# Capstone Project
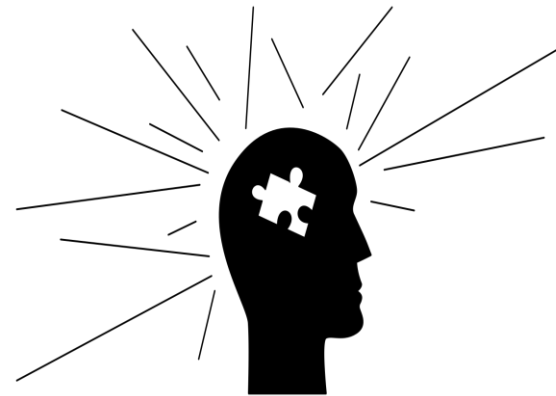## Rossmann Sales Prediction

ROSSMANN

# Team Consists of:

1. Vridhi Parmar
2. Pradip Solanki
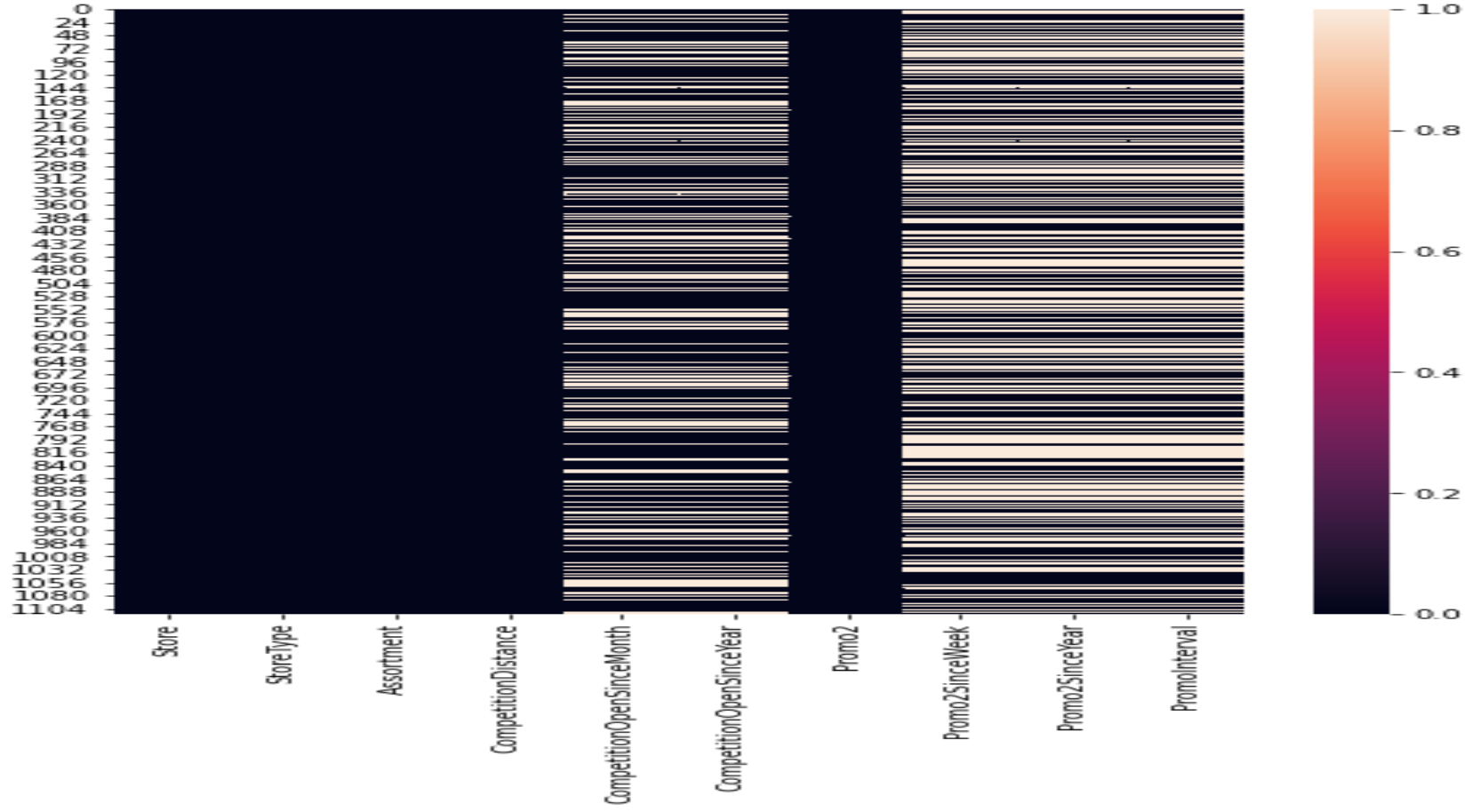3. Hrithik Chourasia
4. Ameen Attar

# Problem Statement:

- Rossmann store managers are tasked with predicting their daily sales for up to six weeks in advance.

