# Capstone Project
## Retail Sales Prediction

**Midhun R**

**AI**

# Points for Discussion

- Business Task

- Data Summary

- Data Cleaning

- Exploratory Data Analysis

- Feature Engineering
- Modelling
- Conclusion

# Business Task

Two datasets are given: one with store data of 1115 stores and the other with historical sales data from January 2013 to July 2015.

The main objective is to analyze the given dataset and build a machine learning model to forecast the sales of all Rossmann stores upto 6 weeks.

This is undertaken as an individual project.

# Data Summary

|  | Store Data | Sales Data |
|---|---|---|
| Number of records (rows) | 1115 | 1017209 |
| Number of features (columns) | 10 | 9 |
| Number of duplicate rows | 0 | 0 |
| Number of columns with missing values | 6 | 0 |
| Number of columns require conversion of data type | 4 | 2 |

- In addition to this, extra columns are needed to be added for easier analysis.
- These irregularities will be handled during data cleaning step and the two datasets will be merged.

# **Data Summary** (Contd.)

## **Store Data:**

1. Store: a unique Id for each store.

2. StoreType: differentiates between 4 different store models: a, b, c, d.

3. Assortment: describes an assortment level: a = basic, b = extra, c = extended.

4. CompetitionDistance: distance in meters to the nearest competitor store.

5. CompetitionOpenSinceMonth: gives the approximate month of the time the nearest competitor was opened.

6. CompetitionOpenSinceYear: gives the approximate year of the time the nearest competitor was opened.

7. Promo2: Promo2 is a continuing and consecutive promotion for some stores: 0 = store is not participating, 1 = store is participating.

8. Promo2SinceWeek: describes the calendar week when the store started participating in Promo2.

9. Promo2SinceYear: describes the year when the store started participating in Promo2.

10. Promo2Interval: describes the consecutive intervals Promo2 is started, naming the months the promotion is started anew. Eg. "Feb,May,Aug,Nov" means each round starts in February, May, August, November of any given year for that store.
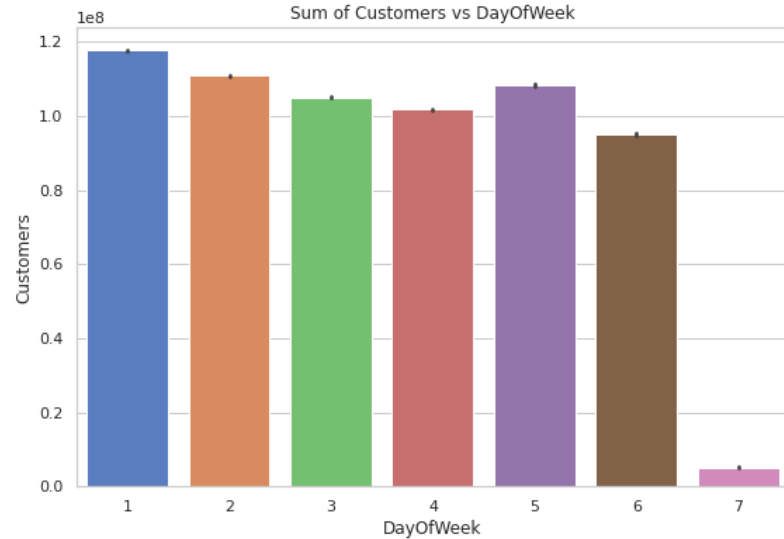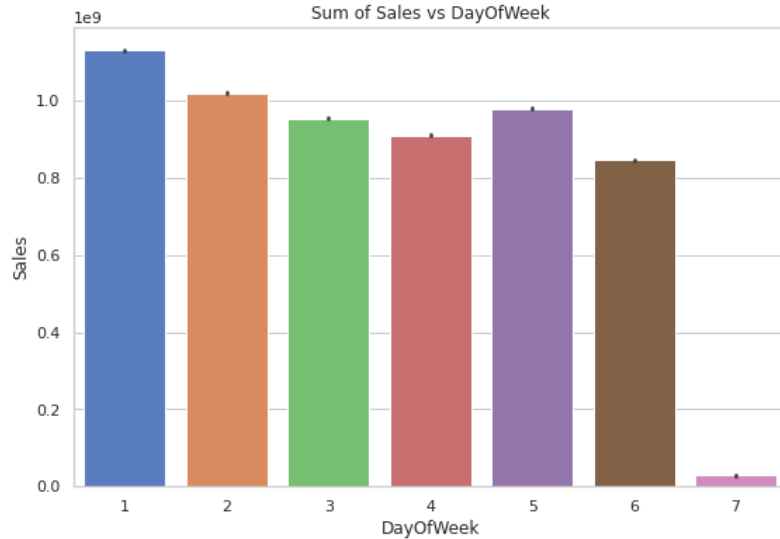
# **Data Summary** (Contd.)

