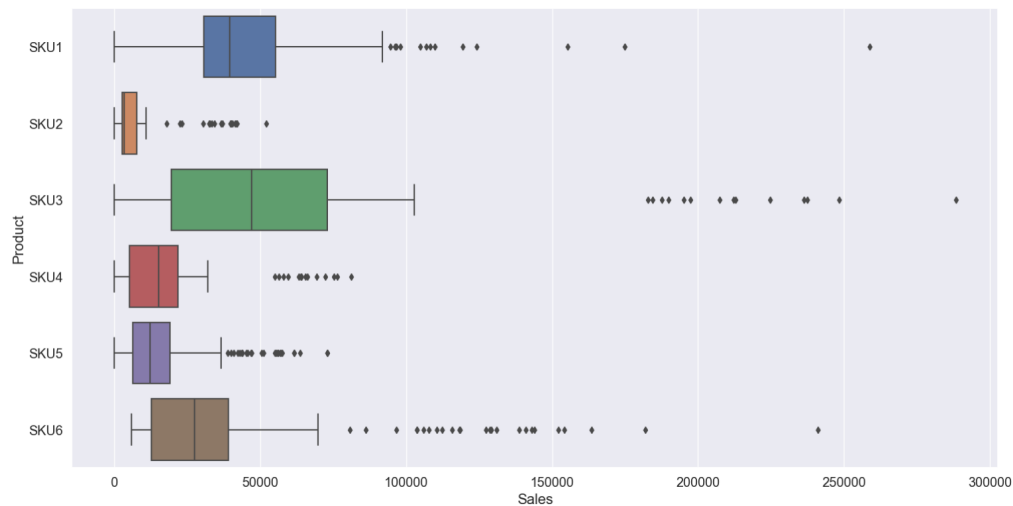
```
Product          object
date             object
Sales            int64
Price Discount (%) object
In-Store Promo   int64
Catalogue Promo  int64
Store End Promo  int64
Google_Mobility  float64
Covid_Flag       int64
V_DAY            int64
EASTER           int64
CHRISTMAS        int64
dtype: object
```

The data types of 'date' and 'Price Discount (%)' should be modified.

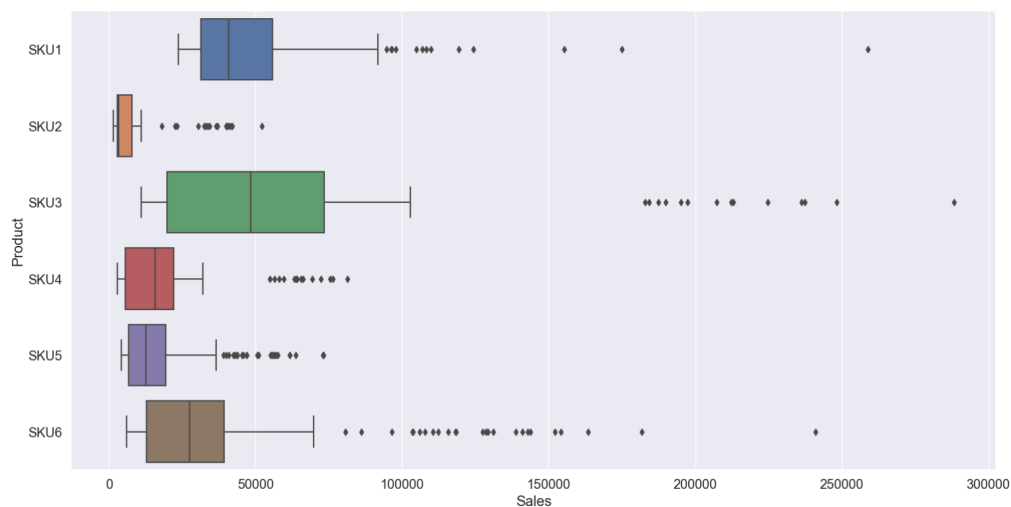
2. Missing values:

There is no missing values in the dataset

3. Outliers

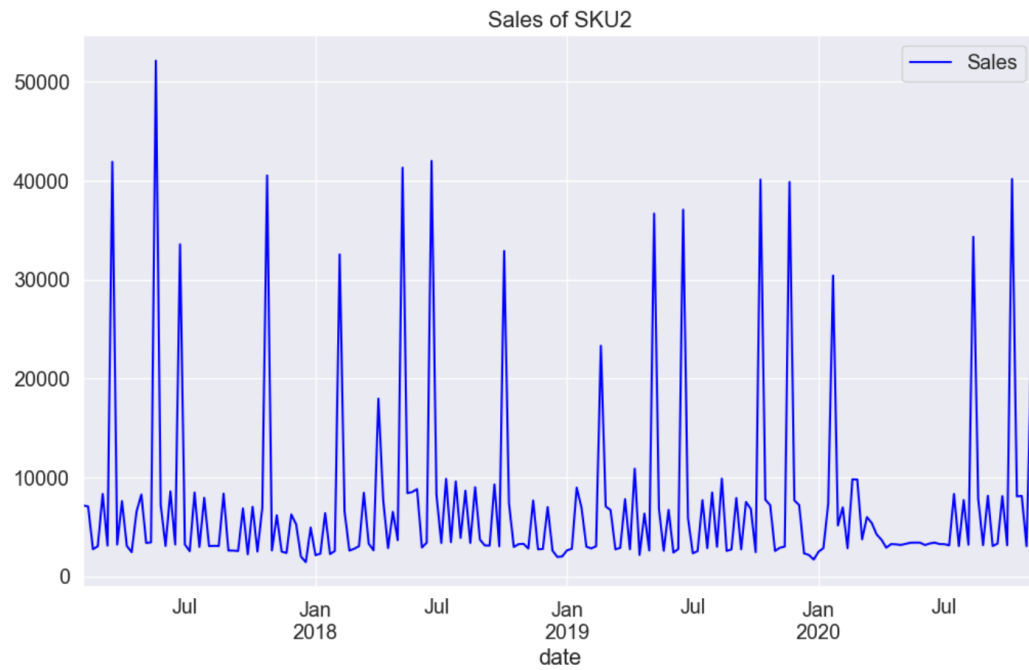
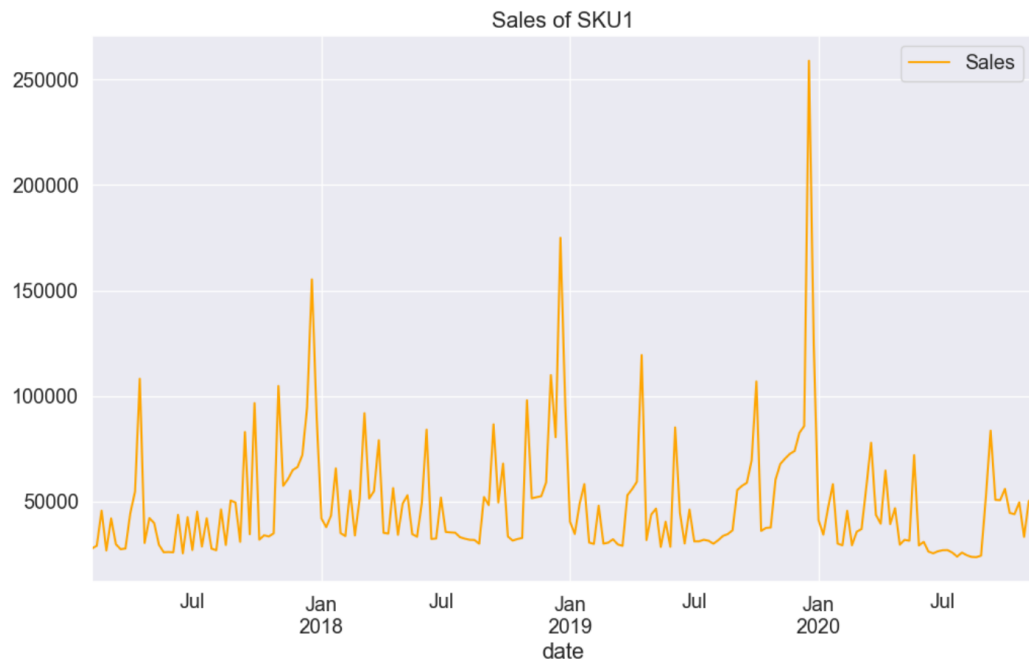