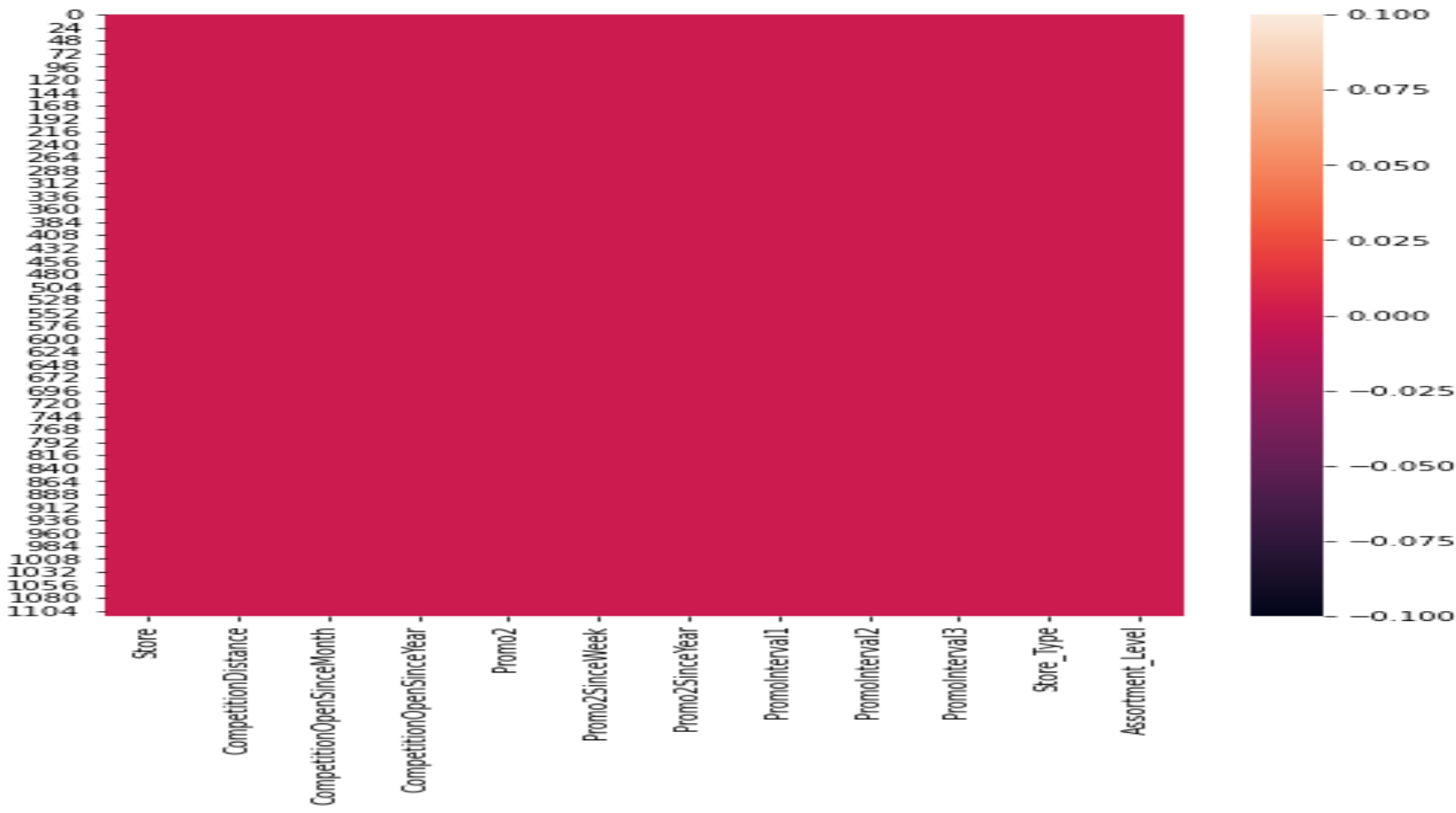- We were provided with historical sales data for 1,115 Rossmann stores and were asked to predict the same

# Description of data provided

- We are provided with 2 data sets:
- **1. Rossmann Stores Data.csv –** This dataset includes the historical data including Sales. This dataset contain features like Sales, Customers, Open, StateHoliday, SchoolHoliday.