## **Sales Data:**

1. Store: a unique Id for each store.

2. DayOfWeek: day of week of sale.

3. Date: date of sale

4. Sales: the turnover for any given day.

5. Customers: the number of customers on a given day.

6. Open: an indicator for whether the store was open: 0 = closed, 1 = open.

7. Promo: indicates whether a store is running a promo on that day.

8. StateHoliday: indicates a state holiday. Normally all stores, with few exceptions, are closed on state holidays. Note that all schools are closed on public holidays and weekends. a = public holiday, b = Easter holiday, c = Christmas, 0 = None.

9. SchoolHoliday: indicates if the (Store, Date) was affected by the closure of public schools.

# Data Cleaning

- 3 rows of data in 'CompetitionDistance' were null, it was imputed with the median of data.

- About 32% of data in 'CompetitionOpenSinceMonth' & 'CompetitionOpenSinceWeek' was null, they were imputed with the mode of data.

- About 49% of data in 'Promo2SinceWeek', 'Promo2SinceYear' & 'PromoInterval' was null in the same rows, where the value of Promo2 is zero. So, the missing values were imputed with 0.

- After the merging of two datasets, 'Date' was converted to datetime datatype and 'CompetitionOpenSinceMonth', 'CompetitionOpenSinceYear', 'Promo2SinceWeek' & 'Promo2SinceYear' to int datatype.

- 'StateHoliday' was converted into a binary feature.

- Five new features 'WeekOfYear', 'Month', 'Year', 'CompetitionOpenNumMonths' and 'Promo2NumWeeks' were added.
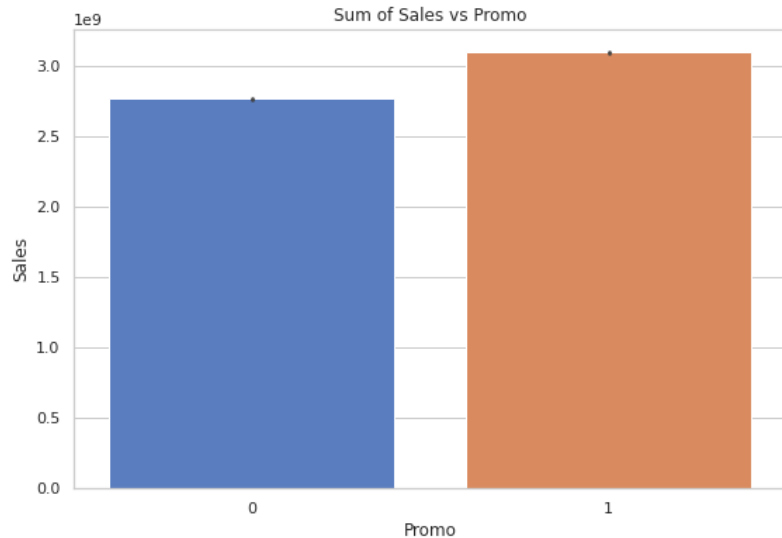
# Exploratory Data Analysis