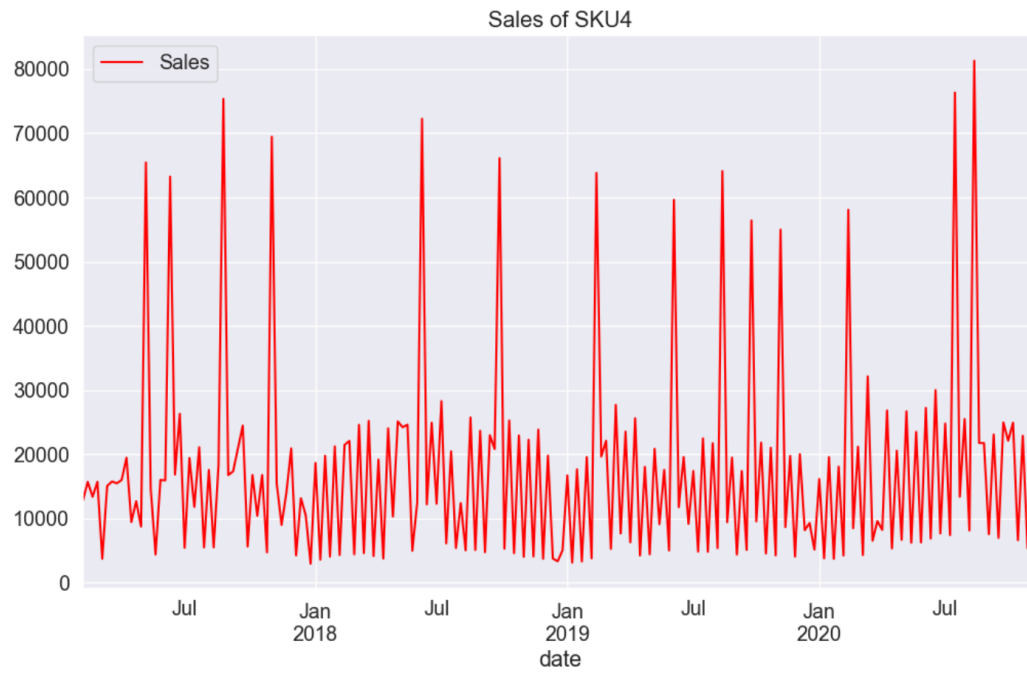
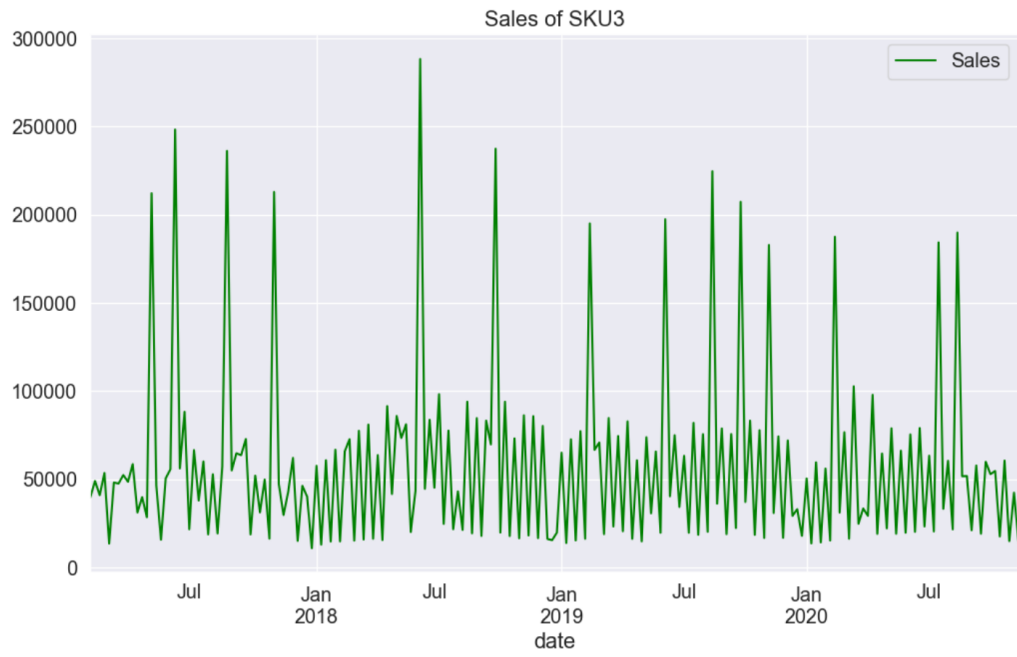
There are weeks with zero Sales in SKU1-5. It is believed that the testing weeks were mixed into the dataset.

