

Hyper Supra 1.0

Supra Without Brakes

Automated EDA and Machine Learning Framework

Developer:

TC.Antony - Data Scientist.

- Manual Book

Purpose:

- Automated Exploratory Data Analysis (EDA) with business insights.
- Machine learning capabilities for classification, regression, and clustering.

Features:

- Automatic Summary Statistics: Provides quick insights into dataset structure.
- Visualization Tools: Includes histograms, scatterplots, and boxplots for data exploration.
- PCA for Dimensionality Reduction: Reduces feature space for better analysis.
- Outlier Detection: Identifies anomalies in the dataset.
- Correlation Insights: Highlights relationships between variables.
- Supervised Learning: Supports classification and regression tasks.
- Unsupervised Learning: Includes clustering algorithms like KMeans, DBSCAN, and Agglomerative Clustering.

- Cross-Validation and Model Evaluation: Ensures robust model performance.
- Feature Engineering and Preprocessing: Handles missing values, encoding, and scaling.
- Business Insights: Provides actionable conclusions for decision-making.
- Downloadable Reports and Datasets: Allows users to export results and cleaned datasets.

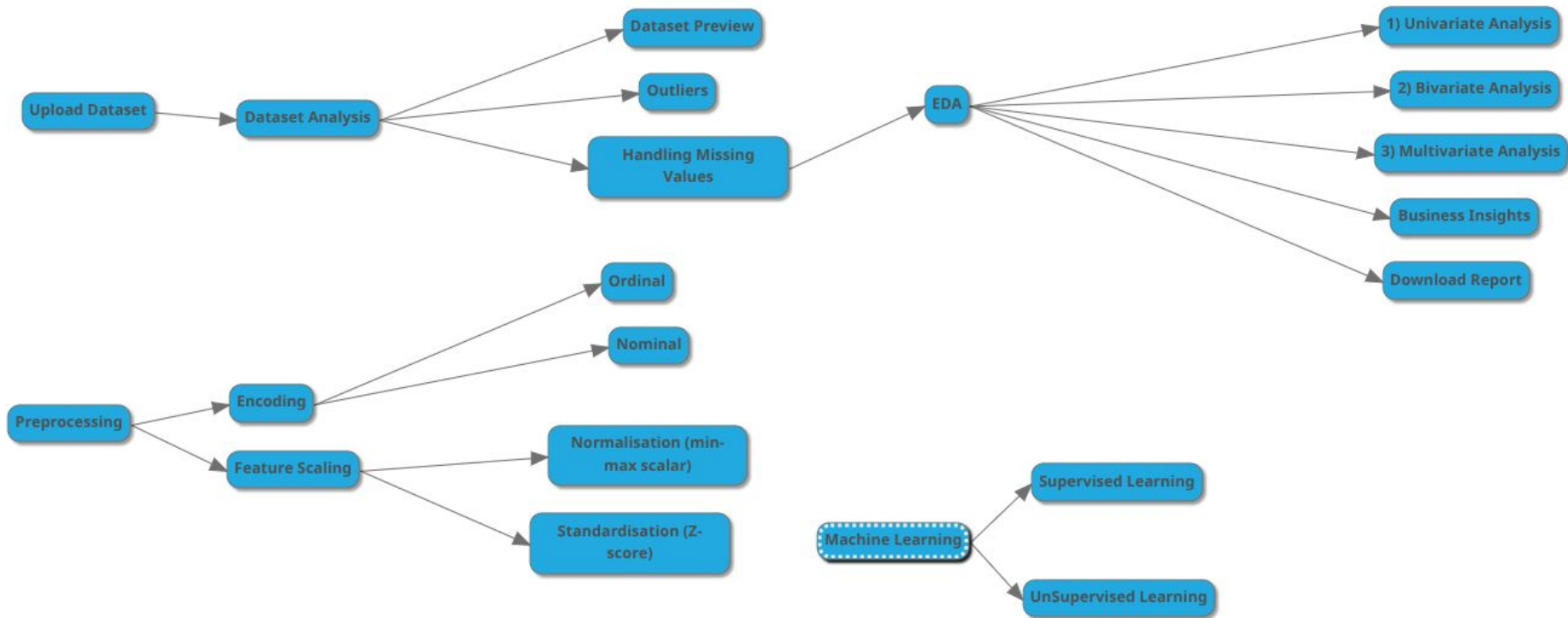
Instructions:

- Upload Dataset: Upload datasets in CSV, Excel, or JSON format.
- Dataset Analysis: Explore summary statistics, outliers, and missing values.
- EDA Process: Perform univariate, bivariate, and multivariate analysis.
- Preprocess: Handle missing values, encode categorical variables, and scale features.
- Machine Learning: Train and evaluate supervised and unsupervised models

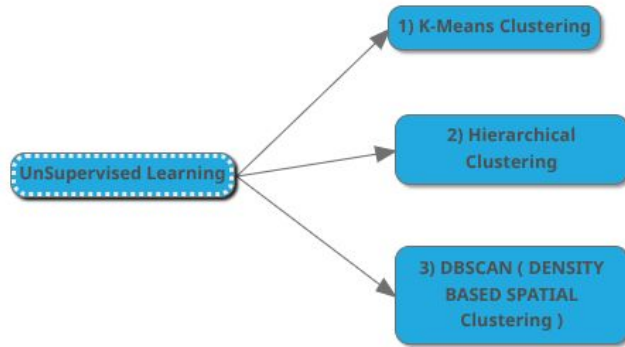
Dependencies:

Streamlit, Pandas, NumPy, Scikit-learn, Seaborn, Matplotlib, XGBoost, LightGBM, Pandas Profiling.

Mind Map:







Step 1: Dataset Upload (CSV, EXCEL,JSON,PICKLE)

Automated EDA and Machine Learning Framework

About Upload Dataset Dataset Analysis EDA Process Preprocess Machine Learning

Upload Your Dataset

Upload CSV, Excel, JSON, or Pickle files



Drag and drop files here

Limit 200MB per file • CSV, XLSX, JSON, PKL

Browse files



taxi_trip_pricing.csv 67.3KB



Step: 2 - Data Preview & Summary

Data Preview

	Trip_Distance_km	Time_of_Day	Day_of_Week	Passenger_Count	Traffic_Conditions	Weather	Base_Fare
0	19.35	Morning	Weekday	3	Low	Clear	
1	47.59	Afternoon	Weekday	1	High	Clear	
2	36.87	Evening	Weekend	1	High	Clear	
3	30.33	Evening	Weekday	4	Low	None	
4	None	Evening	Weekday	3	High	Clear	

Data Summary

	Trip_Distance_km	Passenger_Count	Base_Fare	Per_Km_Rate	Per_Minute_Rate	Trip_Duration_Minute
count	950	950	950	950	950	
mean	27.0705	2.4768	3.503	1.2333	0.2929	62.1
std	19.9053	1.1022	0.8702	0.4298	0.1156	32.1
min	1.23	1	2.01	0.5	0.1	5
25%	12.6325	1.25	2.73	0.86	0.19	35.8
50%	25.83	2	3.52	1.22	0.29	62
75%	38.405	3	4.26	1.61	0.39	89.1

Step 3: Dataset Analysis

Automated EDA and Machine Learning Framework

[About](#) [Upload Dataset](#) [Dataset Analysis](#) [EDA Process](#) [Preprocess](#) [Machine Learning](#)

Dataset Analysis for taxi_trip_pricing.csv

[Dataset Preview](#) [Outliers](#) [Standard Deviation](#) [Multicollinearity](#) [Missing Values](#) [Preprocessing Preview](#)

Dataset Preview

Total Rows: 1000

Total Columns: 11



Step 4: EDA (Exploratory Data Analysis)

Automated EDA and Machine Learning Framework

[About](#) [Upload Dataset](#) [Dataset Analysis](#) [EDA Process](#) [Preprocess](#) [Machine Learning](#)

Numeric Columns: Trip_Distance_km, Passenger_Count, Base_Fare, Per_Km_Rate, Per_Minute_Rate, Trip_Duration_Minutes, Trip_Price

Categorical Columns: Time_of_Day, Day_of_Week, Traffic_Conditions, Weather

[Univariate Analysis](#) [Bivariate Analysis](#) [Multivariate Analysis](#) [Business Insights](#) [Download Report](#)