Highest sales are recorded on Mondays and lowest sales are recorded on Sundays. This may be because most of the shops are closed on Sundays and this leads to higher demand on the next day, which is Monday.
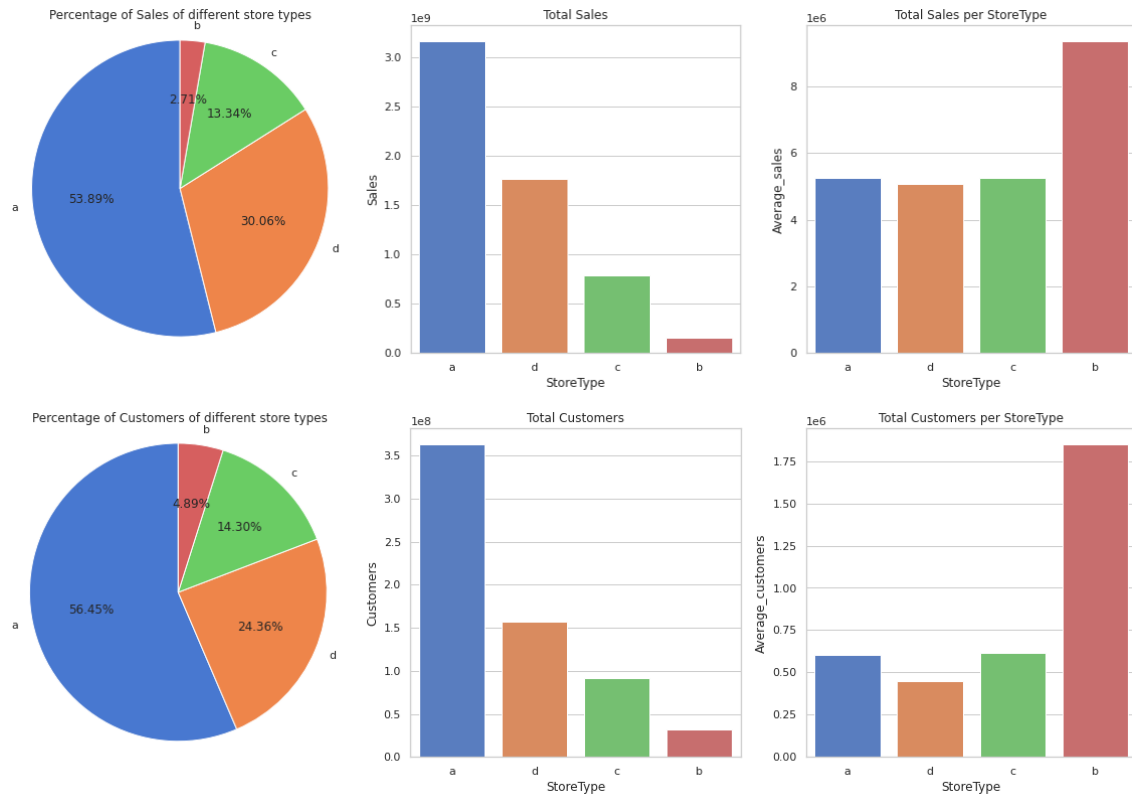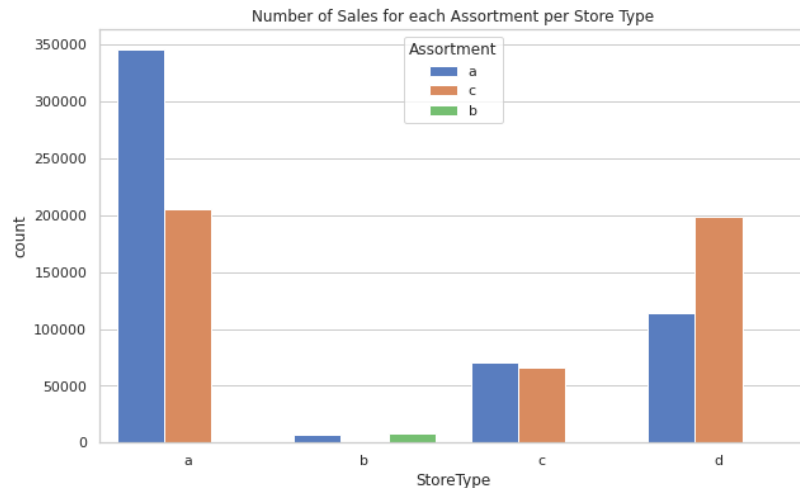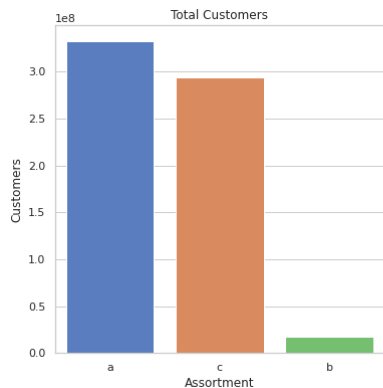
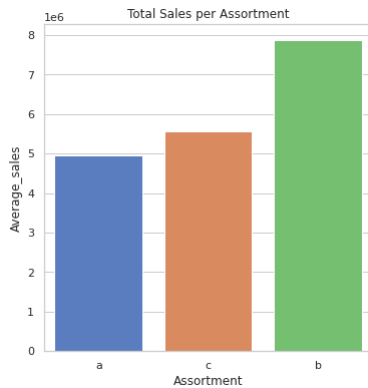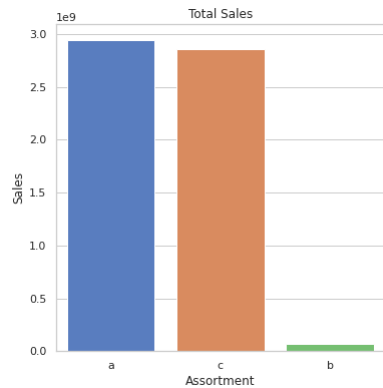# Exploratory Data Analysis (Contd.)