- **2. store.csv –** This includes supplemental information about the stores. This dataset contain features like Assortment, CompetitionDistance, CompetitionOpenSince[Month/Year], Promo, Promo2, Promo2Since[Year/Week]
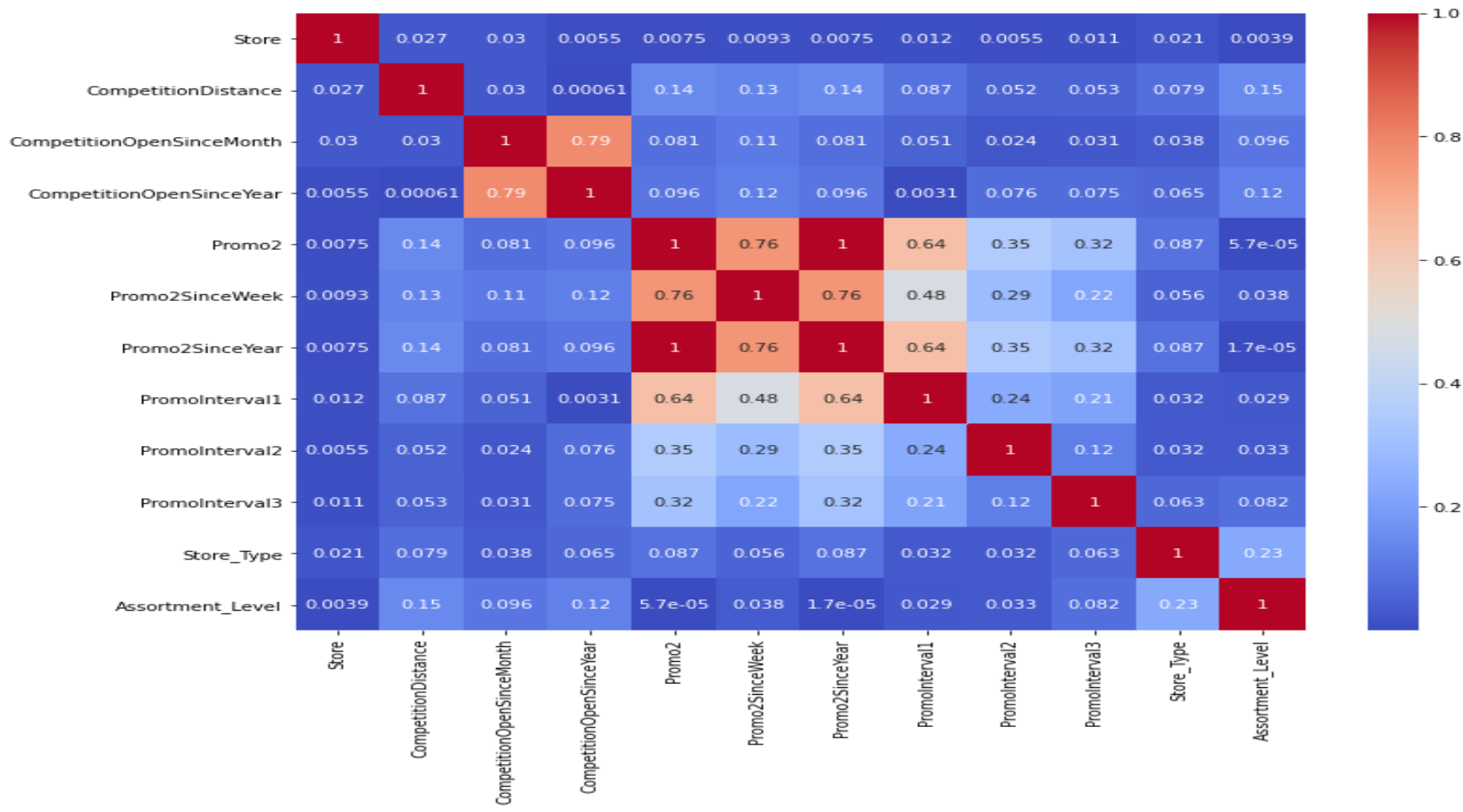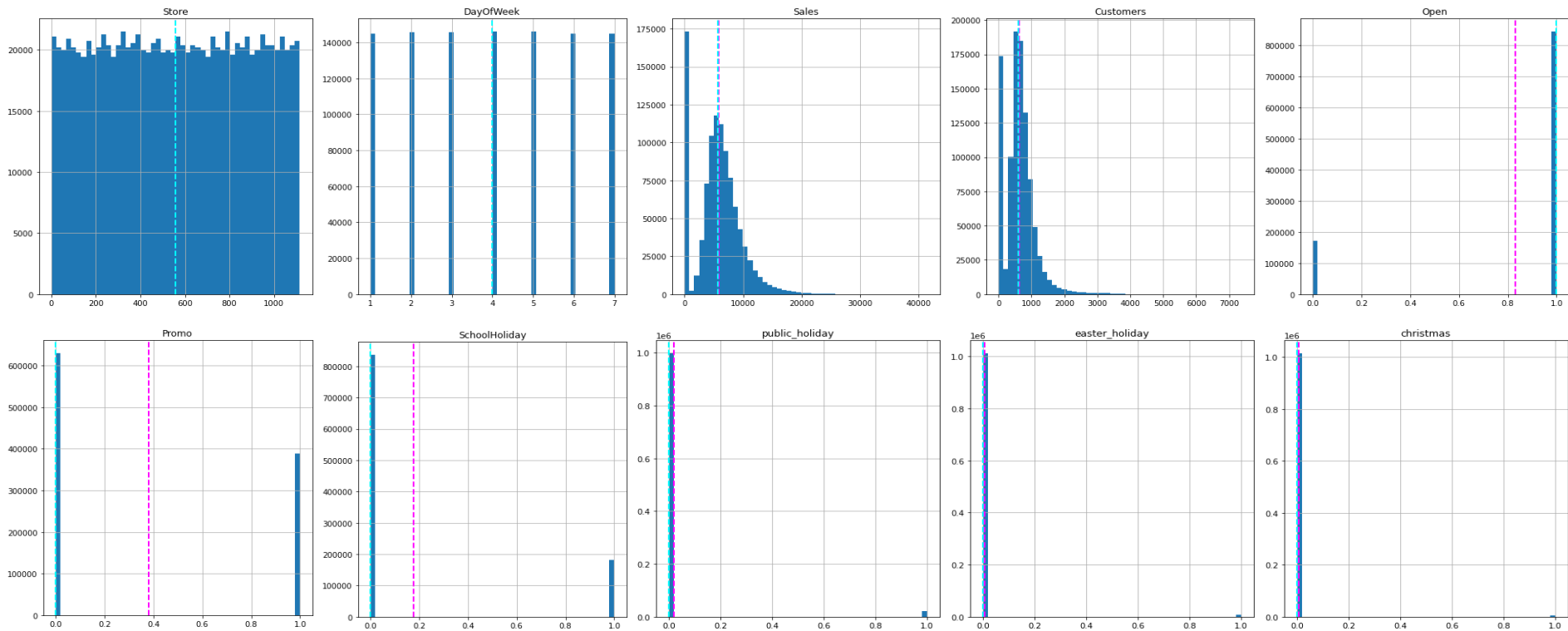
Missing Values in the Datasets
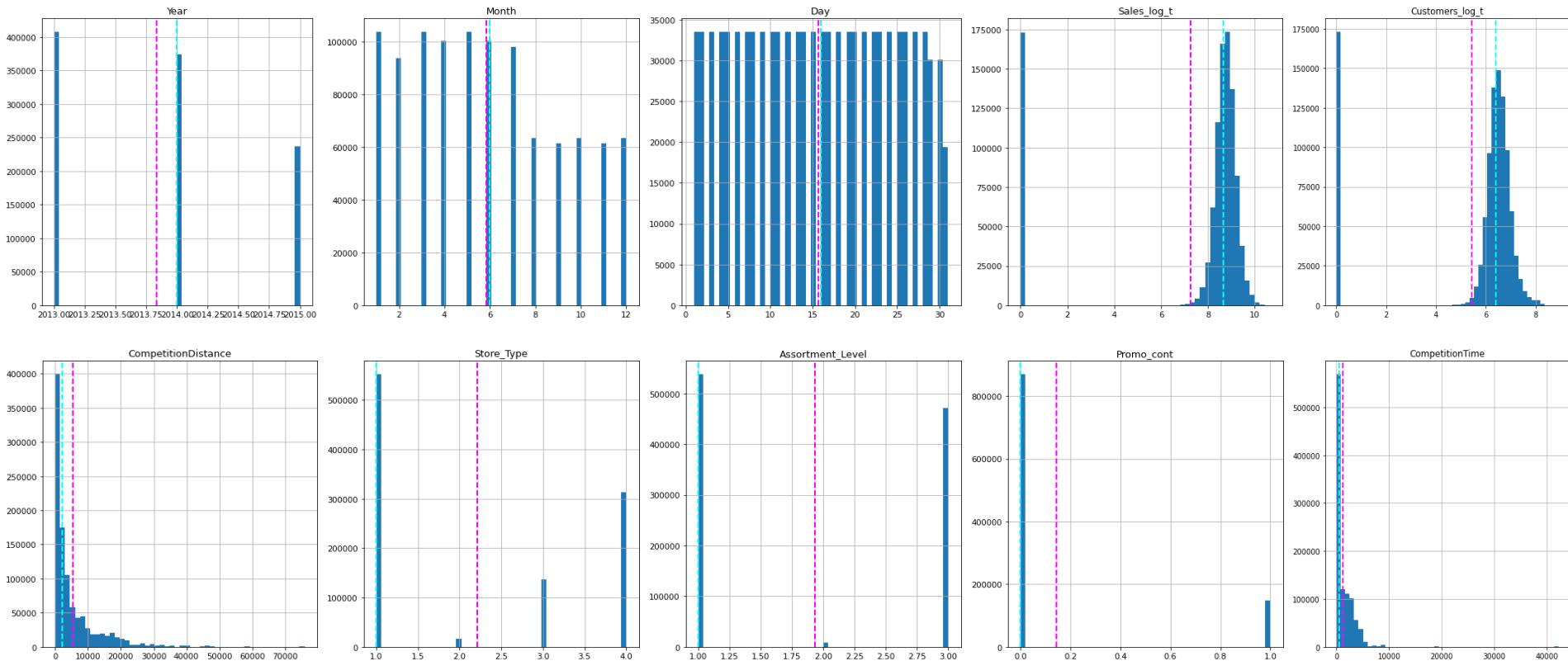
