

Data Analysis and Modeling Report.....	1
1. Introduction.....	1
2. Data Exploration.....	1
3. Data Preprocessing .....	2
4. Modeling.....	3
5. Model Evaluation .....	3
6. Insights.....	5
7. Conclusion .....	5

# Data Analysis and Modeling Report

## 1. Introduction

This report shows the result of a thorough data analysis and modeling exercise done on a dataset of 1000 rows and 17 columns which includes features such as identification, timestamps, numerical, and categorical variables. The main objective of this analysis was to develop two predictive models to understand the relationships between the input features and the target variable, and to predict future results with reasonable accuracy.

The analysis was carried out in this order data exploration, preprocessing, and the development of multiple regression models, including Support Vector Regression (SVR) and XGB Regressor. The models were validated using standard evaluation metrics like R-squared, Mean Absolute Error (MAE), and Mean Squared Error (MSE).

## 2. Data Exploration

I started the analysis by exploring the data in the dataset to understand the structure of the dataset, the variable distributions, and the relationships between different features. The dataset contains both numerical and categorical features. Notable observations from this step include:

**Distribution of Numerical Features:** I made use of histograms to visualize the distribution of numerical features. Some features exhibited properties indicating imbalance, this confirms the need for transformation during preprocessing.

**Correlation Analysis:** A correlation matrix was generated to identify the relationships between numerical features. This analysis displayed strong correlations between certain pairs of features, which could influence the modeling process.

**Categorical Features:** The frequency distribution of categorical variables was analyzed to find out their impact on the target variable the target column in my dataset it's the column named Rating. Some categories were dominant, while others were not.

### 3. Data Preprocessing

Preprocessing of data was necessary to make sure that the data quality was good enough before making use of it in training the predictive models. These steps in preprocessing can be summarized as follows:

**Dealing with Missing Values:** Missing values were detected and addressed using imputation techniques. For numerical variables, missing values were replaced with the mean or median, while for categorical variables, the mode was used.

**Feature Scaling:** Since SVR and other models used are sensitive to the scale of the data, feature scaling was applied to standardize the numerical variables.

**Encoding Categorical Variables:** Categorical variables were encoded using techniques like one-hot encoding, which converts categories into binary features that can be used in regression models.

## 4. Modeling

Two primary models were developed and tested in this analysis:

**Support Vector Regression (SVR):** SVR is a regression model that attempts to find a hyperplane in a high-dimensional space that best fits the data. The model was trained using the training data, and hyperparameters were tuned to optimize performance.

**Training and Evaluation:** The SVR model was trained on a subset of the data, and predictions were made on the test set. The model's performance was evaluated using R-squared, MAE, and MSE.

**XGB Regressor:** XGB Regressor is an advanced ensemble method that uses gradient boosting for regression tasks. It was selected for its ability to handle large datasets and capture complex patterns in the data.

**Training and Evaluation:** Similar to SVR, the XGB Regressor model was trained on the data, and its performance was evaluated using the same metrics.

## 5. Model Evaluation

The performance of the models was compared using several evaluation metrics:

**R-Squared:** R-squared measures the proportion of variance in the dependent variable that is predictable from the independent variables. A higher R-squared indicates better model performance.

**SVR Performance:** The SVR model achieved an R-squared value of approximately -0.011, indicating that the model was not able to explain much of the variance in the data.

**XGB Regressor Performance:** The XGB Regressor model, however, performed better with an R-squared value that indicated a stronger fit to the data.

**Mean Absolute Error (MAE):** MAE measures the average magnitude of errors in a set of predictions, without considering their direction. It is a linear score, meaning all individual differences are weighted equally.

**SVR Performance:** The MAE for SVR was around 1.52, indicating a moderate level of prediction error.

**XGB Regressor Performance:** XGB's MAE was lower than that of SVR, suggesting it made more accurate predictions on average.

**Mean Squared Error (MSE):** MSE measures the average of the squares of the errors—that is, the average squared difference between the estimated values and the actual value.

**SVR Performance:** The MSE for SVR was approximately 3.10, indicating higher variance in the prediction errors.

**XGB Regressor Performance:** The MSE for XGB Regressor was significantly lower, reinforcing its superiority over SVR in this analysis.

## 6. Insights

Several unique insights were derived from the analysis:

**Model Selection:** The comparison between SVR and XGB Regressor highlighted the importance of model selection. While SVR is useful in certain scenarios, XGB Regressor demonstrated a much stronger ability to model the data accurately in this case.

**Feature Importance:** The XGB Regressor model also provided insights into feature importance, revealing which variables were most influential in predicting the target variable. This information is crucial for feature engineering and further model refinement.

**Data Quality:** The analysis underscored the impact of data quality on model performance. Even small amounts of missing data or unscaled features could drastically reduce model accuracy.

**Correlation Effects:** High correlation between certain features suggested potential multicollinearity, which could be addressed in future analyses to improve model robustness.

## 7. Conclusion

The analysis and modeling exercise provided valuable insights into the dataset and demonstrated the power of advanced regression techniques like XGB Regressor. The models developed offer a solid foundation for predicting the target variable, though there is room for further refinement. Future work could involve exploring additional models, tuning hyperparameters more extensively, and addressing potential issues like multicollinearity. The findings from this analysis can inform decision-making and guide future data-driven strategies.

This report presents a comprehensive overview of the data analysis process, from exploration to model evaluation, and highlights key insights that were derived. The models developed, particularly XGB Regressor, show promise for accurate prediction and offer significant potential for real-world application.