

Restaurant Revenue Prediction

Using Regression Techniques

Kavya Mistry, MS in Data Science, Rowan University

1. Introduction

Anything that comes under food industry plays a vital role for the development of a country. Especially restaurants have large amount of share in the societal growth and economic boost. According to Forbes, with approximately 11 million jobs and over 4% of the GDP attributable to the restaurant industry, one could make a very well-founded argument that “as restaurants go, so goes the economy.”

Restaurants being such crucial into business sector comes with its own challenges, like choosing location, type of restaurant, tough competitions in market, varying cultures, subjective personal judgement, inventory decisions, and other demographic, commercial and real estate factors. Fortunately, most of these challenges could be tackled using developing technologies in Machine learning and Artificial Intelligence. The focus of this project/paper is to explore with various methods to create an efficient Revenue prediction model for a restaurant using Data Mining applications.

2. Business Problem and Tasks

The main aspect of this challenge is to determine revenue for a restaurant and obtain significant insights on whether that restaurant will be profitable or will it incur loss of money in the near future.

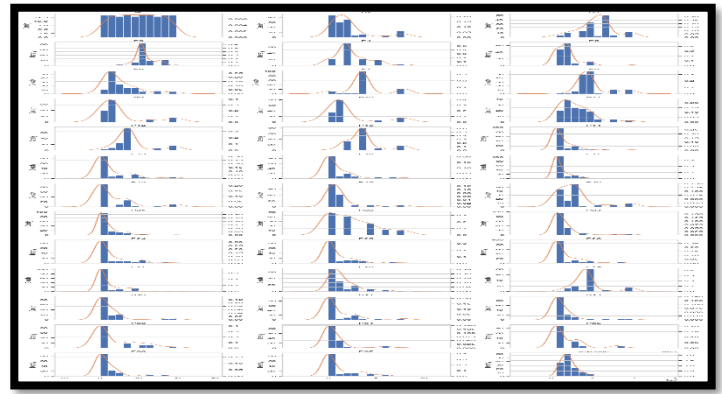
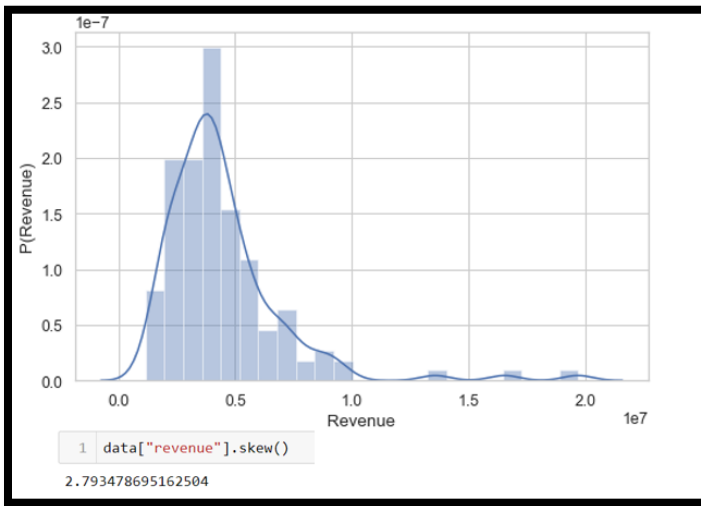
The steps involved to achieve the above goal are Data-preprocessing, Exploratory Data Analysis (EDA), Data Visualization, Feature Engineering, Training and Testing of machine learning model. As this is supervised regression problem, the paper also talks about evaluation of various Regression machine learning models to check

their performance for the dataset.