Presence of promos increases sales, but it doesn't help much in generating new customers.
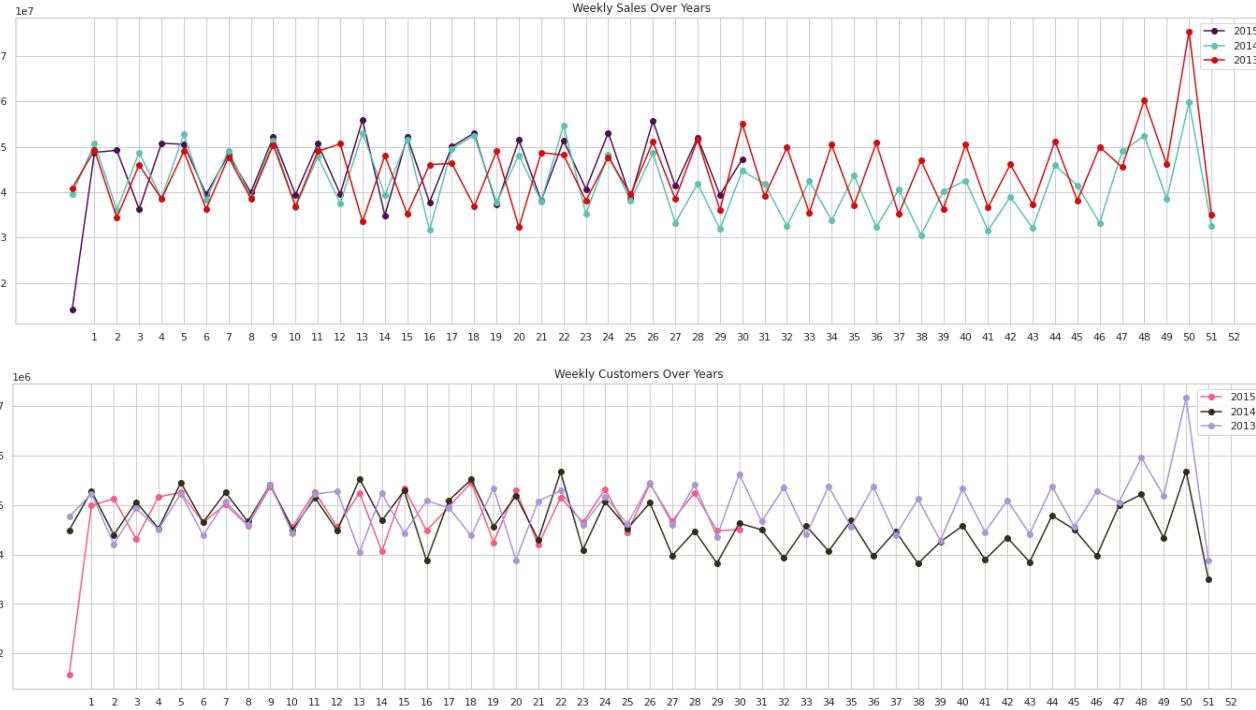
# Exploratory Data Analysis (Contd.)



Even though the volume of sales and customers is low, store type b has the highest average sales and customers. This means store type b is more preferred by customers.
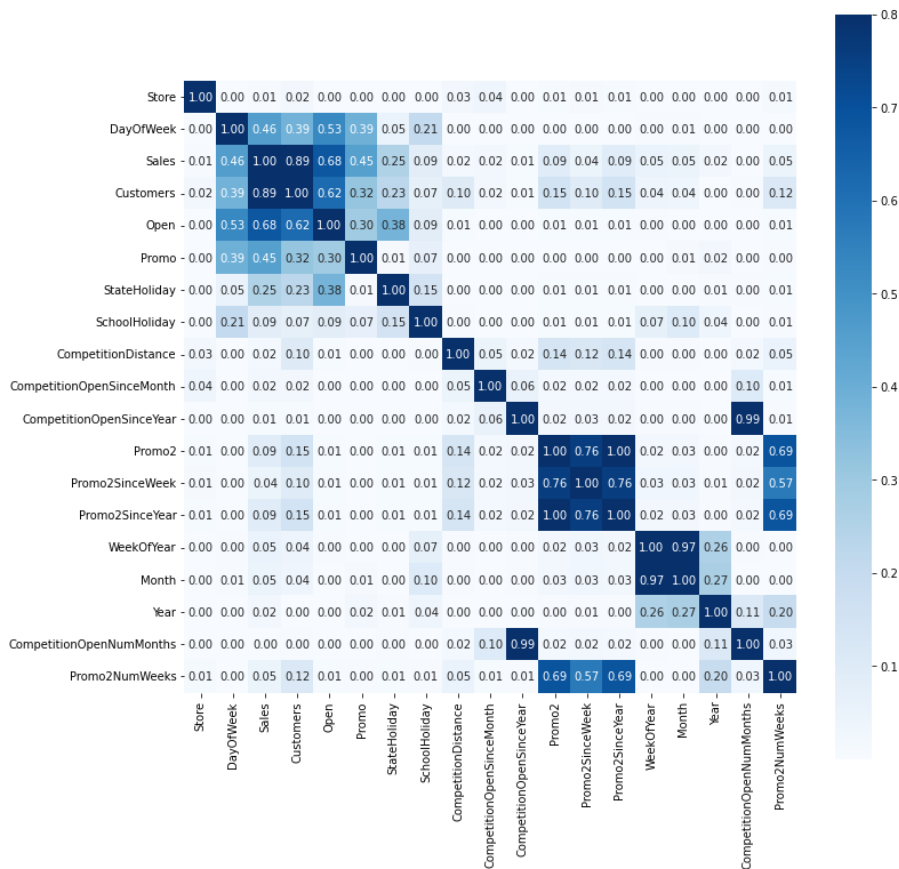
# Exploratory Data Analysis (Contd.)