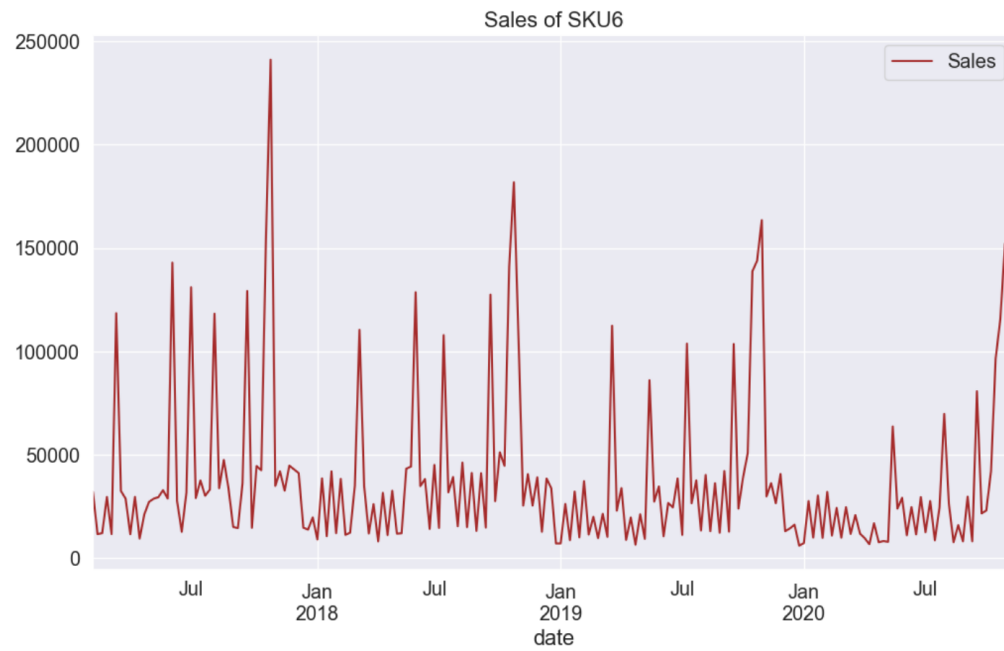
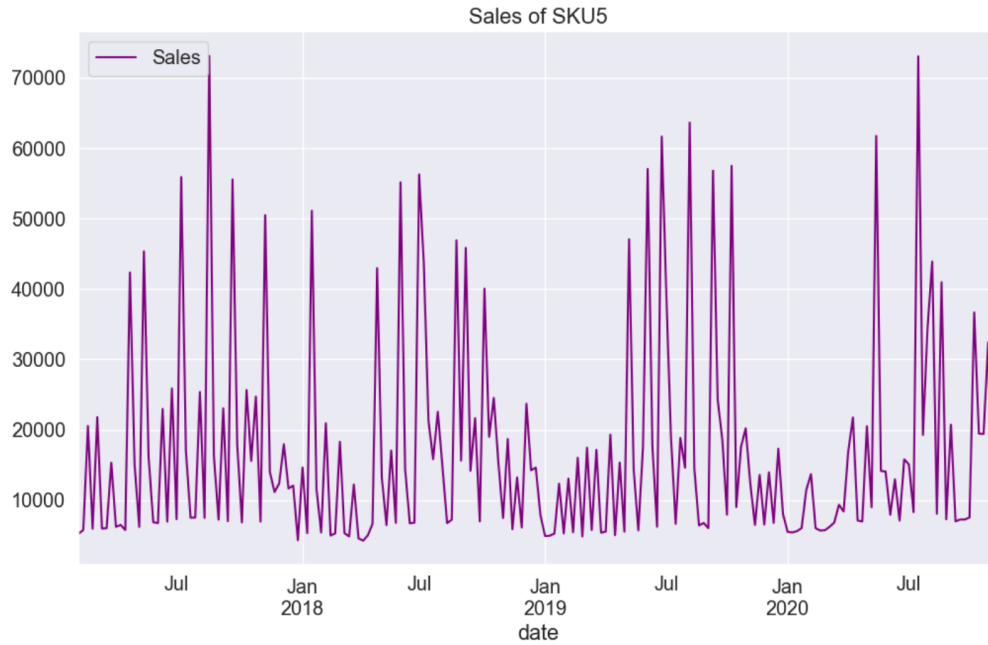