1) Univariate Analysis

Select Column for Univariate Analysis

Trip_Distance_km





Step 5: Preprocessing

Learning Framework

About Upload Dataset Dataset Analysis EDA Process **Preprocess** Machine Learning

Preprocessing Tab

Encoding Recommendations Feature Scaling Recommendations

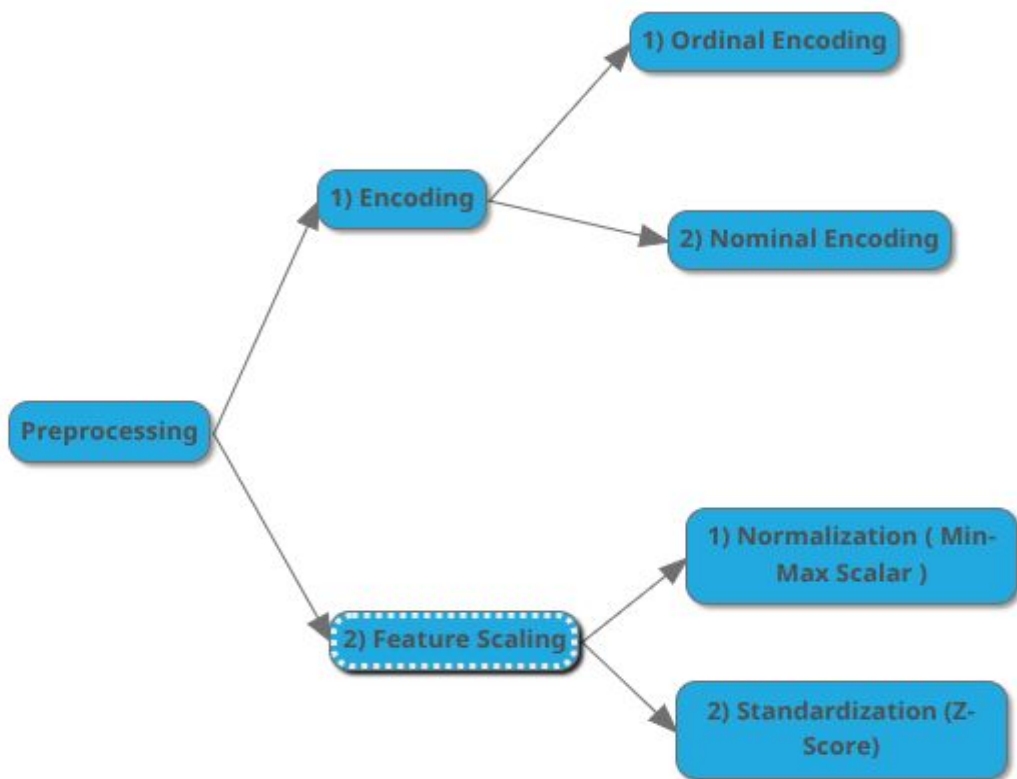
Encoding Recommendations

Columns recommended for Ordinal Encoding:

```
▼ [  
  0 : "Time_of_Day"  
  1 : "Day_of_Week"  
  2 : "Traffic_Conditions"  
  3 : "Weather"  
]
```

No columns recommended for nominal encoding.

Recommended columns for encoding:



Step 6: Supervised Learning (Regressor)

Automated EDA and Machine Learning Framework

About Upload Dataset Dataset Analysis EDA Process Preprocess **Machine Learning**

Supervised Learning Score Possibility Download Train and Test Datasets with Results Future Prediction Unsupervised

Supervised Learning

Select the target column:

Trip_Price



Target Column Data Type: Continuous

Inferred Task Type: Regression

Select the problem type manually:

☐ Classification

☒ Regression

Step 7: Algorithm Selection

Time_of_Day × Day_of_Week × Traffic_Conditions × Base_Fare × Weather ×

Select Algorithm

Linear Regression

Linear Regression

Polynomial Regression

Ridge Regression

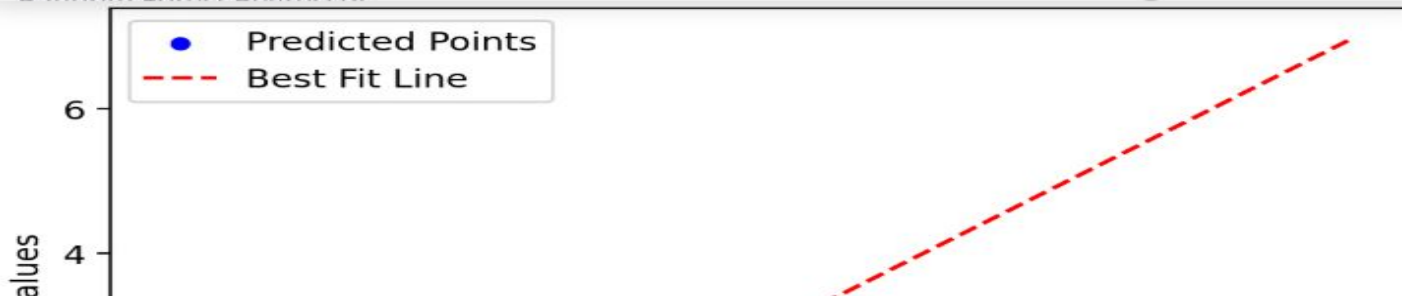
Lasso Regression

ElasticNet Regression

SVM Regression

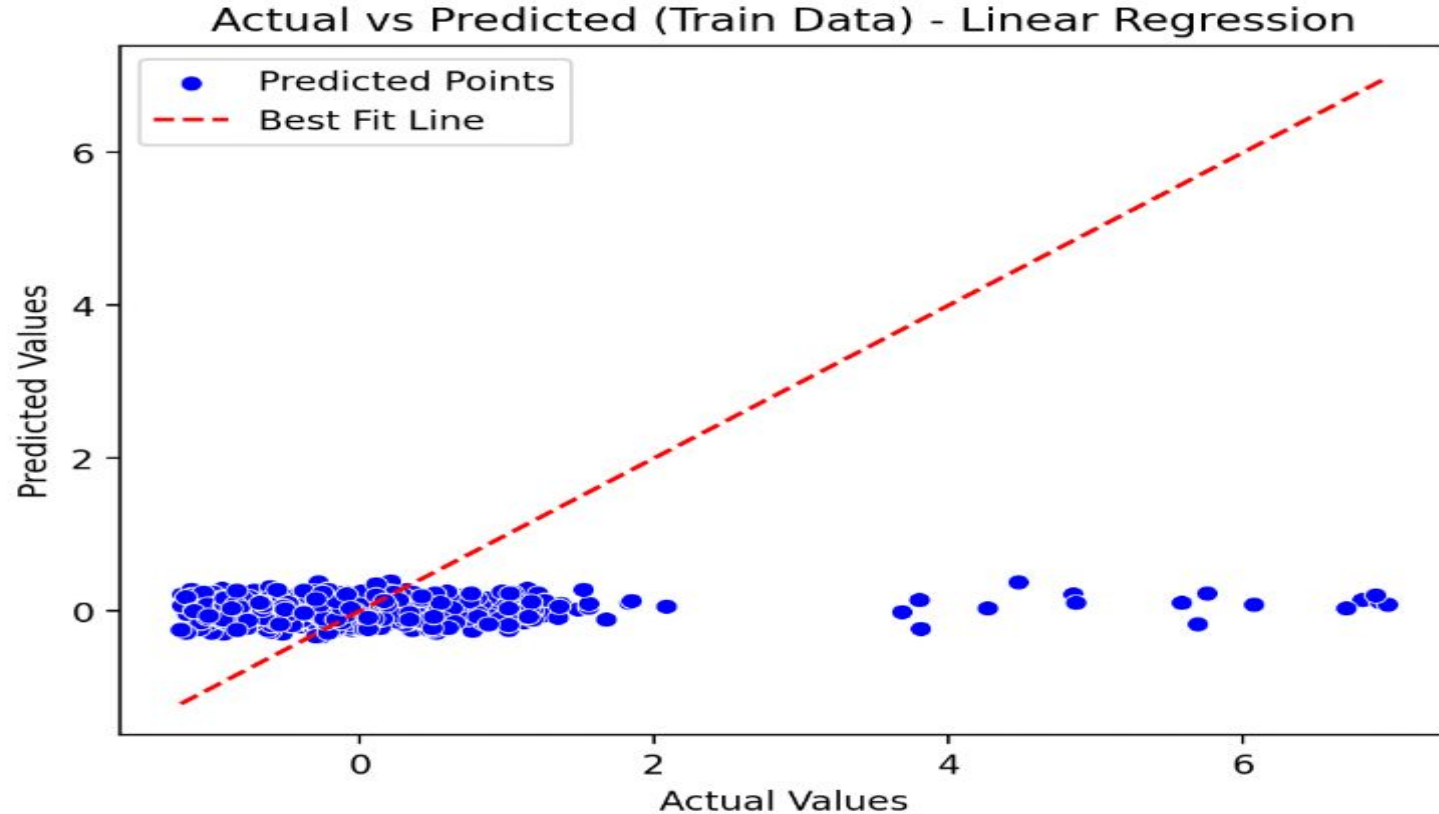
Decision Tree Regressor

Random Forest Regressor



PCA 2d Chart based on Algorithm (Target & Feature)

