

# Restaurant revenue prediction using deep learning.

Mokhutli Letsae



# Abstract

Predictive analytics is a powerful tool that enables businesses to forecast metrics such as future revenue and profitability by analyzing historical data and market trends. This capability allows businesses to anticipate upcoming financial performance and proactively address potential profit shortfalls. For instance, they can launch targeted marketing campaigns to boost revenue and mitigate losses. Ultimately, predictive analytics helps businesses optimize their operations and make informed decisions to maximize profitability.

This write-up documents a data science project aimed at predicting monthly revenue for restaurants using the CRISP-DM (Cross-Industry Standard Process for Data Mining) framework. The project utilized a deep learning model, specifically a multi-layer perceptron neural network, to build the predictor. The training data for the neural network model was sourced from Kaggle's "Restaurants Revenue Prediction" dataset. The project demonstrates how advanced machine learning techniques can be effectively applied to real-world business problems, providing actionable insights that enhance decision-making and operational efficiency.

By following the CRISP-DM methodology, the project involved several key phases:

**Business Understanding:** Defining the project objectives and requirements from a business perspective.

**Data Understanding:** Collecting and exploring the dataset to gain insights and identify patterns.

**Data Preparation:** Cleaning and transforming the data to make it suitable for modeling.

**Modeling:** Building and training the multi-layer perceptron neural network using the prepared data.

**Evaluation:** Assessing the model's performance to ensure it meets the business objectives.

**Deployment:** Implementing the model in a real-world setting to start generating predictions.

The results of this project highlight the effectiveness of deep learning models in predicting restaurant revenue, demonstrating significant potential for improving financial forecasting and strategic planning in the food-selling industry.

# Introduction

This project implemented in python 3.8 on Jupyter notebooks using CRISP-DM aims to predict monthly restaurant revenue using a variety of predictors such as review count, marketing spend, and average menu price, with the goal of helping restaurant owners anticipate and mitigate potential losses through targeted campaigns hence maximizing profit. Utilizing a dataset from Kaggle, a Multi-Layer Perceptron (MLP) neural network was implemented due to its robustness and superior accuracy compared to other models like ensemble methods, XGBoost, and Random Forests. Through careful exploratory data analysis (EDA), data preprocessing, model training, and evaluation, the MLP demonstrated its efficacy in forecasting revenue, making it the preferred choice for this regression task. The project's findings underscore the potential of data science to enable proactive business strategies in the restaurant industry.

## 1.Methodology

The effort to predict monthly revenue followed a structured approach known as the CRISP-DM methodology as mentioned in the introduction. It all began with a clear understanding of the business goal, which was to accurately forecast monthly revenue, crucial for making informed decisions. The first step involved gathering historical revenue data and relevant features that could impact revenue generation. Then, we carefully examined this data to understand its quality, structure, and any underlying patterns.

Next up was getting the data ready for analysis. This meant cleaning it up and transforming it so that it could be used effectively by modeling techniques. We fixed any errors, dealt with missing information, and standardized the format of the data. With the data prepped, it was time to try out different regression models to predict monthly revenue. We tested three models: Random Forest, XGBoost, and MLP Regressors.

When it came to evaluating these models, we looked at metrics like Mean Absolute Error (MAE) and  $R^2$  value to see how well they performed. What we found was that Random Forest and XGBoost models seemed to be overfitting, meaning they performed much better on the training data than on new, unseen data. However, the MLP Regressor showed consistent performance on both the training and test datasets, making it the most reliable choice. The table below shows the metrics mentioned above.