Assortment b has the highest average sales and customers. This means assortment b is in high demand. Only store type b sells assortment b.

# Exploratory Data Analysis (Contd.)



Weekly Sales and Customers are showing almost similar trends. Both peak at mid-December.

# Exploratory Data Analysis (Contd.)



- 'Customers' and 'Sales' are highly correlated.
- We can see that 'WeekOfYear' and 'Month' are also highly correlated.
- 'Open' is moderately correlated with 'Sales' and 'Customers'.
- 'Promo2' is correlated with 'Promo2SinceWeek' and 'Promo2SinceYear'.
- 'CompetitionOpenNumMonths' is highly correlated with 'CompetitionOpenSinceYear'.
- 'Promo2NumWeeks' is moderately correlated with 'Promo2', 'Promo2SinceWeek' and 'Promo2SinceYear'.

# Feature Engineering

- 'Store' is dropped since sales can be predicted through store type, assortment, etc.

- 'Date' is also dropped since there is already day of week and week of year features in the dataset.

- 'CompetitionOpenSinceMonth' and 'CompetitionOpenSinceYear' are dropped as the information provided by them can be obtained from 'CompetitionOpenNumMonths'.

- 'Promo2', 'Promo2SinceWeek' and 'Promo2SinceYear' are also dropped as the information provided by them can be obtained from 'Promo2NumWeeks'.

- 'Month' is dropped since we get the same information from 'WeekOfYear'.

- 'Year' is also dropped as we have already established in EDA that it's not the year that influence the sales but 'DayOfWeek' and 'WeekOfYear'.

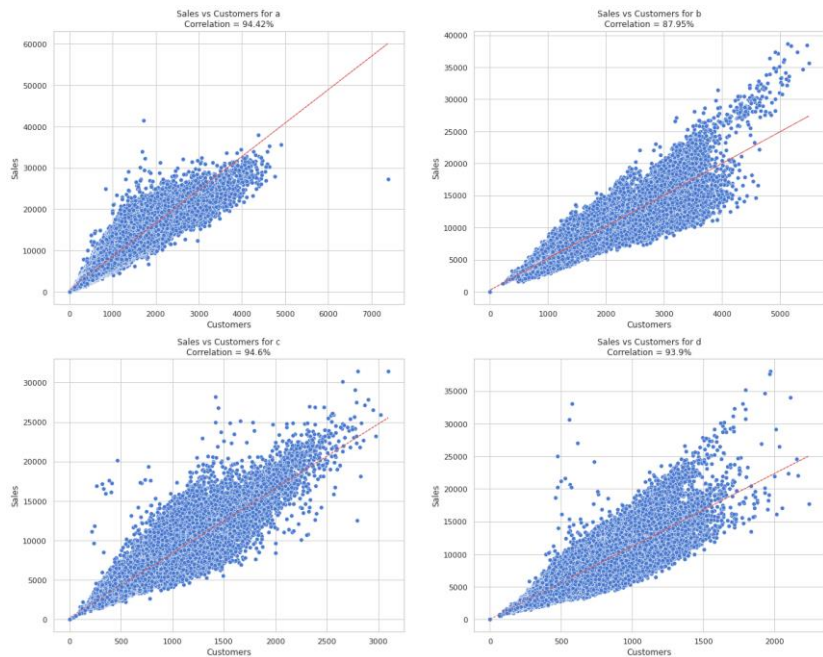# Feature Engineering (Contd.)

- The variance inflation factor (VIF) of all numerical features except 'Sales' is calculated in order to remove highly correlated features.

- Features having VIF greater than 5 should be eliminated.