Data has more than 2 features. Reducing dimensions using PCA...



Polynomial , L1 & L2 (Hyperparameter Option)

Select Algorithm

Polynomial Regression



Select Degree for Polynomial Regression

2



2

10

Training Time: 0.01 seconds

Training Time: 0.02 seconds

Select Algorithm

Ridge Regression



Select Alpha for Ridge Regression

1.00



0.01

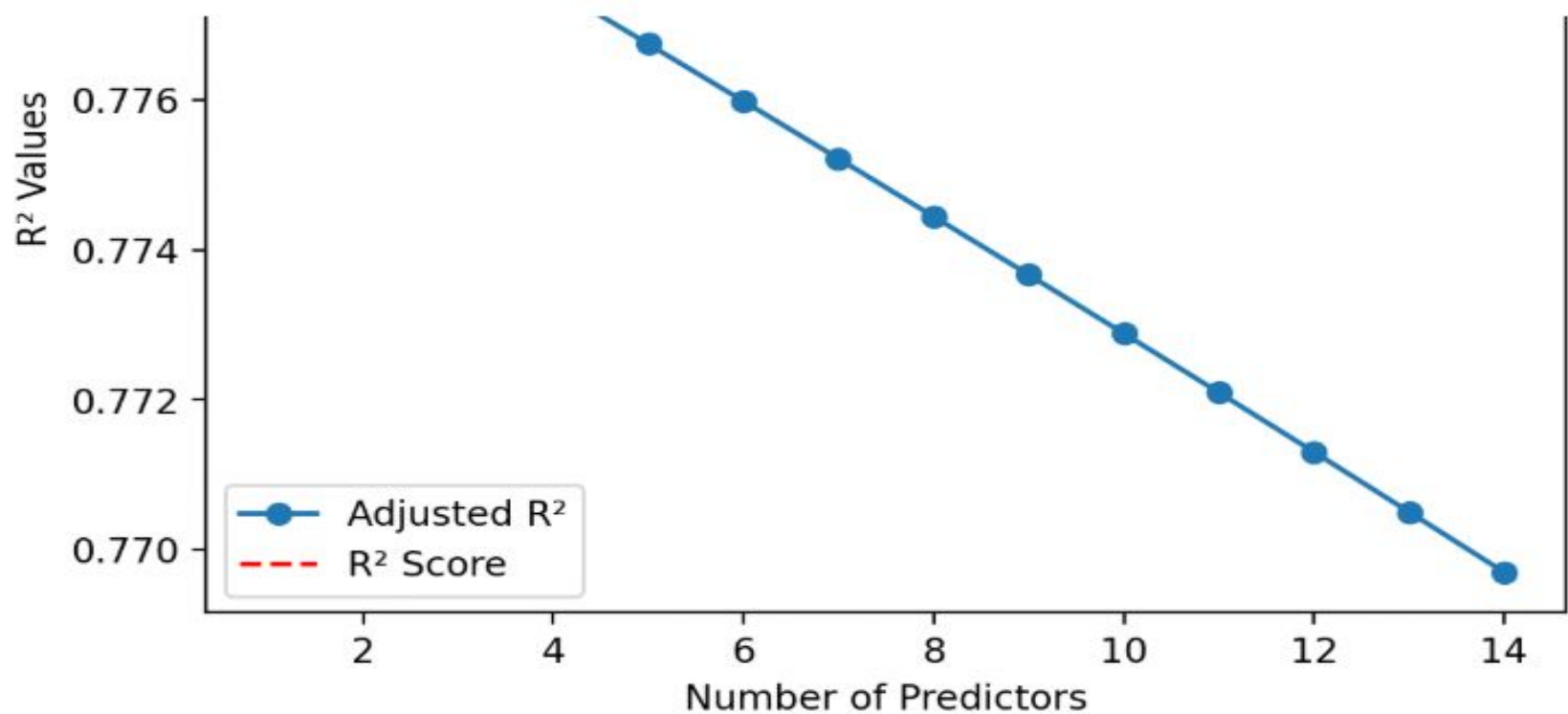
10.00

Training Time: 0.01 seconds

Training Time: 0.01 seconds

Step 8 : Scores for Regressor (Both Train & Test Date) with Chart

- 1) MSE (Mean Square Error)
- 2) MAE (Mean Absolute Error)
- 3) R²
- 4) Adjusted R²



Adjusted R^2 Score on Training Data: 0.8512

Adjusted R^2 Score on Test Data: 0.7729

Adjusted R² Score on Test Data: 0.7729

Training vs Test Accuracy (R² Score)

Training R² Score: 0.8534

Test R² Score: 0.7805

The model may be overfitting, as the training accuracy is significantly higher than the test accuracy.

Training vs Test Accuracy (Adjusted R² Score)