Model	Train MSE	Test MSE	Train RMSE	Test RMSE	Train MAE	Test MAE	Train R2	Test R2
MLP	3214.18	3465.08	56.69	58.86	45.33	47.33	0.70	0.63
RandomForest	530.60	3983.61	23.03	63.12	18.47	50.71	0.95	0.57
XGBoost	16.26	4497.92	4.03	67.07	2.75	54.64	1.00	0.52

## 1. Dataset Overview

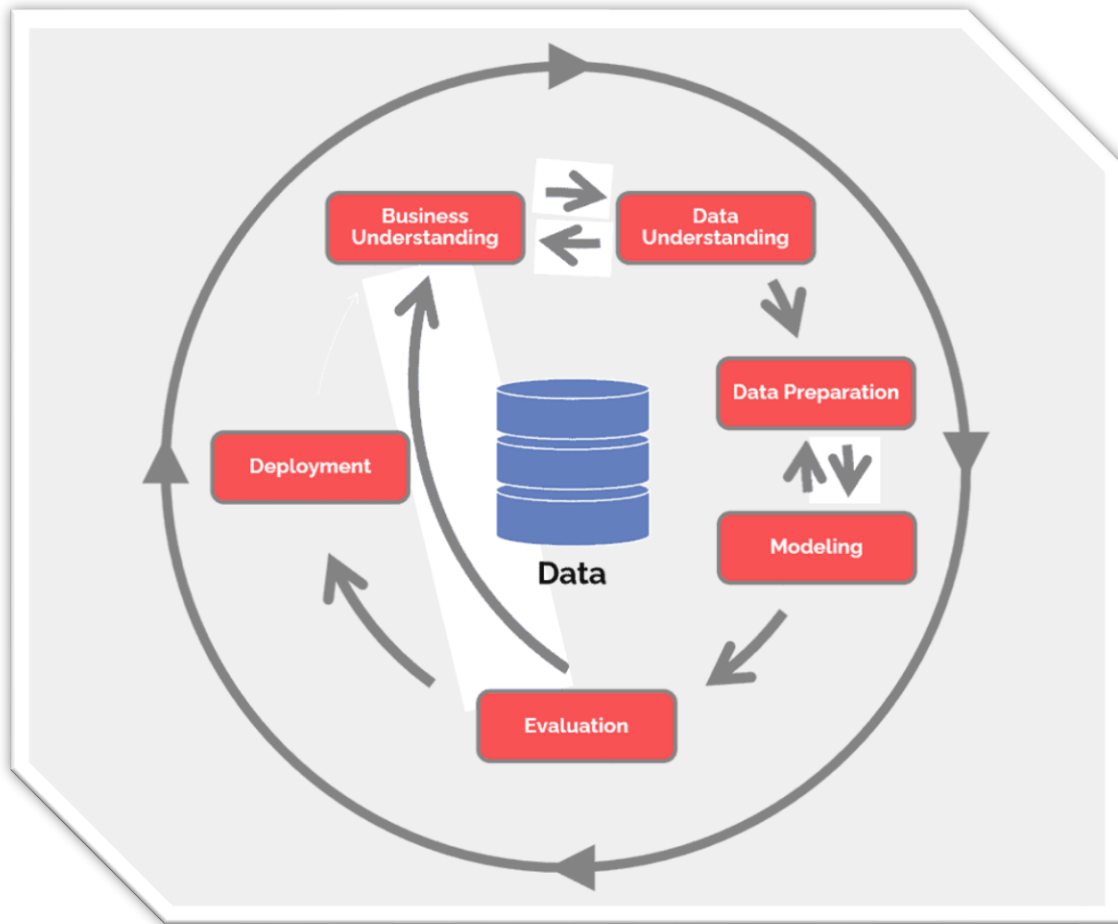
The dataset with usability of 10 was downloaded in a **csv** format from Kaggle (<https://www.kaggle.com/datasets/mrsimple07/restaurants-revenue-prediction> ).

The file had 1,000 rows and 8 columns and no missing values. The 8 columns' roles, data types, value count, and definitions are listed on the below table.

Attribute	ML-model data role	data type	value count	Attribute description
<b>Number_of_Customers</b>	Predictor	int	1000	The count of customers visiting the restaurant
<b>Menu_Price</b>	Predictor	float	1000	Average menu prices at the restaurant
<b>Marketing_Spend</b>	Predictor	float	1000	Expenditure on marketing activities
<b>Cuisine_Type</b>	Predictor	string	1000	The type of cuisine offered (Italian, Mexican, Japanese, American).
<b>Average_Customer_Spending</b>	Predictor	float	1000	Average spending per customer
<b>Promotions</b>	Predictor	float	1000	Binary indicator (0 or 1) denoting whether promotions were conducted
<b>Reviews</b>	Predictor	float	1000	Number of reviews received by the restaurant
<b>Monthly_Revenue</b>	Target	float	1000	Simulated monthly revenue, the target variable for prediction

## 2. CRISP-DM

**Crisp-DM** framework was used as a methodology for this project. The framework acted as a guideline for optimal attempt of the project.