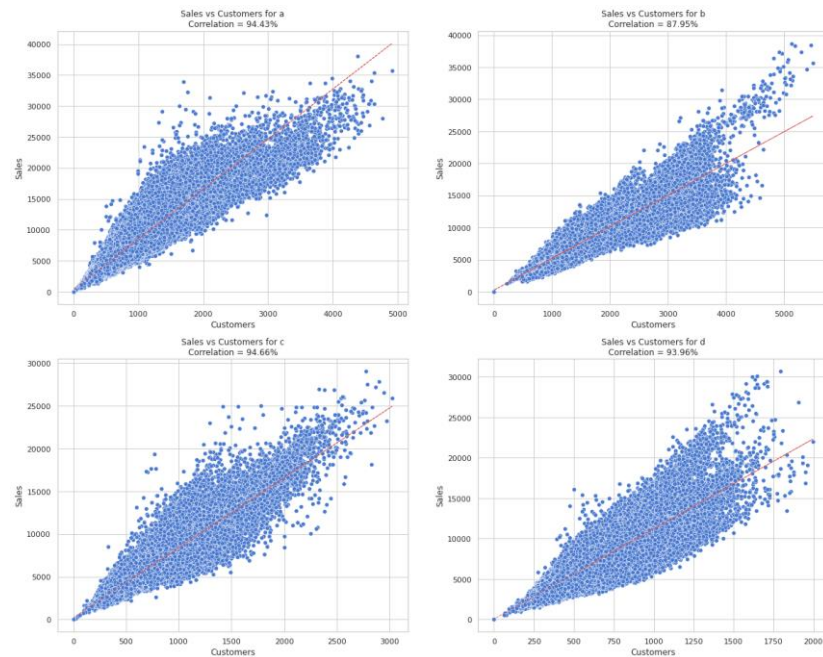| | Feature | VIF |
|---|---|---|
| 0 | Customers | 1.568147 |
| 1 | CompetitionDistance | 1.306160 |
| 2 | CompetitionOpenNumMonths | 1.525657 |
| 3 | Promo2NumWeeks | 1.235789 |

- All features have VIF less than 5.

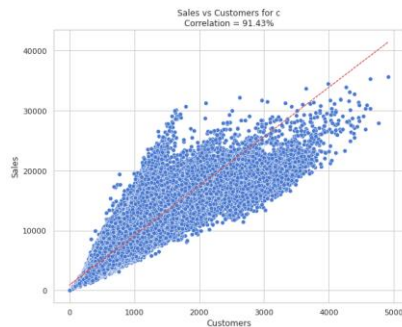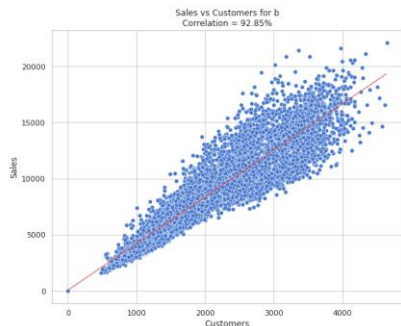# Feature Engineering (Contd.)

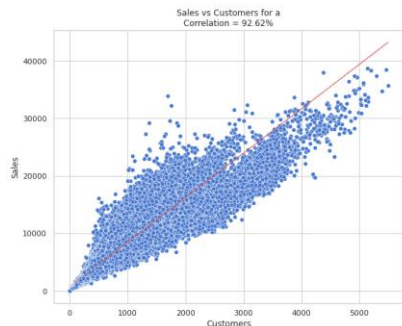**Before removing outliers in 'StoreType'**
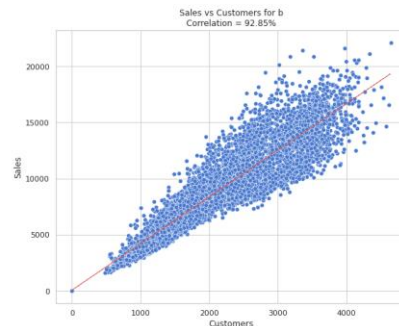
**After removing outliers in 'StoreType'**