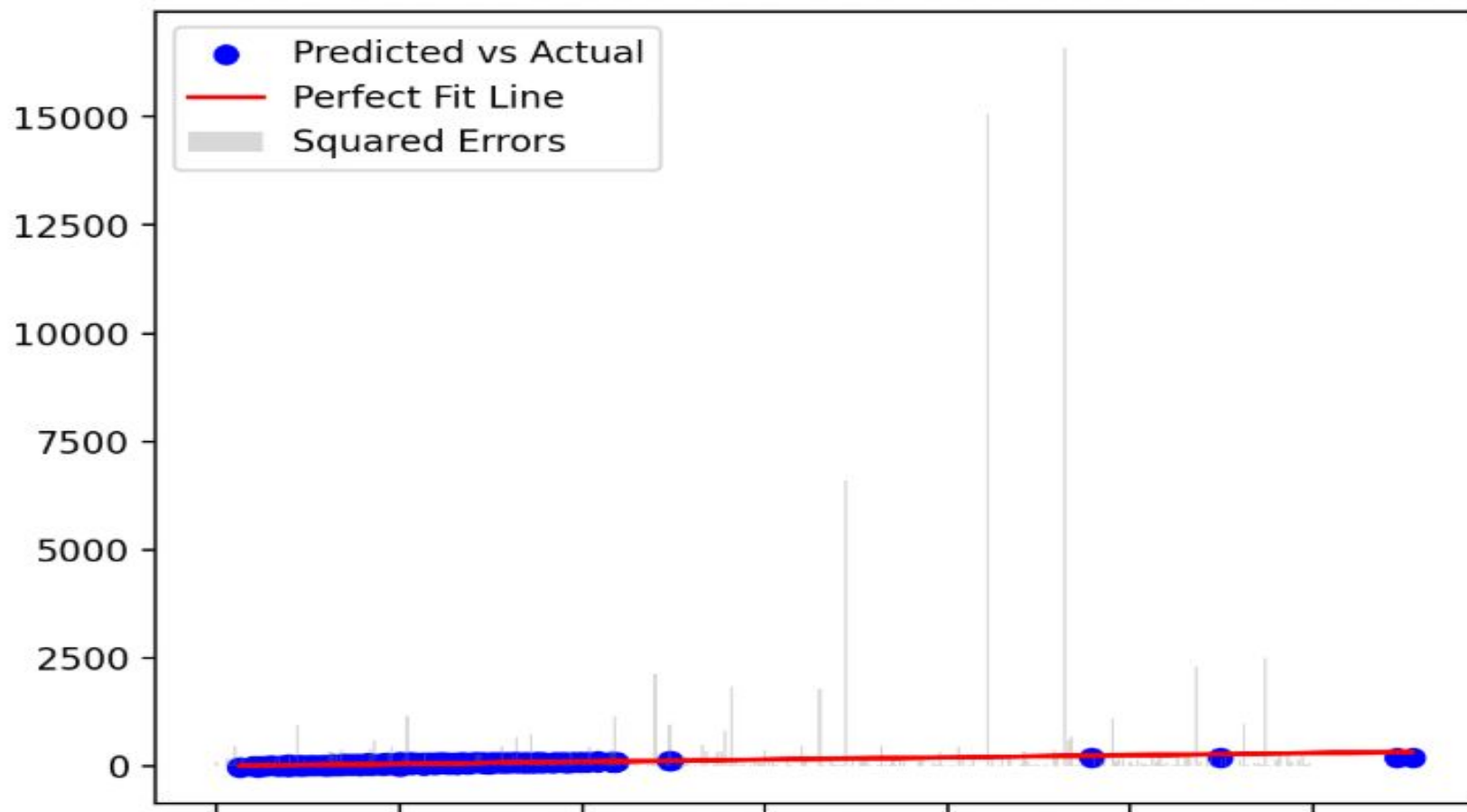
Training Adjusted R² Score: 0.8512

Test Adjusted R² Score: 0.7729

The model may be overfitting, as the training accuracy is significantly higher than the test accuracy.

Definition: Mean Squared Error (MSE) measures the average squared difference between actual and

MSE Visualization



Step 9: Cross Validation with hyperparameter option

Cross-Validation

Shape of dataset: (1000, 11)

Select Cross-Validation Method

K-Fold



Select the number of folds:



Cross-Validation Accuracy (K-Fold): 0.8219 ± 0.0306

Train Accuracy: 0.8361

Model Fit Status: Generalizing Well

Recommendation

Recommendation

- **K-Fold:** Suitable for general cases with balanced datasets.
- **Stratified K-Fold:** Best for imbalanced classification problems.
- **Holdout:** Quick but less reliable for small datasets.
- **Leave-One-Out:** Best for very small datasets but computationally expensive.
- **Leave-P-Out:** Rarely used due to computational cost.

Cross-Validation Best Practices ⇄

	Method	Best for Classifier	Best for Regressor
0	K-Fold	Yes	Yes
1	Stratified K-Fold	Yes	No
2	Holdout (80/20 Split)	No	Yes
3	Leave-One-Out	Yes	Yes
4	Leave-P-Out	No	No

Step 10: Supervised Learning (Classification)

About Upload Dataset Dataset Analysis EDA Process Preprocess **Machine Learning**

Supervised Learning Score Possibility Download Train and Test Datasets with Results Future Prediction Unsup

Supervised Learning

Select the target column:

Dataset



Target Column Data Type: Discrete

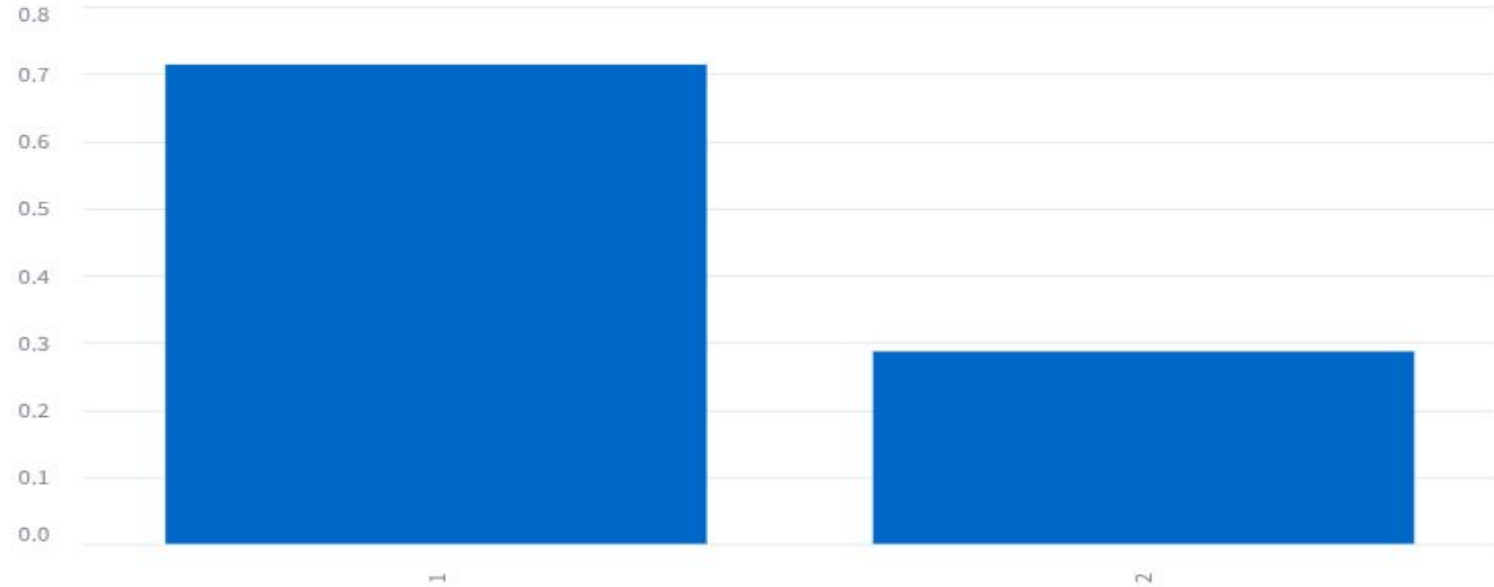
Inferred Task Type: Classification

Dataset Balance Status: Unbalanced

Class Distribution Status

Dataset Balance Status: Unbalanced

Class Distribution:



Suggested Binary Classification Algorithms:

- Logistic Regression
- SVM Classifier

Select feature columns

Age ×

Gender ×

Total_Bilirubin ×

Direct_Bilirubin ×



Select Algorithm

Logistic Regression



Training Time: 0.03 seconds

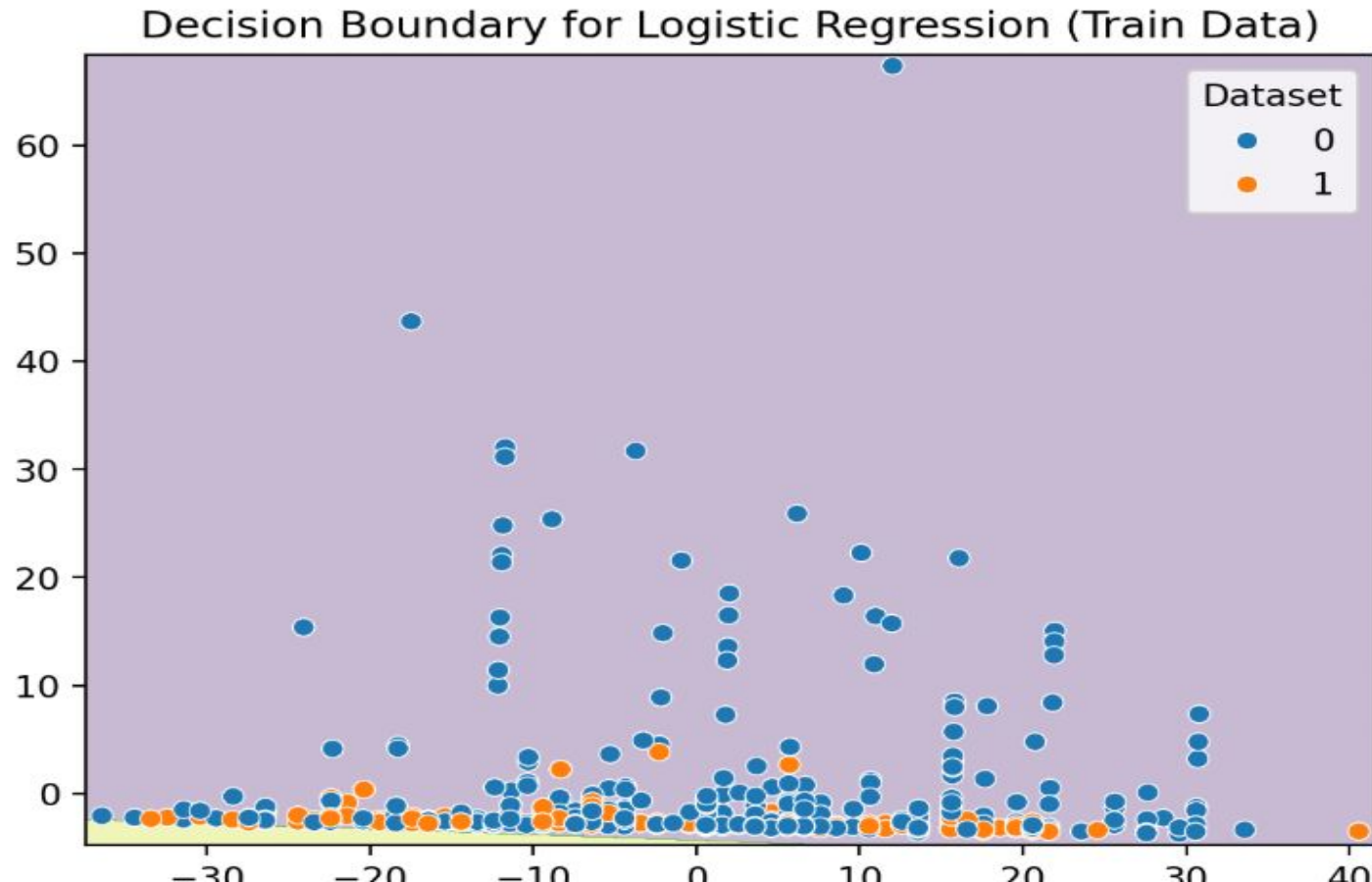
Training Time: 0.02 seconds

Select data for visualization:

☒ Train Data

☐ Test Data

PCA Chart for Classification



Score for Train & Test Data

Score

Train Metrics

Accuracy: 0.6985

Precision: 0.4524

Recall: 0.4972

F1 Score: 0.4190

Test Metrics

Accuracy: 0.7257

Precision: 0.3649

Recall: 0.4961

F1 Score: 0.4205

Model Fit Status

Status: Good Fit: Model generalizes well to new data.

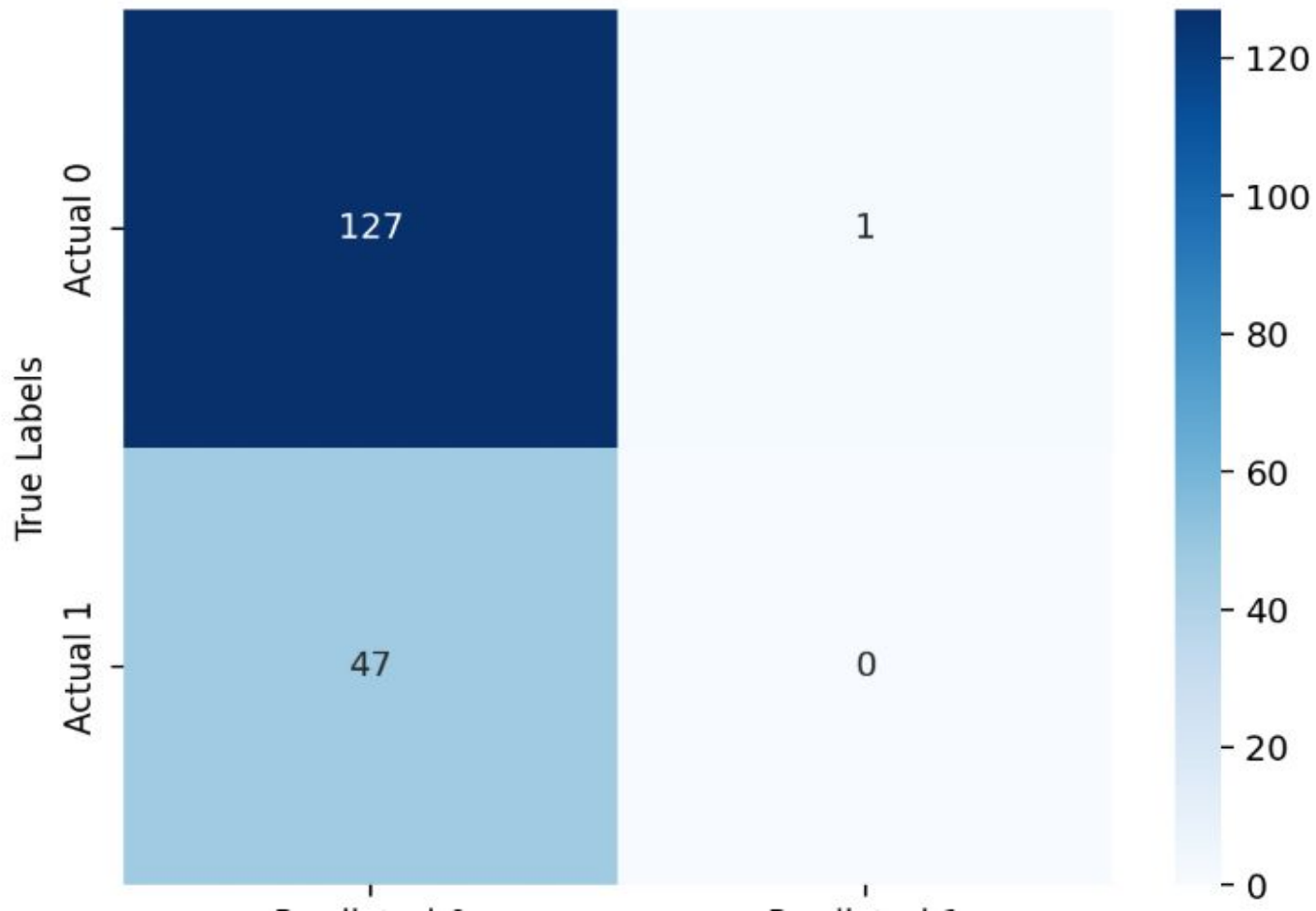
Precision = $TP / (TP + FP)$

Recall = $TP / (TP + FN)$

Classification Report (Test Data)

	precision	recall	f1-score	support
0	0.7299	0.9922	0.8411	128.0000
1	0.0000	0.0000	0.0000	47.0000
accuracy	0.7257	0.7257	0.7257	0.7257
macro avg	0.3649	0.4961	0.4205	175.0000
weighted avg	0.5339	0.7257	0.6152	175.0000

Confusion Matrix with Labels



Confusion Matrix (Test Data)