**Figure 1**

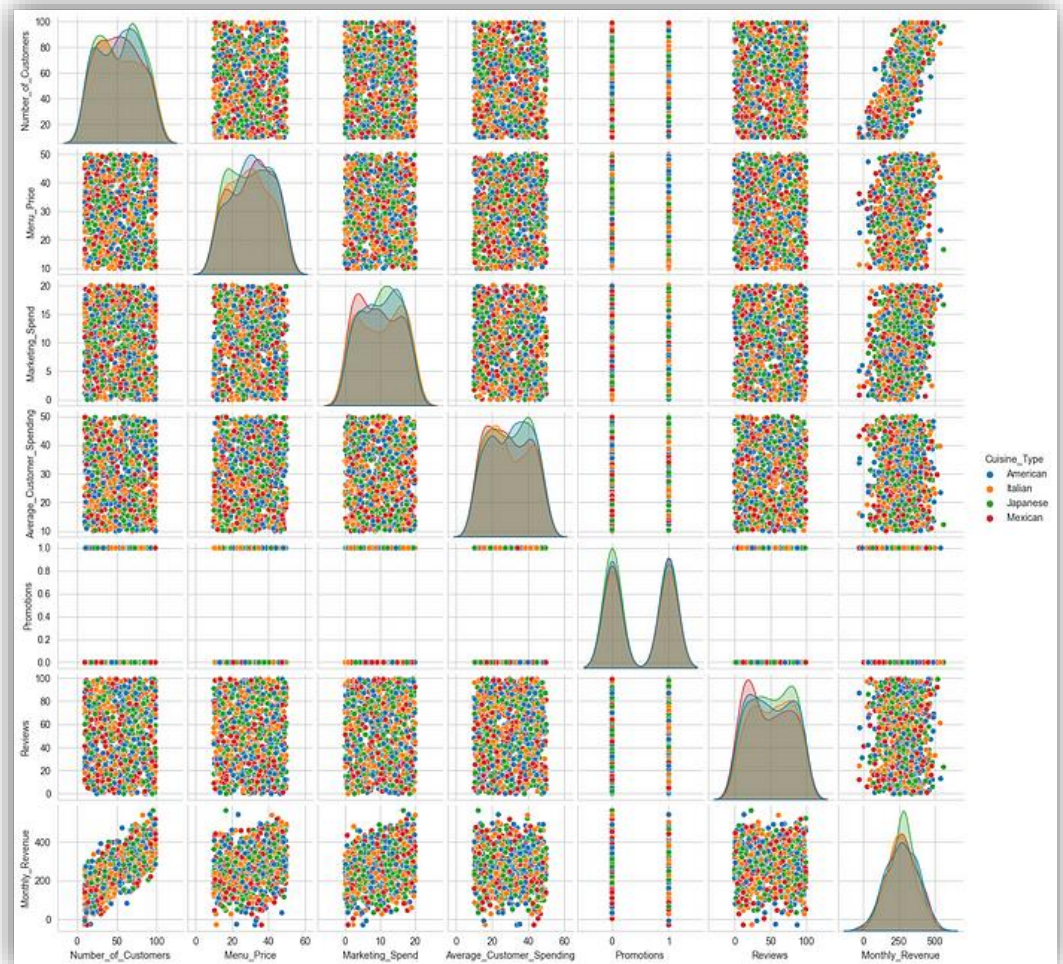
**Figure 1** shows the stages of the framework, from the figure you can see that the framework follows a cyclic flow implying agility. The project was built in an agile manner.

- a. **Business understanding:** The business problem in this project is predicting monthly revenue for a restaurant. The value of forecasting revenue is that the restaurant owner will be able to avoid losses and maximize profits by putting in place targeted campaigns in seeing that the revenue forecast will produce a loss.
- b. **Data understanding (EDA):** While establishing facts and insight, below is what was extracted from the data.