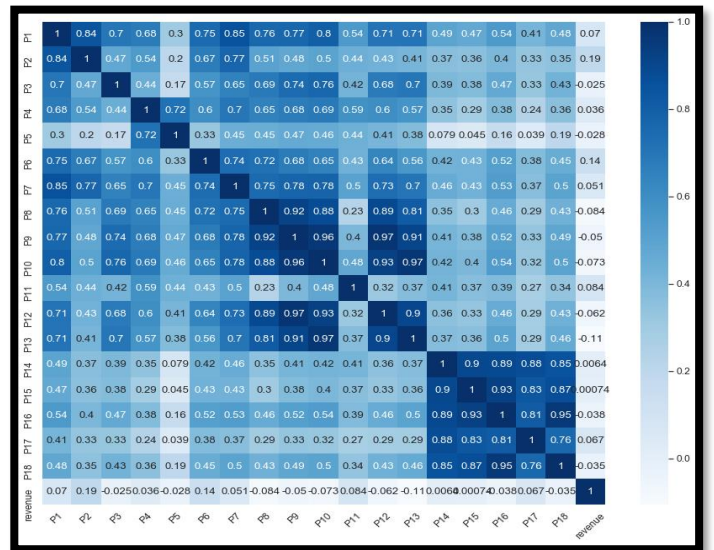
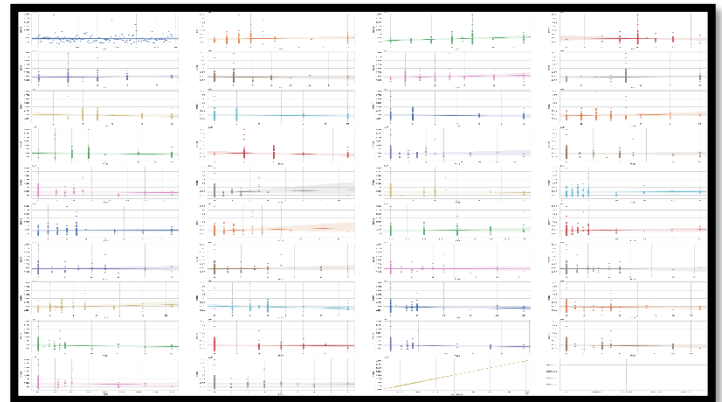
3. Dataset

For this project, I took the dataset from one of the Kaggle competition. With over 1,200 quick service restaurants across the globe, TFI is the company behind some of the world's most well-known brands: Burger King, Sbarro, Popeyes, Usta Donerci, and Arby's. TFI has provided a dataset with 137 restaurants in the training set, and a test set of 100000 restaurants. The description of the data fields is stated below:

- Id: Restaurant id.
- Open Date: opening date for a restaurant
- City: City that the restaurant is in.
- City Group: Type of the city. Big cities, or Other.
- Type: Type of the restaurant. FC: Food Court, IL: Inline, DT: Drive Thru
- P1, P2 - P37: There are three categories of these obfuscated data. Demographic data are gathered from third party providers with GIS systems. These include population in any given area, age and gender distribution, development scales. Real estate data mainly relate to the m2 of the location, front facade of the location, car park availability. Commercial data mainly include the existence of points of interest

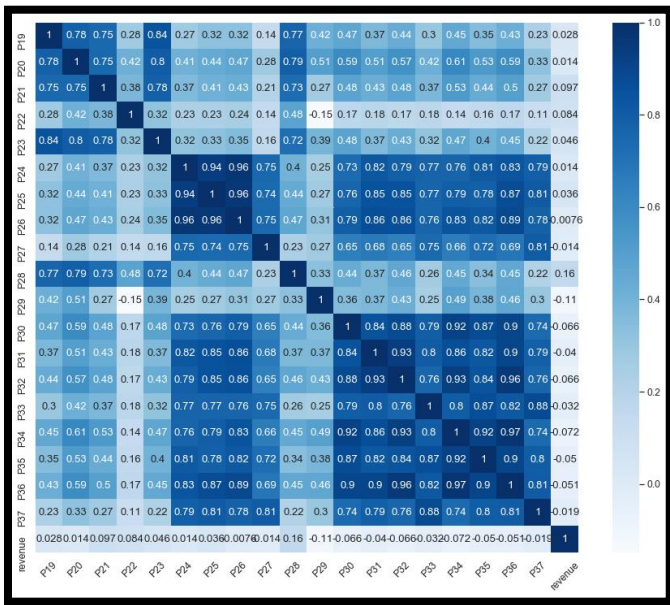


Next, to analyze the relationship between Revenue and these numerical attributes, I plotted a graph using Regression plot and Correlation plot.