True Positive (TP): 0 - Predicted 1 and Actual 1

True Negative (TN): 127 - Predicted 0 and Actual 0

False Positive (FP): 1 - Predicted 1 but Actual 0

False Negative (FN): 47 - Predicted 0 but Actual 1

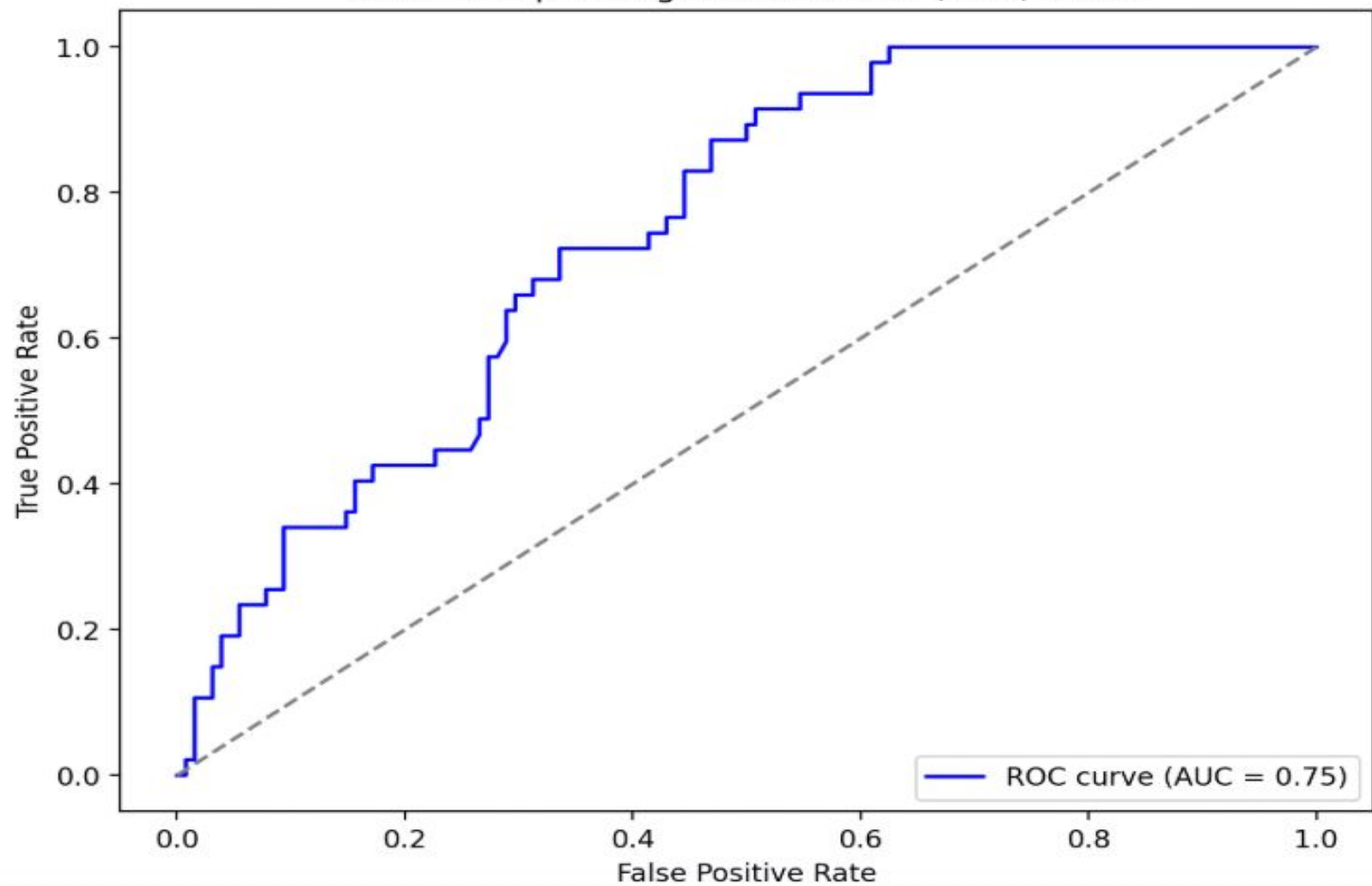
True Positive: The test correctly predicts a positive outcome when the actual outcome is positive.

True Negative: The test correctly predicts a negative outcome when the actual outcome is negative.

False Positive: The test incorrectly predicts a positive outcome when the actual outcome is negative.

False Negative: The test incorrectly predicts a negative outcome when the actual outcome is positive.

Receiver Operating Characteristic (ROC) Curve



Can Fix Threshold Option

- **Diagonal line (Random Classifier):** The line represents the performance of a random classifier. A good classifier should have its ROC curve above this line.

AUC (Area Under the Curve): AUC quantifies the overall ability of the model to discriminate between positive and negative cases.

- An AUC of 0.5 means the model is no better than random guessing.

- An AUC of 1.0 indicates a perfect model.

Adjust Threshold for Classification



Adjusted Metrics at Threshold 0.5:

Accuracy: 0.7257

Precision: 0.3649

Recall: 0.4961

F1 Score: 0.4205

Cross Validation for Classification (Stratified K-Fold)

Reduce Bias - Hyperparameter Option Select value for unique classes

Shape of dataset: (583, 12)

Select Cross-Validation Method

Stratified K-Fold



Based on the number of unique values in the target column (2 unique values), we recommend a maximum of 2 splits.

Select the number of splits for Stratified K-Fold (max 2 based on unique classes):



Custom Class Distribution

Select the ratio for Class 1 (e.g., 75% for Class 1 and 25% for Class 2):



Class 1 will have 75% of the samples, and Class 2 will have 25%.

Applying Stratified K-Fold with 75% Class 1 and 25% Class 2.

To get more Accuracy and prevent Overfit or Underfit

Class 1 will have 75% of the samples, and Class 2 will have 25%.

Applying Stratified K-Fold with 75% Class 1 and 25% Class 2.

Cross-Validation Accuracy (Stratified K-Fold): 0.7221 ± 0.0094

Train Accuracy: 0.7187

Model Fit Status: Generalizing Well

Download Train and Test Dataset

About Upload Dataset Dataset Analysis EDA Process Preprocess **Machine Learning**

Supervised Learning Score Possibility **Download Train and Test Datasets with Results** Future Prediction Unsup

Download Train and Test Datasets with Results

Download Train Dataset with Results ⇄

Download Train Dataset with Results (CSV)

Download Test Dataset with Results

Download Test Dataset with Results (CSV)

Step 11: Future Prediction

[About](#) [Upload Dataset](#) [Dataset Analysis](#) [EDA Process](#) [Preprocess](#) [Machine Learning](#)

[Supervised Learning](#) [Score](#) [Possibility](#) [Download Train and Test Datasets with Results](#) [Future Prediction](#) [Unsupervised Learning](#)

Future Prediction

Please input values for the following features to make a prediction:

Enter value for Trip_Distance_km:

12

Enter value for Time_of_Day:

1

Enter value for Day_of_Week:

2

Enter value for Passenger_Count:

4

2

Enter value for Passenger_Count:

4

Enter value for Traffic_Conditions:

2

Enter value for Weather:

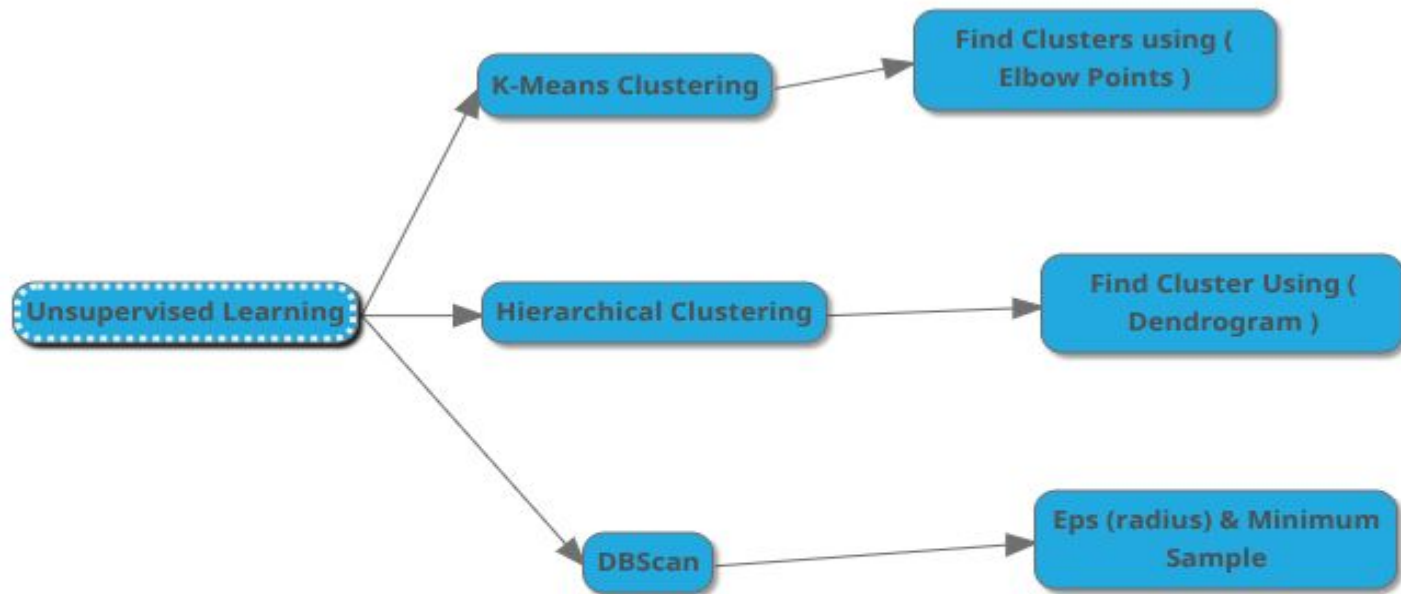
3

Enter value for Base_Fare:

12

Make Prediction

Prediction Result: 36.4206042870544



Step 12 : Unsupervised Learning

About Upload Dataset Dataset Analysis EDA Process Preprocess **Machine Learning**

ning Score Possibility Download Train and Test Datasets with Results Future Prediction **Unsupervised Learning**

Unsupervised Learning Recommendations

Choose Clustering Method

KMeans



KMeans Clustering

Select feature columns for clustering

Trip_Distance_km ×

Time_of_Day ×

Day_of_Week ×

Passenger_Count ×

Traffic_Conditions ×

Weather ×

Base_Fare ×

Per_Km_Rate ×

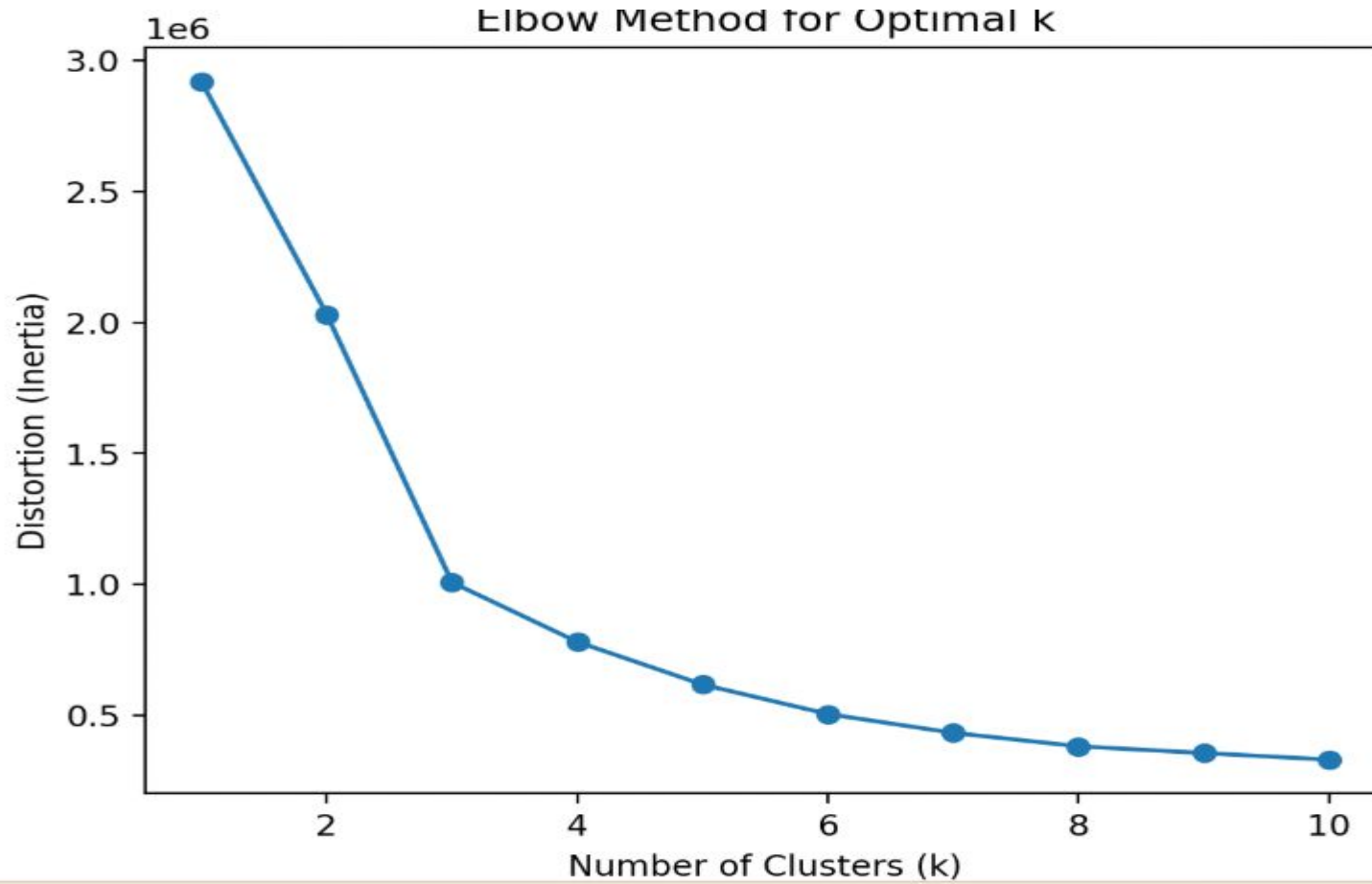


Per_Minute_Rate ×

Trip_Duration_M... ×

Trip_Price ×

K-Means Algorithm (Use Clusters based on Elbow Points)