1. From the table below, we observe that on average, 53 customers visit the restaurant per cuisine. Our dataset includes various types of cuisines, and we have expressed all statistical metrics per cuisine. This means that the metrics in this table represent averages across all cuisines.

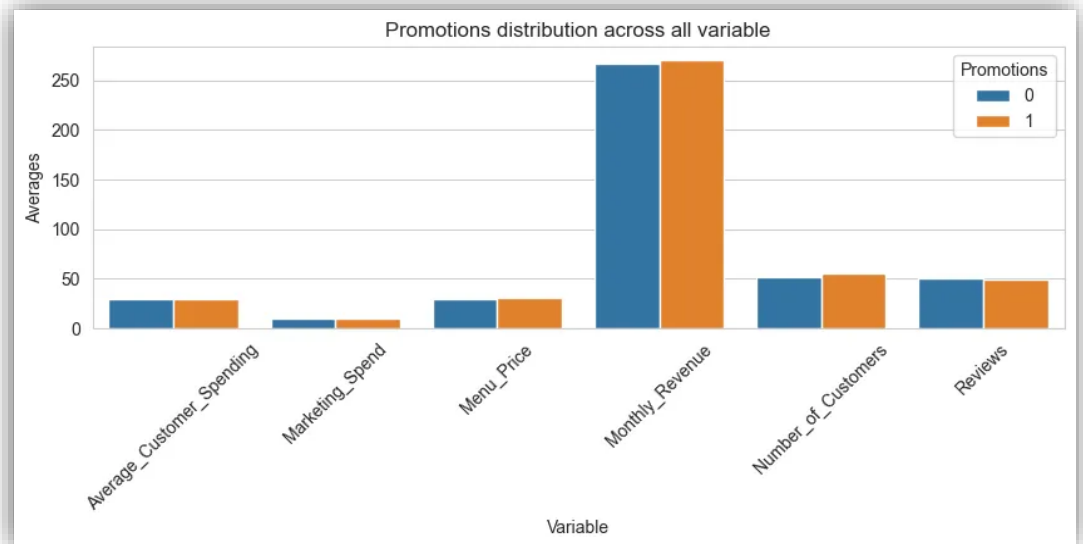
	Number_of_Customers	Menu_Price	Marketing_Spend	Average_Customer_Spending	Promotions	Reviews	Monthly_Revenue
count	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000
mean	53.271000	30.219120	9.958726	29.477085	0.497000	49.837000	268.724172
std	26.364914	11.278760	5.845586	11.471686	0.500241	29.226334	103.982950
min	10.000000	10.009501	0.003768	10.037177	0.000000	0.000000	-28.977809
25%	30.000000	20.396828	4.690724	19.603041	0.000000	24.000000	197.103642
50%	54.000000	30.860614	10.092047	29.251365	0.000000	50.000000	270.213964
75%	74.000000	39.843868	14.992436	39.553220	1.000000	76.000000	343.395793
max	99.000000	49.974140	19.994276	49.900725	1.000000	99.000000	563.381332

2. Another important aspect of understanding our data would be checking the relationships of our variables and their distributions, to achieve this we can use a pair-plot.