# After filling NaN Values in the Datasets
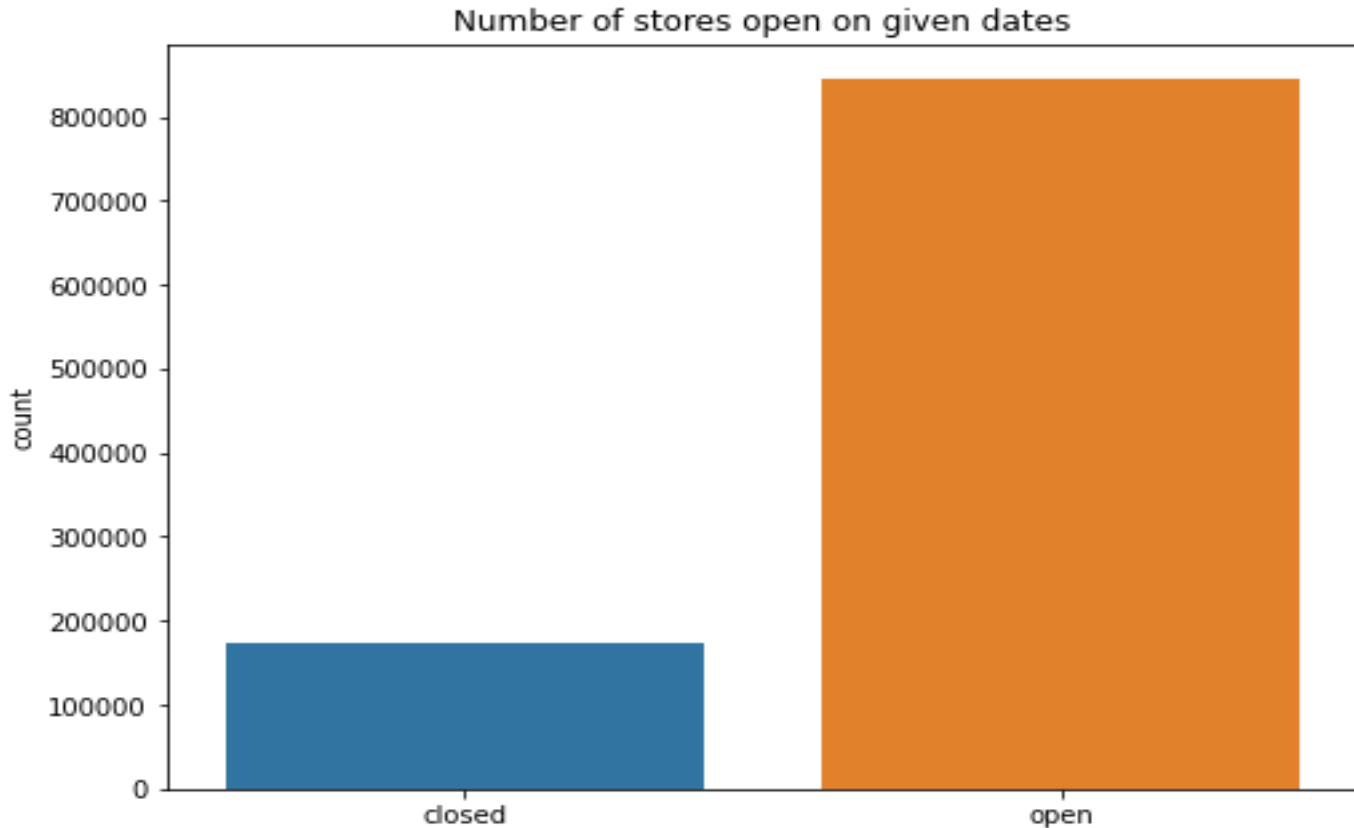
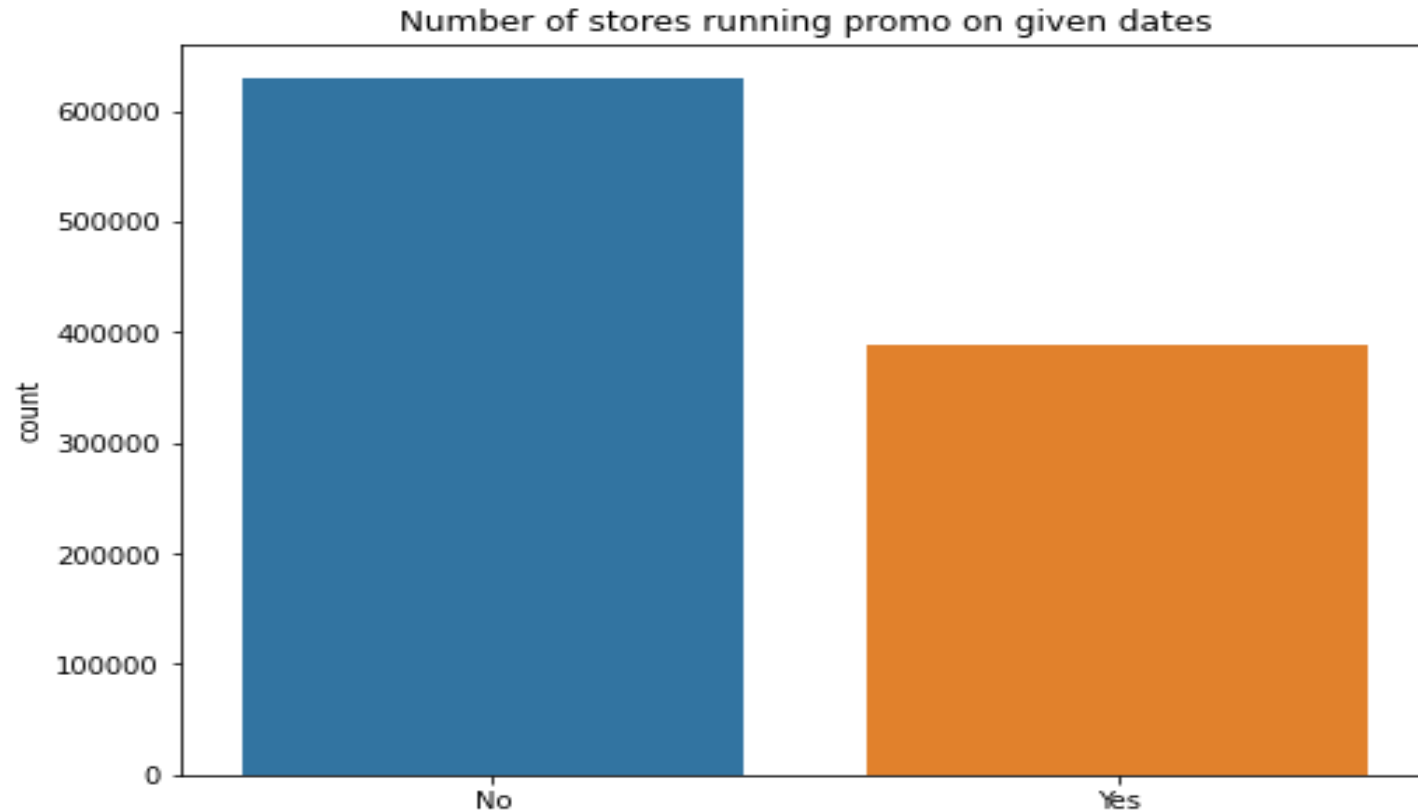# Correlation of variables with each others