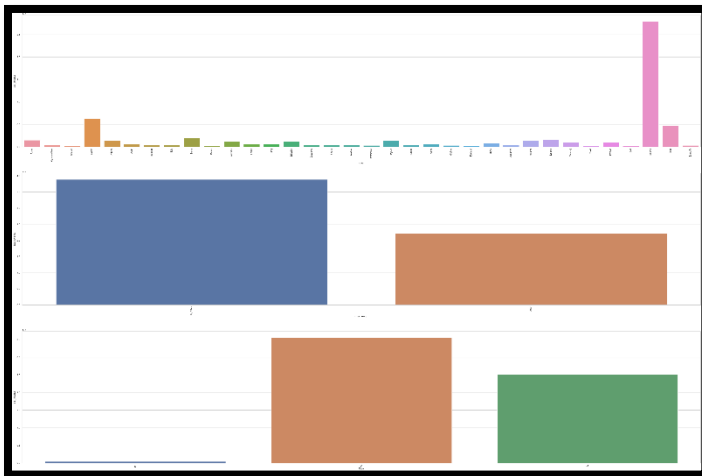
Moreover, I also checked for the distribution of P1 to P37 via Histogram and from below distributions it can be observed that there are clear anomalies in some of the features in the numerical aspect of the data, the suitable method is to study these anomalies with an expert in the field, also the distributions describe certain biases that made that data distribution either not fair or skewed for some reason.

For instance, if we look at P37 value zero is more likely to show among the feature values, possibly this is a data entry error or is there something important that number zero represents?



As the heatmap and the scatter distribution clearly shows that there is no obvious relation between numerical variables independently with the target.

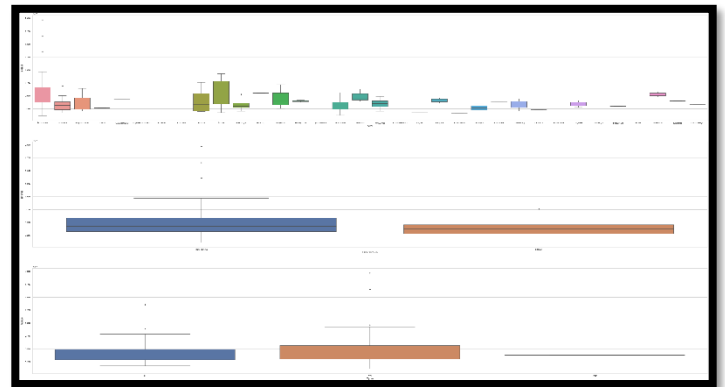
5.2 Analyzing Categorical Features Vs Revenue



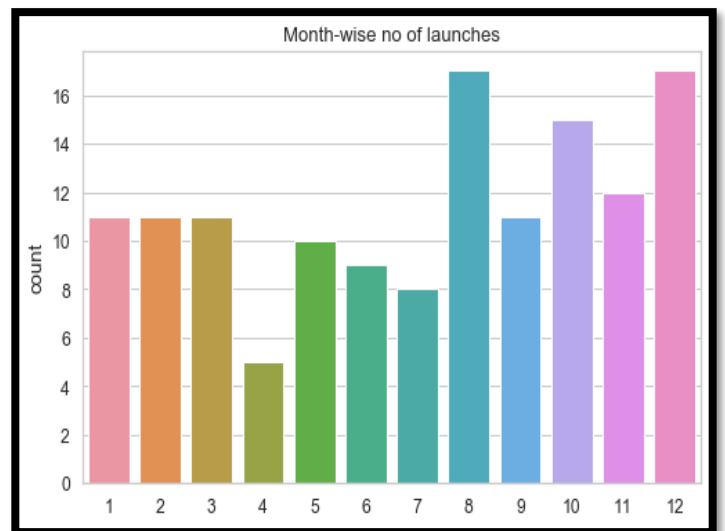
Using Matplotlib library, insights obtained from categorical variables charts shows that Istanbul is making the most revenues compared to other cities, hence opening a restaurant in Istanbul could be a nice idea. It is also visible that big cities are making more revenues than smaller cities. It can also be noted that