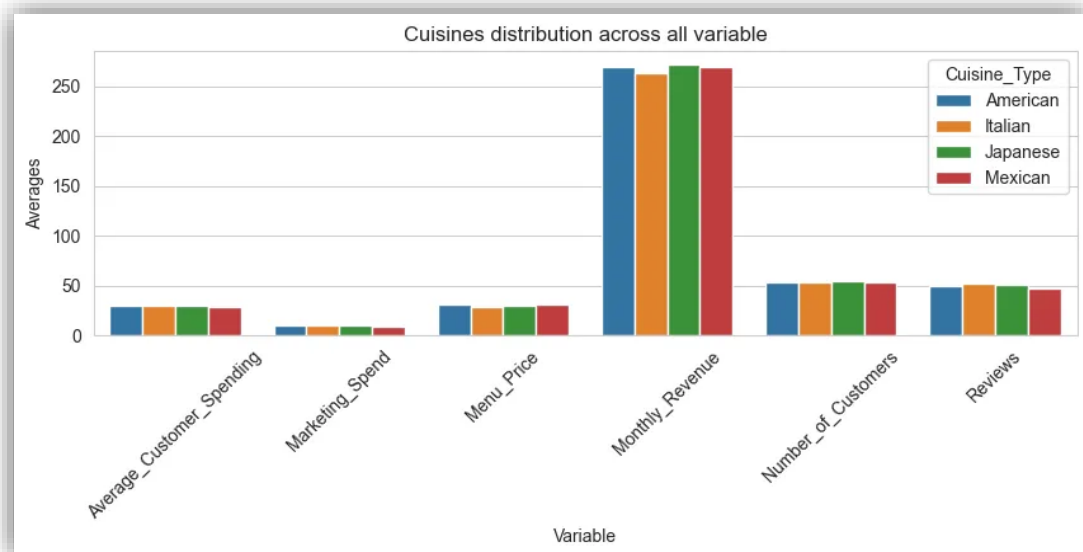


***What we see from the above pair-plot is that no relationships exhibit linear relationships except for number of customers and monthly revenue of which is also not totally linear but forms a linear-like. The relationships are cluster relationships.***

3. With the below histogram, we check if promotions have any impact on monthly revenue. We can extract that promotions seem to have a very little impact on monthly revenue as seen from the histogram, there's a very little difference between revenue made when items are promoted and not. To make promotion have impact we can increase marketing spend on promotions as marketing is the heart of promotions.



4. Lastly, there seems to be very small differences in pricings of cuisines hence the near the same monthly revenues. Japanese and Italian cuisines should have a higher pricing compared to the other 2 as they known to be high-end and expensive.





- c. **Data Preparation:** For this stage, basic processes for data processing and preparation and cleaning are performed:

**Handling missing values** - There were no missing values from the dataset hence there was no imputations or dropping of values performed.

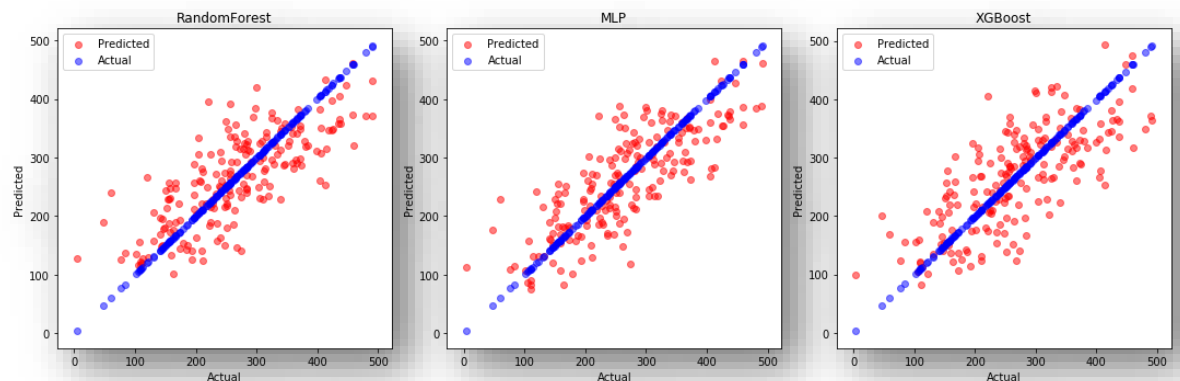
**Addressing outliers** - The data did not have may outliers, but IQR was performed on monthly revenue to remove outliers, and there were rows with negative revenue which were dropped.

**Feature engineering** – One-hot encoding was used on nominal variables in which this case we had one variable being cuisine type.

**Feature selection** - Variance and correlation threshold techniques were used to select features **Number\_of\_Customers**, **Menu\_Price**, **Marketing\_Spend**.

- d. **Model Building:** Two models were developed(scikit-learn) for comparison: a Deep Learning model (multi-layer perceptron) and XGBoost regressor. The multi-layer perceptron regressor outperformed the XGBoost regressor and it was used for predictions.
- e. **Evaluation:** Evaluation of the model's performance metrics.