# Plot of all the variables:

# Plot of all the variables:

# Numbers of stores open on a given Dates



Number of stores open on given dates

# Number of stores running promo on given Dates



Number of stores running promo on given dates
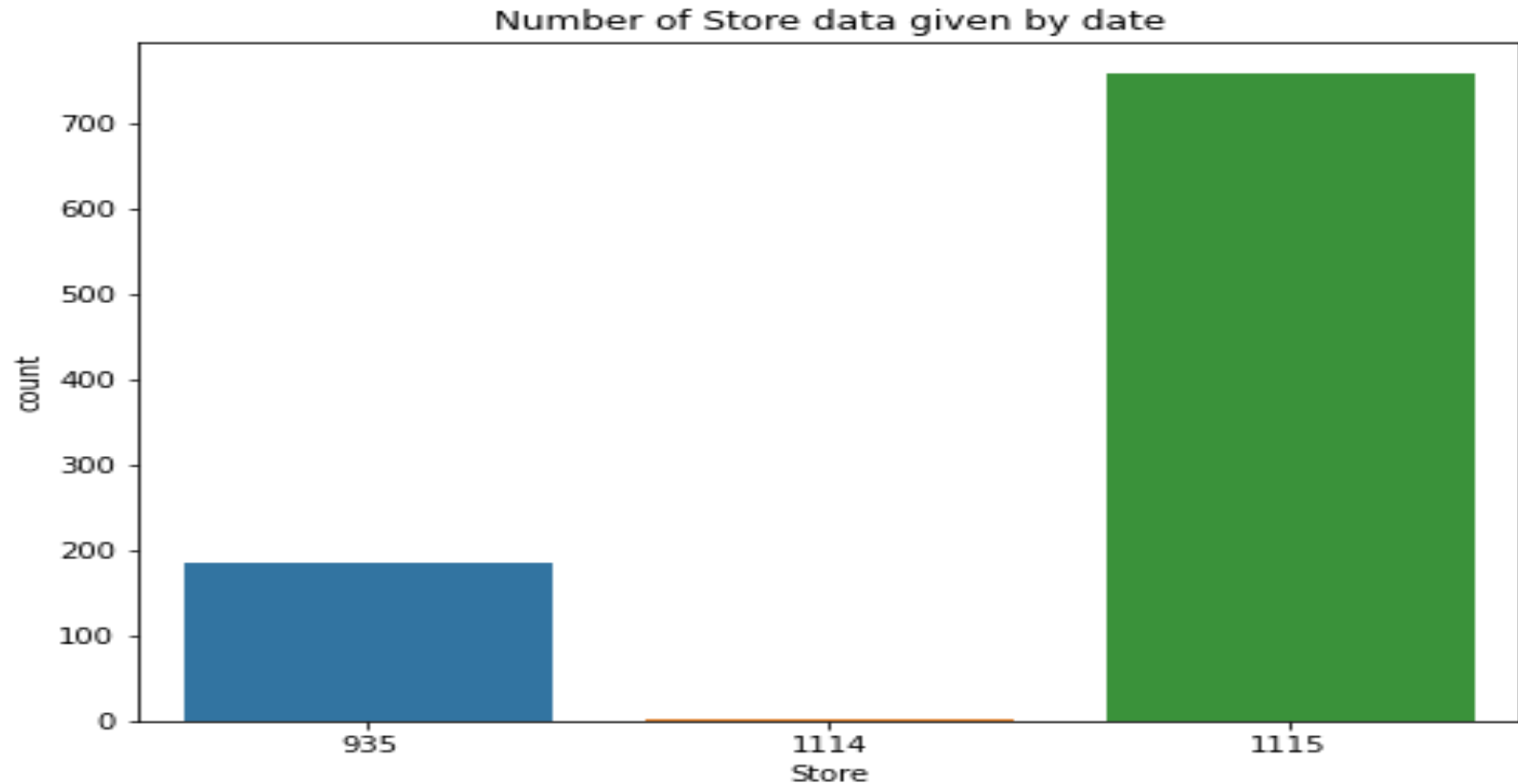
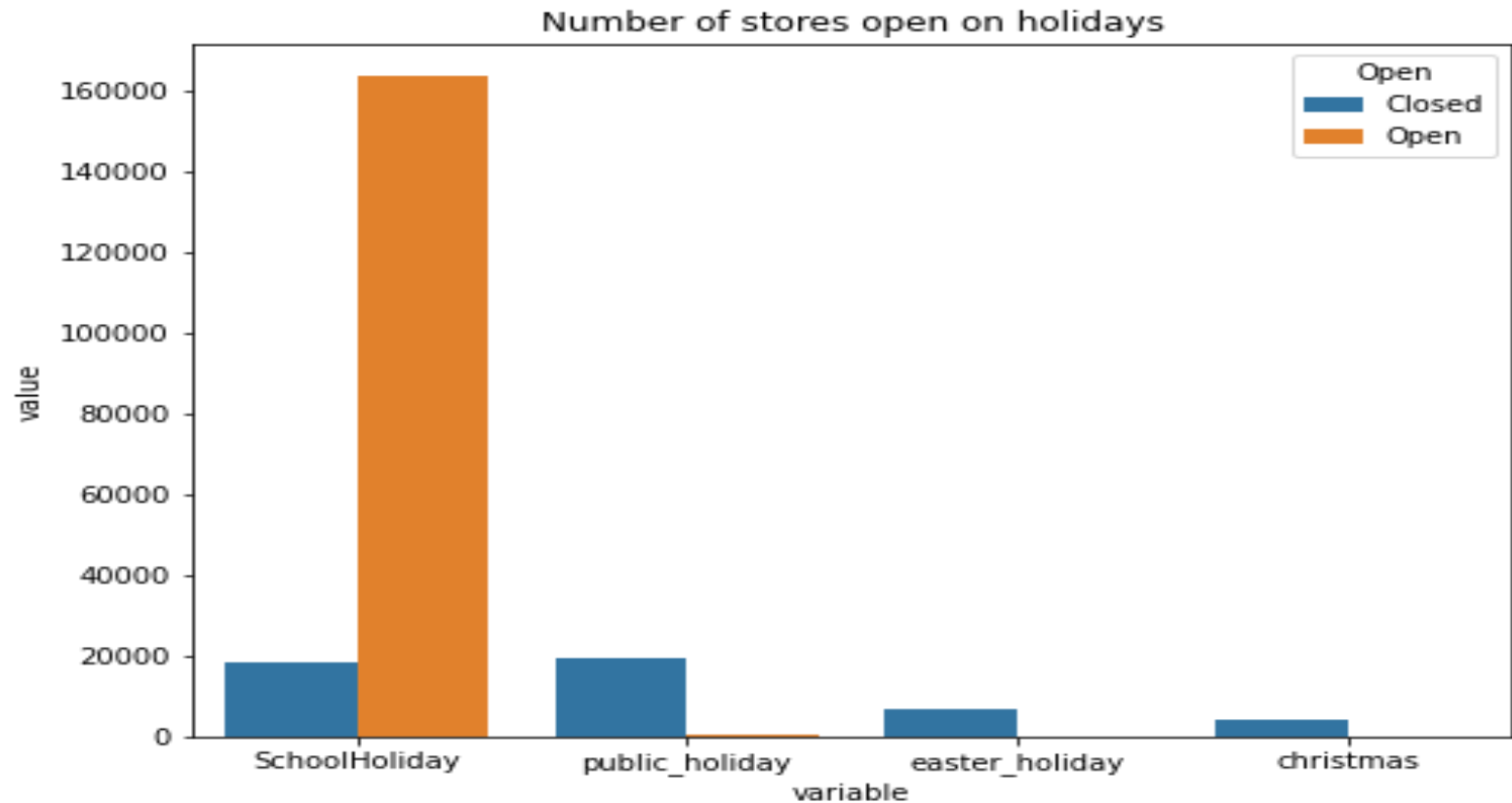# Number of school Holidays