Enter the optimal k value based on the Elbow method

3

- +

Silhouette Score: 0.38

KMeans Clusters

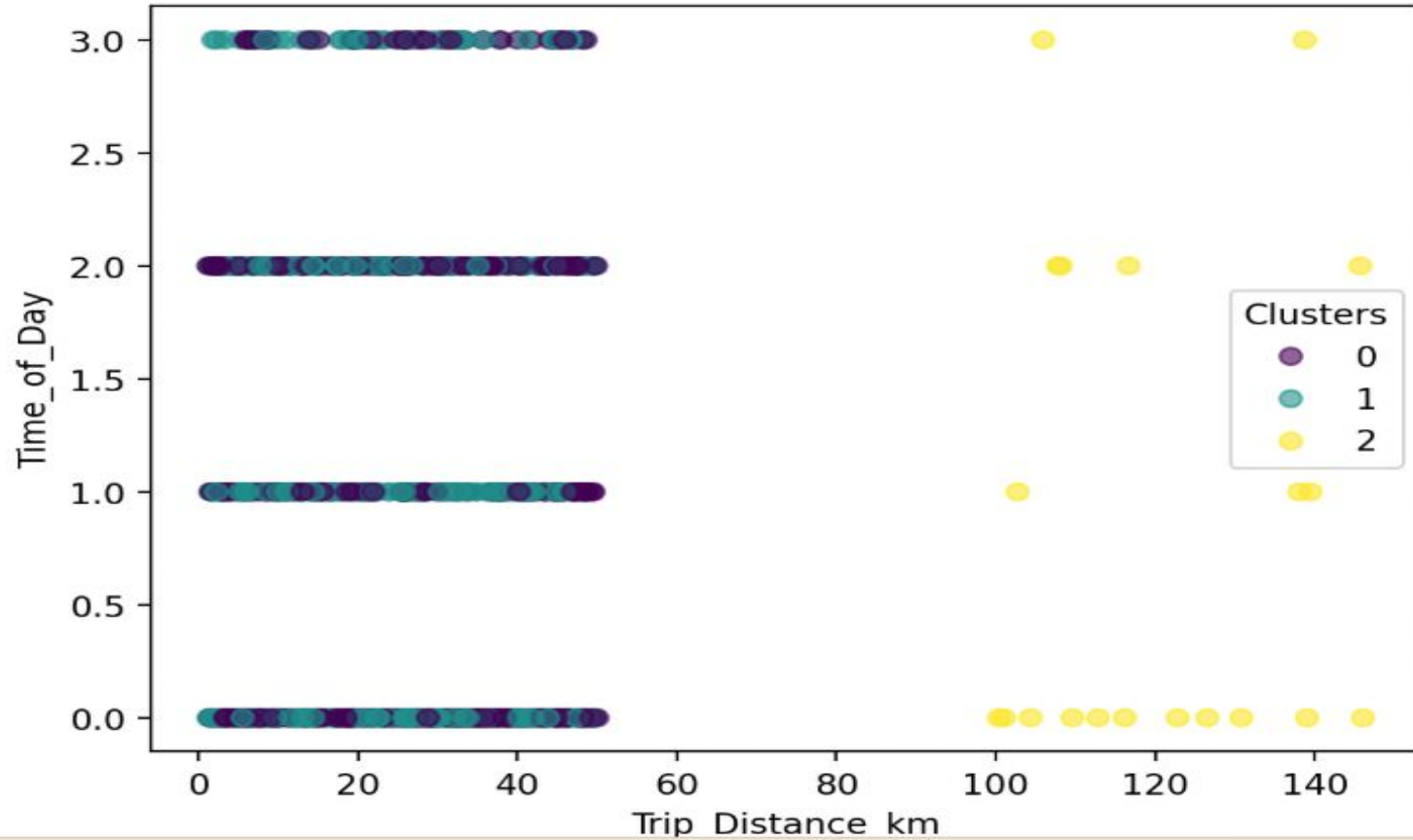
Cluster 0: 467 points

Cluster 1: 513 points

Cluster 2: 20 points

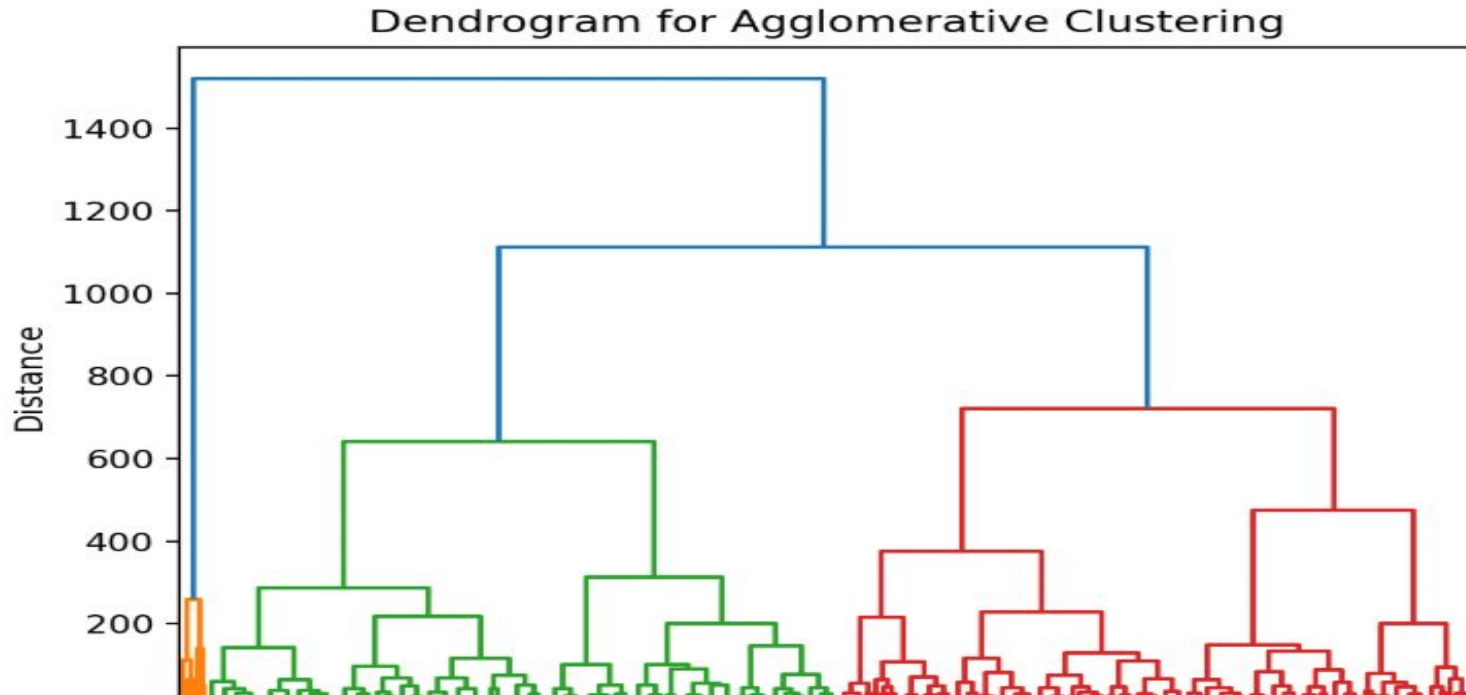


This is not actual Cluster its Just Example to Show

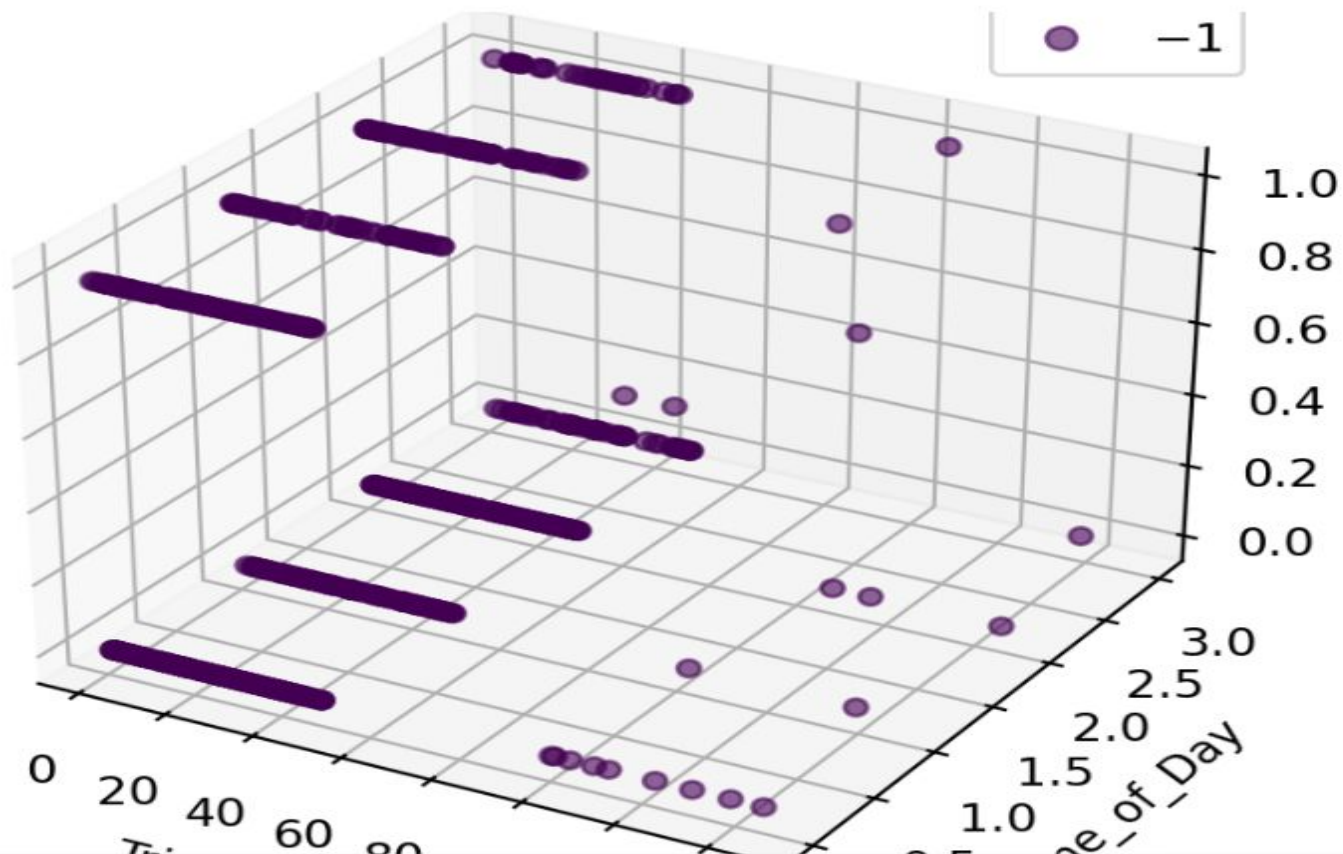


Hierarchical Clustering (Dendrogram)

Dendrogram Method: Find Optimal Number of Clusters



This is not Actual Clustering Point Chart for only Understanding



Thank You

Developer:

TC Antony - [Data Scientist]