Looking at the below graphs, we realize that some of the predictions deviate away from the actual values on the below scatter plots..



On the below table we can see that the MLP regressor has the lowest average error (MAE) and the highest accuracy(R2) hence it was used as our model.

Model	Train MSE	Test MSE	Train RMSE	Test RMSE	Train MAE	Test MAE	Train R2	Test R2
MLP	3214.18	3465.08	56.69	58.86	45.33	47.33	0.70	0.63
RandomForest	530.60	3983.61	23.03	63.12	18.47	50.71	0.95	0.57
XGBoost	16.26	4497.92	4.03	67.07	2.75	54.64	1.00	0.52

MLP Regressors, or Multi-Layer Perceptrons, are highly flexible neural networks capable of modeling intricate non-linear relationships due to their deep and wide architecture. This flexibility allows them to approximate any continuous function accurately, given sufficient data and appropriate network design. In contrast, XGBoost and Random Forest Regressors, which are tree-based ensemble methods, may struggle with very complex non-linear relationships, especially if the trees are shallow or if intricate interaction terms are not explicitly modeled.

In addition, Neural networks, like MLP Regressors, implicitly perform feature engineering through their hidden layers, automatically learning useful data representations and optimizing for the final prediction goal in an end-to-end manner. This ability is particularly advantageous when dealing with raw data. On the other hand, XGBoost and Random Forest Regressors typically require more explicit feature engineering. Although they can handle a variety of input formats and can benefit from well-engineered features, they do not inherently learn new feature representations and may need manual preprocessing to achieve optimal performance. Hence MLP regressor yielded better performance than XGBoost and Random Forest regressors.

Hyper-parameter tuning was performed, and the tuned models were over-fitting due to the parameter searched through making the model complex. Grid-search approach was used on all the 3 models. Below are the results of tuned modes, we see the overfitting on the significant metrics differences in training and test. None of the tuned model was used but the non-tuned MLP model was used. It seemed that the non-tuned models were not complex.

Model	Best_Params	Train_MSE	Test_MSE	Train_RMSE	Test_RMSE	Train_MAE	Test_MAE	Train_R2	Test_R2
MLP	{'activation': 'relu', 'hidden_layer_sizes': (...	3059.99	3668.18	55.32	60.57	44.31	48.82	0.71	0.61
XGBoost	{'learning_rate': 0.1, 'max_depth': 3, 'n_esti...	2398.39	3715.90	48.97	60.96	39.09	49.75	0.77	0.60
RandomForest	{'max_depth': 10, 'n_estimators': 200}	718.90	3937.02	26.81	62.75	21.98	50.36	0.93	0.58

- f. **Deployment:** The model artifacts can be deployed on on-prem or cloud platforms, and it can be utilized using a web pass, desktop app or just an environment it is used in.

## 2.Conclusion

To have better performance of our model, augmentation of training data and more features is necessary. The models' results are acceptable but can be improved. The project was undertaken as a capstone project for completing a nanodegree from Udacity. It was undertaken by Mokhutli Letsae( <https://www.linkedin.com/in/mokhutliletsae/> ), has the code repo on GitHub( [https://github.com/letsaemokhutli/capstone\\_project.git](https://github.com/letsaemokhutli/capstone_project.git) ). The MLP neural network model demonstrates good performance with an acceptable level of accuracy. The metrics suggest that the model does not overfit, as evidenced by the minimal difference between training and test metrics. While the error percentage indicates there is room for improvement, the model's predictions are reliable enough for practical use in decision-making processes for the restaurant.

## 3.Future Work

Augmenting data (features and row count) will be the next step to improve the model performance metrics. Using a cloud environment is also part of future intentions.

## 4.Acknowledgements

I would like to acknowledge my boss, Itumeleng Senekane (Executive Head: CVM and Big data), Vodacom Lesotho for giving me a chance to take this Udacity course. Chatgpt has been a very useful resource in the completion of the project hence it's owners I would like to acknowledge. I would also like to finally Acknowledge the whole Vodacom group for giving us a chance into a more than a glimpse of data science.