**AI**

Number of days school are closed

# Number of Holidays celebrated by stores

# Number of store data given by Date



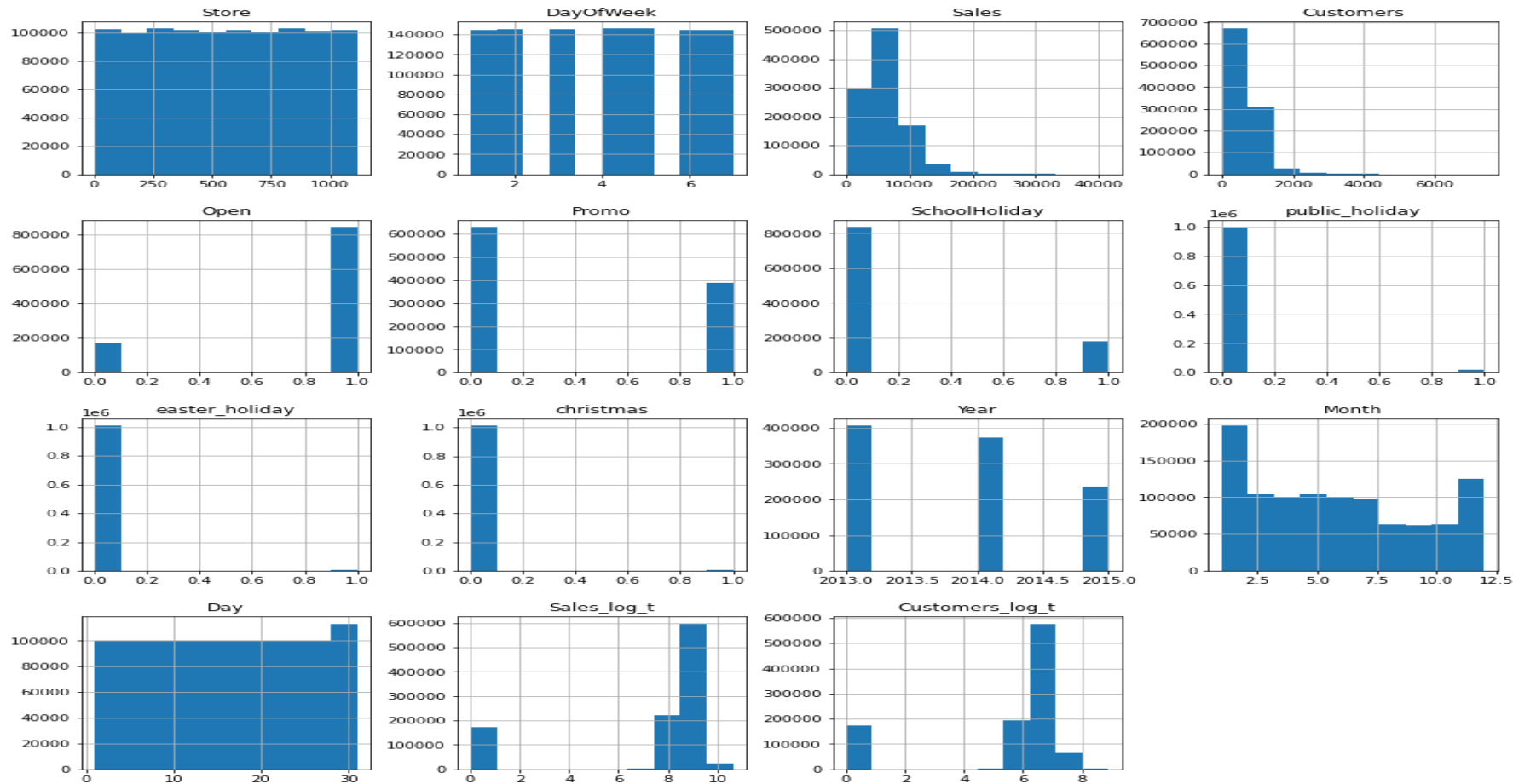Number of Store data given by date

# Number of stores open on Holidays

# After removing the skewness

# Correlation matrix of Rossmann's Dataset



Correlation matrix

# Correlation after merging the two dataset



Correlation matrix

# The scatter plot for sales_log_t with other Variables:

# The scatter plot for sales_log_t with other Variables:

# The scatter plot for sales_log_t with other Variables:

# The scatter plot for sales_log_t with other Variables:

# Plot for first 100 observations between actual and predicted after building the linear regression model:

# We are using linear regression model so, we need to check 4 basic assumptions of linear regression.

1. There should to be linear relationship between independent and dependent variables.
2. The sum of residuals/error should be near to 0.
3. There should not be multicollinearity.
4. There should not be heteroscedasticity.

# Assumptions:

1. The linear relationship between independent and dependent variables.



2. The sum of residuals/error should be near to 0, and the sum of residuals we have got is –0.0007.

# 3. There should not be multicollinearity.

| | | |
|---|---|---|
| 1 | DayOfWeek | 4.452981 |
| 2 | Open | 6.191915 |
| 3 | Promo | 1.872410 |
| 4 | SchoolHoliday | 1.342025 |
| 5 | public_holiday | 1.099311 |
| 6 | easter_holiday | 1.074031 |
| 7 | christmas | 1.072576 |
| 8 | Month | 3.904310 |
| 9 | Day | 3.965944 |
| 10 | CompetitionDistance | 1.530442 |
| 11 | Store_Type | 3.705991 |
| 12 | Assortment_Level | 4.908106 |
| 13 | Promo_cont | 1.176941 |
| 14 | CompetitionTime | 1.408883 |

# Assumptions:

4. Checking the Heteroscedasticity.

# Evaluation metrics for our multiple regression models used while model building:

**AI**

**Linear Regression model:**
Mean Squared Error : 0.11
Root Mean Squared Error : 0
Mean Absolute Error : 0.25
Mean Absolute Percentage Error : 3.41 %
R-Square : 0.99
Adjusted R-Square :  0.99

**Lasso Regression with cross validation:**
Mean Squared Error : 0.11
Root Mean Squared Error : 0
Mean Absolute Error : 0.25
Mean Absolute Percentage Error : 3.41 %
R-Square : 0.99
Adjusted R-Square :  0.99

**ElasticNet Regression:**
Mean Squared Error : 0.11
Root Mean Squared Error : 0
Mean Absolute Error : 0.25
Mean Absolute Percentage Error : 3.41 %
R-Square : 0.99
Adjusted R-Square :  0.99

**Ridge Regression with cross validation:**
Mean Squared Error : 0.11
Root Mean Squared Error : 0
Mean Absolute Error : 0.25
Mean Absolute Percentage Error : 3.41 %
R-Square : 0.99
Adjusted R-Square :  0.99

# Stack Model:

Stacking is a good way to combine all the predictions from different models into one. We can adjust weights for each model in stacking.