# **Feature Engineering** (Contd.)

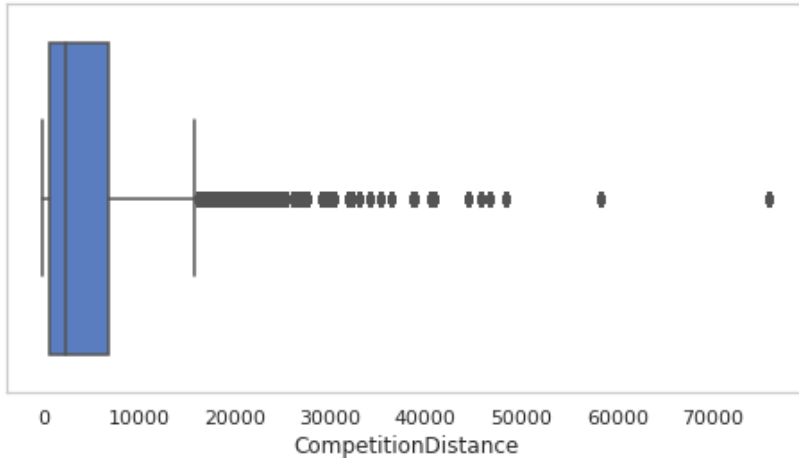**Before removing outliers in 'Assortment'**
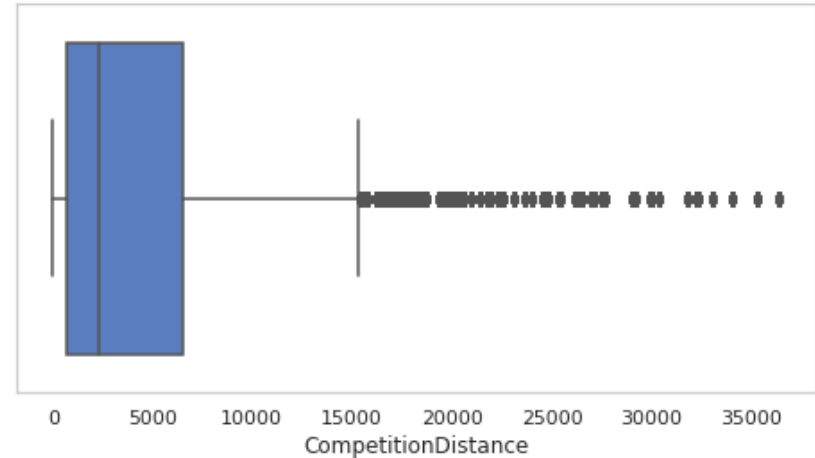
**After removing outliers in 'Assortment'**

# Feature Engineering (Contd.)
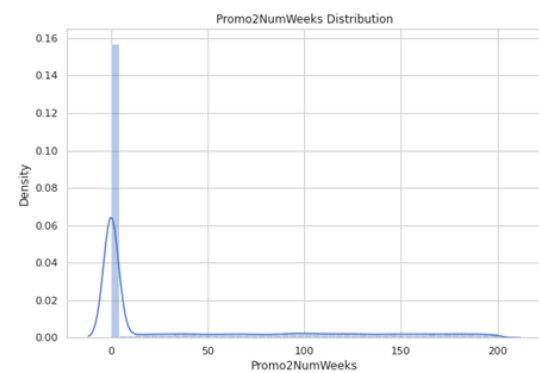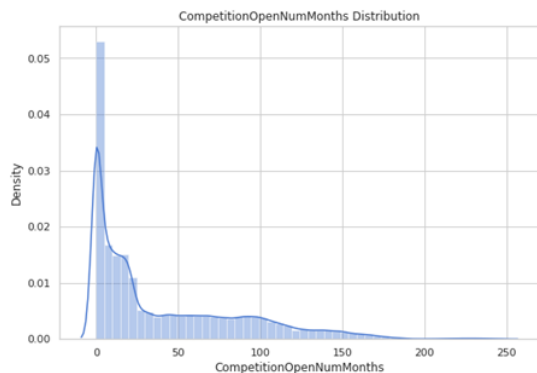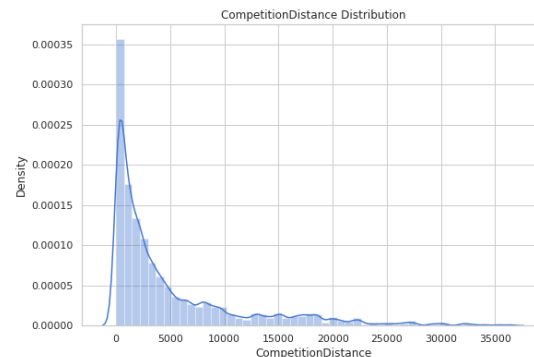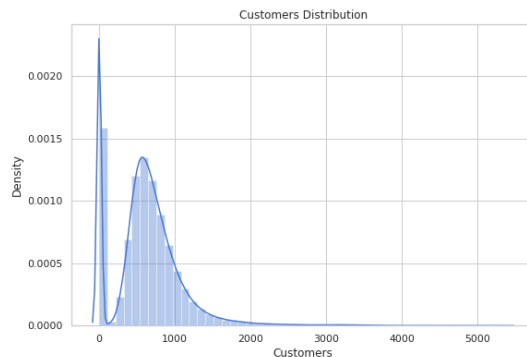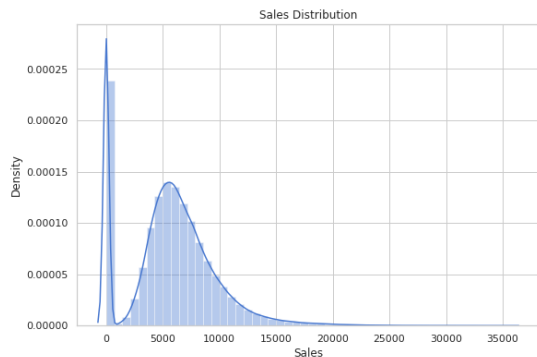
**Before removing outliers in 'CompetitionDistance'**