the gap between them is not an obvious one that is much bigger, may be because even small cities have good tourist attraction. Moreover, if you are about to open a restaurant it better be either FC (Food court) or IL (inline).

I also made a boxplot to visualize the anomalies in the categorical attributes.

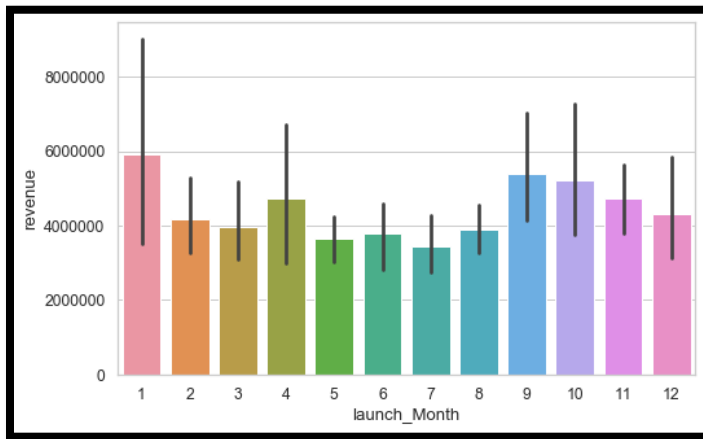


5.3 Analyzing Open Date Column



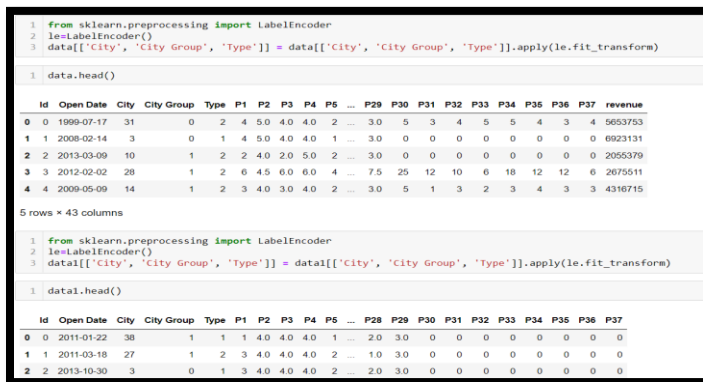
Using Seaborn library, From the count plot in the above figure, seems like the first three months (January to March) and ending months (August to December) have more launches.

7. Model Building & Evaluation



From the bar plot in the right figure, it can be seen that the months that are having more revenue are: January, September, October, November and December. Hence, it's safe to say that one can launch their restaurants in these months to have a better chance of generating profits!

6. One Hot Encoding



Before feeding the data to the model, the categorical data should be converted to numbers, so that we can feed these variables in machine learning algorithms using sklearn. For this I imported Label Encoder and did one-hot encoding on City, City type and Type for both train and test data.

For this purpose, we have training data and testing data separate files. The model is trained on train dataset and to predict the revenue we are using testing dataset. There is another sample file which already has results according to the independent variables. The output obtained from testing data is stored in data frame so that one can compare those predicted results and sample revenue values. Through this we can know about the accuracy of the model, and what fits best to the model.