**Evaluation metrics of our Stack model:**

Mean Squared Error : 0.11

Root Mean Squared Error : 0

Mean Absolute Error : 0.25

Mean Absolute Percentage Error : 3.41 %

R-Square : 0.99

Adjusted R-Square : 0.99
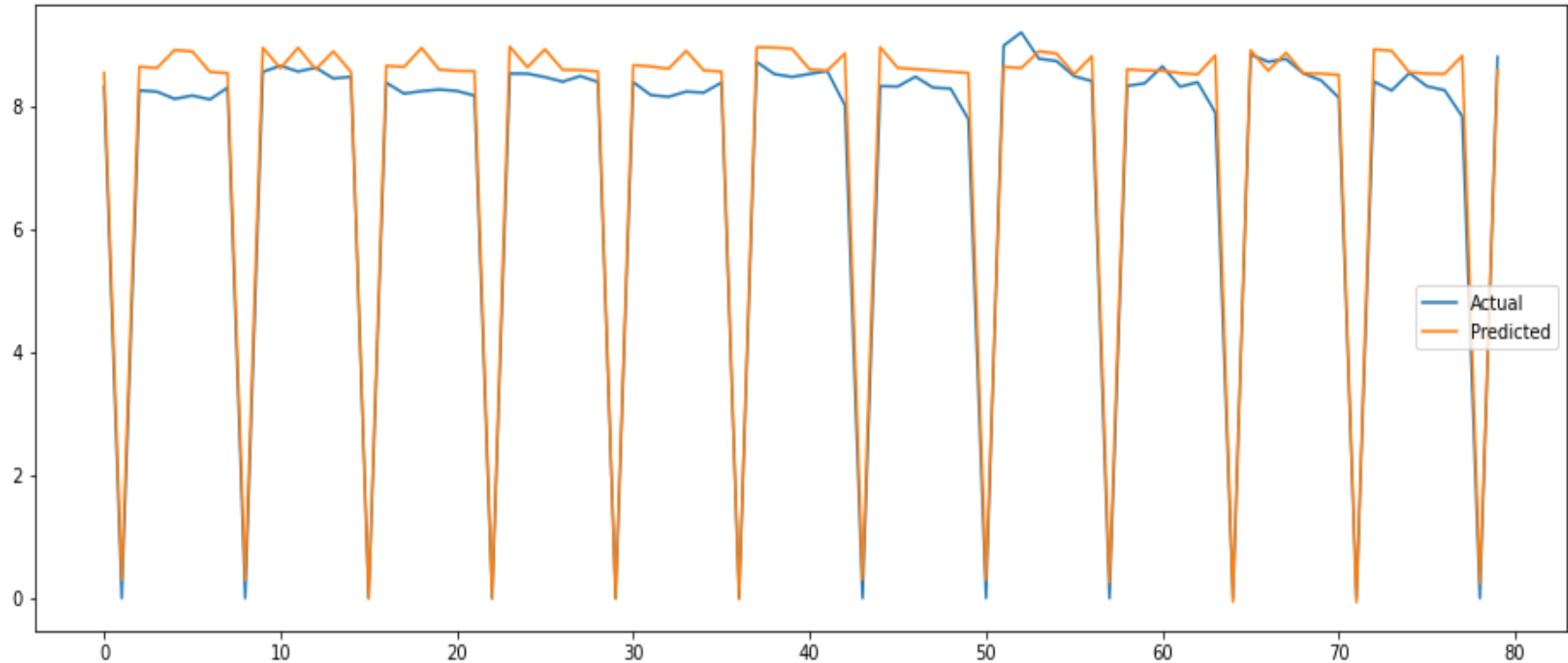
# Graph for Stack Model:

# Predicting Sales for next six weeks:

Now that we have linear regression models performing well.
we need to predict Sales for next 6 weeks in Advance.
This is what our problem statement states and hence we have
considered the last 42 days as our test data.

# Predicting Sales for next six weeks:

# Challenges:

1. Here we had to deal with the larger datasets. We had to merge two datasets into one to get the meaningful insights from data.

2. We did a lot of feature engineering in order to get useful information from multiple columns.

3. We also applied log transform to a lot of columns to avoid skewness which helped us to predict our values well.

4. We generated the data for next six weeks using the test data which was the complex part in our project.

# Conclusion:

The Rossmann store sales prediction is very engrossing data science problem to solve. We noticed that the problem is more concentrated towards the feature engineering and the feature selection part than on model selection. We had to spend around 60-70% of our time on analyzing data for trends in order to make our feature selection easier.

As we are building a liner regression model we emphasized on the basic 4 assumptions of a liner regression model.

Con**clus**ion...

# Conclusion:

- We also have applied regularization techniques like Lasso, Ridge and Elastic Net to Avoid Overfitting.

- We also used Stacking to make predictions that have better performance than any single model.

- Most important feature came out to be customers, where sales is directly related to number of customers.

- The graph above is what our final model looks like at the end. As per our model prediction we can conclude that the Total revenue of Predicted Sales in next 6 weeks will be: 250,776,406 euros By 1115 Rossmann stores.