**After removing outliers in 'CompetitionDistance'**



- Almost 10% of data are outliers and removing them will cause the loss of useful information.
- So, only the outliers having values above 37500 were removed.

# Feature Engineering (Contd.)



All features are positively skewed but linear regression assumes normal distribution.

# Feature Engineering (Contd.)



'Sales', 'Customers', 'CompetitionOpenNumMonths' and 'Promo2NumWeeks' are square root transformed and 'CompetitionDistance' is log transformed.

# Feature Engineering (Contd.)

- Most algorithms cannot handle the categorical variables unless they are converted into a numerical value.

- 'StoreType' and 'Assortment' were encoded using one hot encoder while 'PromoInterval' was dummified.

- To overcome dummy variable trap, one resultant feature from each encoded feature must be removed. Correlation matrix was used to decide which features to remove.

- 'PromoInterval_Jan,Apr,Jul,Oct', 'StoreType_c' & 'Assortment_b' have the least correlation to Sales and they were removed.

# Modelling

- Input and target data were separated, and both were split into training and test data with 25% test data.

- The feature 'Customers' was also removed since the number of customers for the period, which is under consideration for forecasting, won't be available until the mentioned period is over.

- Training and test data of independent features were scaled using standardization.

- Model training was done with these data using 7 different algorithms:

    1. Linear regression
    2. Ridge regression
    3. Lasso regression
    4. Elastic net regression
    5. Decision tree regression
    6. Random forest regression
    7. XGBoost regression

# Modelling (Contd.)

| | Regression_Model | Train_R2 | Test_R2 | Train_RMSE | Test_RMSE | Train_RMSPE | Test_RMSPE |
|---|---|---|---|---|---|---|---|
| 0 | Linear | 0.834639 | 0.833716 | 14.050363 | 14.062944 | 20.774380 | 20.743497 |
| 1 | Ridge | 0.834639 | 0.833716 | 14.050363 | 14.062943 | 20.774380 | 20.743495 |
| 2 | Lasso | 0.834639 | 0.833716 | 14.050363 | 14.062944 | 20.774380 | 20.743496 |
| 3 | Elastic Net | 0.834639 | 0.833716 | 14.050363 | 14.062943 | 20.774380 | 20.743495 |
| 4 | Decision Tree | 0.848202 | 0.847136 | 13.461839 | 13.483542 | 19.904209 | 19.888852 |
| 5 | Random Forest | 0.872353 | 0.870884 | 12.344587 | 12.391997 | 18.252278 | 18.278772 |
| 6 | XGBoost | 0.960901 | 0.958972 | 6.832110 | 6.985427 | 10.101721 | 10.303830 |



- Evaluation metrics like R-squared, RMSE, RMSPE, etc. were calculated for each model.
- R2 score can be used to compare different models and find out which one gives higher accuracy. Higher the R2 score, higher the accuracy.
- The model built using XGBoost algorithm has the highest R2 score with ~96%, followed by the one using random forest and decision tree.

# Conclusion

## EDA Conclusions

- Mondays have most sales since most of the Sundays are closed.
- Promotions seem to have a significant effect on sales but not for the number of customers.
- Store type b has higher sales and customers per store than other store types.
- Assortment b is available only at store type b and it has more average sales and customers than any other assortment.
- Weekly sales and customers peak at the mid-December. Best guess is that people buy drugs in advance just before the shops close for the holidays.

## Modelling Conclusions

- The model built using XGBoost algorithm gives unusually high accuracy. This may lead to overfitting. Therefore, it is advisable to not use this model.
- Among the remaining, the model built using random forest algorithm is the most accurate one.
- If model interpretability is more important than accuracy, model built using decision tree algorithm should be chosen over the one using random forest algorithm.
- Decision tree based algorithms are slightly more accurate than linear regression based algorithms.