To train the dataset, numerous regression models were used. For evaluation purpose, Root Mean Square Error (RMSE) was used. Firstly, logistic regression was trained with the below results:

```

1 logisticRegr.score(x_train,y_train)

1.0

1 from sklearn.metrics import mean_squared_error
2 from math import sqrt
3 label_list=y_pred_label['Prediction'].tolist()
4 #print('Root Mean squared error {}'.format(sqrt(mean_squared_error(label_list, predictions))))
5 rmse1=round(sqrt(mean_squared_error(label_list, predictions)),2)
6 print(rmse1)

2366286.68

```

Secondly, Gradient Boosting Regression was applied and below were the results:

observed:

```
1 gbRegr.score(x_train, y_train)

0.9403860943246145

1 from sklearn.metrics import mean_squared_error
2 from math import sqrt
3 label_list=y_pred_label['Prediction'].tolist()
4 #print('Root Mean squared error {}'.format(sqrt(mean_squared_error(label_list, predictions))))
5 rmse2=round(sqrt(mean_squared_error(label_list, prediction_rr)),2)
6 print(rmse2)

655991.78
```

```
1 neigh.score(x_train, y_train)

0.49045279578697765

1 from sklearn.metrics import mean_squared_error
2 from math import sqrt
3 label_list=y_pred_label['Prediction'].tolist()
4 #print('Root Mean squared error {}'.format(sqrt(mean_squared_error(label_list, predictions))))
5 rmse5=round(sqrt(mean_squared_error(label_list, prediction_knn)),2)
6 print(rmse5)

2346324.81
```

Afterwards, with Linear Regression below results were seen:

```
1 reg.score(x_train, y_train)

0.3260403831258316

1 from sklearn.metrics import mean_squared_error
2 from math import sqrt
3 label_list=y_pred_label['Prediction'].tolist()
4 #print('Root Mean squared error {}'.format(sqrt(mean_squared_error(label_list, predictions))))
5 rmse3=round(sqrt(mean_squared_error(label_list, prediction_lr)),2)
6 print(rmse3)

218217353.19
```

With Random Forest Regressor below was the output:

```
1 regr.score(x_train, y_train)

0.30770677147819403

1 from sklearn.metrics import mean_squared_error
2 from math import sqrt
3 label_list=y_pred_label['Prediction'].tolist()
4 #print('Root Mean squared error {}'.format(sqrt(mean_squared_error(label_list, predictions))))
5 rmse4=round(sqrt(mean_squared_error(label_list, prediction_rf)),2)
6 print(rmse4)

1538357.74
```

Lastly, using Neighbors Regressor, below output was

7.1 Stacking for Regression

Stacked Generalization or “Stacking” for short is an ensemble machine learning algorithm. It involves combining the predictions from multiple machine learning models on the same dataset, like bagging and boosting. The architecture of a stacking model involves two or more base models, often referred to as level-0 models, and a meta-model that combines the predictions of the base models, referred to as a level-1 model.

Level-0 Models (Base-Models): Models fit on the training data and whose predictions are compiled.

Level-1 Model (Meta-Model): Model that learns how to best combine the predictions of the base models.

A box-and-whisker plot is then created comparing the distribution negative MAE scores for each model.

KNRegressor, Decision Tree Regressor and Support Vector Regressor were used for Level-0 that is as base models for stacking. After implementing Cross validation and model fit, the below figure represents the output:

8. Conclusion

From the tasks I did in this project, I got the knowledge about best city & city type, best months and the relationship between the numerical attributes and target variable. Gradient Boosting seems to work best for the dataset as it has good score and pretty reasonable root mean square error value. When done stacking, the negative MSE of stacking and the MSE for support vector regressor doesn't show much difference, hence both can fit well for the dataset.

For level-1, `get_stacking` function defines the `StackingRegressor` model by first defining a list of tuples for the three base models, then defining the Linear regression meta-model to combine the predictions from the base models using 5-fold cross-validation. Implementing this, we get the below outputs:

```
>knn -1981639.171 (466618.059)
>cart -2300354.596 (752088.098)
>svm -1611450.553 (481780.547)
>stacking -1719066.360 (452842.453)
```