After modifying the data, there are still some outliers in the Sales feature. But since we do not have enough information on the components for the sales, it is not appropriate to treat it as an outlier.

4. Data Visualisation







5. Correlation

	Sales	Price Discount (%)	In-Store Promo	Catalogue Promo	Store End Promo	Google_Mobility	Covid_Flag	V_DAY	EASTER	CHRISTMAS
Sales	1.000000	0.432887	0.252071	-0.124398	0.234811	0.044885	-0.047748	-0.011653	-0.013495	-0.013893
Price Discount (%)	0.432887	1.000000	0.225429	-0.091492	0.234464	-0.207491	0.265120	-0.042953	0.003640	-0.035163
In-Store Promo	0.252071	0.225429	1.000000	-0.488728	0.367410	0.060471	-0.038873	0.020951	0.020951	0.021550
Catalogue Promo	-0.124398	-0.091492	-0.488728	1.000000	0.124778	0.075202	-0.098381	-0.045466	-0.045466	0.036468
Store End Promo	0.234811	0.234464	0.367410	0.124778	1.000000	0.082501	-0.067667	0.019489	-0.068211	0.009620
Google_Mobility	0.044885	-0.207491	0.060471	0.075202	0.082501	1.000000	-0.764376	0.076392	-0.111869	0.048450
Covid_Flag	-0.047748	0.265120	-0.038873	-0.098381	-0.067667	-0.764376	1.000000	0.015213	0.015213	-0.063385
V_DAY	-0.011653	-0.042953	0.020951	-0.045466	0.019489	0.076392	0.015213	1.000000	-0.020619	-0.017810
EASTER	-0.013495	0.003640	0.020951	-0.045466	-0.068211	-0.111869	0.015213	-0.020619	1.000000	-0.017810
CHRISTMAS	-0.013893	-0.035163	0.021550	0.036468	0.009620	0.048450	-0.063385	-0.017810	-0.017810	1.000000