# Capstone Project - 2

## Team 2

## Taxi Mobility Surge Price Prediction

**Team Members**

**Bhanu Pratap Shahi**
**Sanchita Paul**
**Shubham Kumar**
**Sumanta Muduli**

# Content

## Title

- Problem Statement
- Data Summary
- Comparing features with Surge Pricing Type
- Feature Selection
- Models used
- Which model did we choose and why?
- Challenges
- Conclusion

**AI**

# Problem Statement

- The goal is to build a predictive model which can help Sigma Cabs in predicting Surge Pricing Type proactively.
- This will help them in matching the right priced cabs with the right customers quickly and efficiently.

**AI**

# Data Summary:

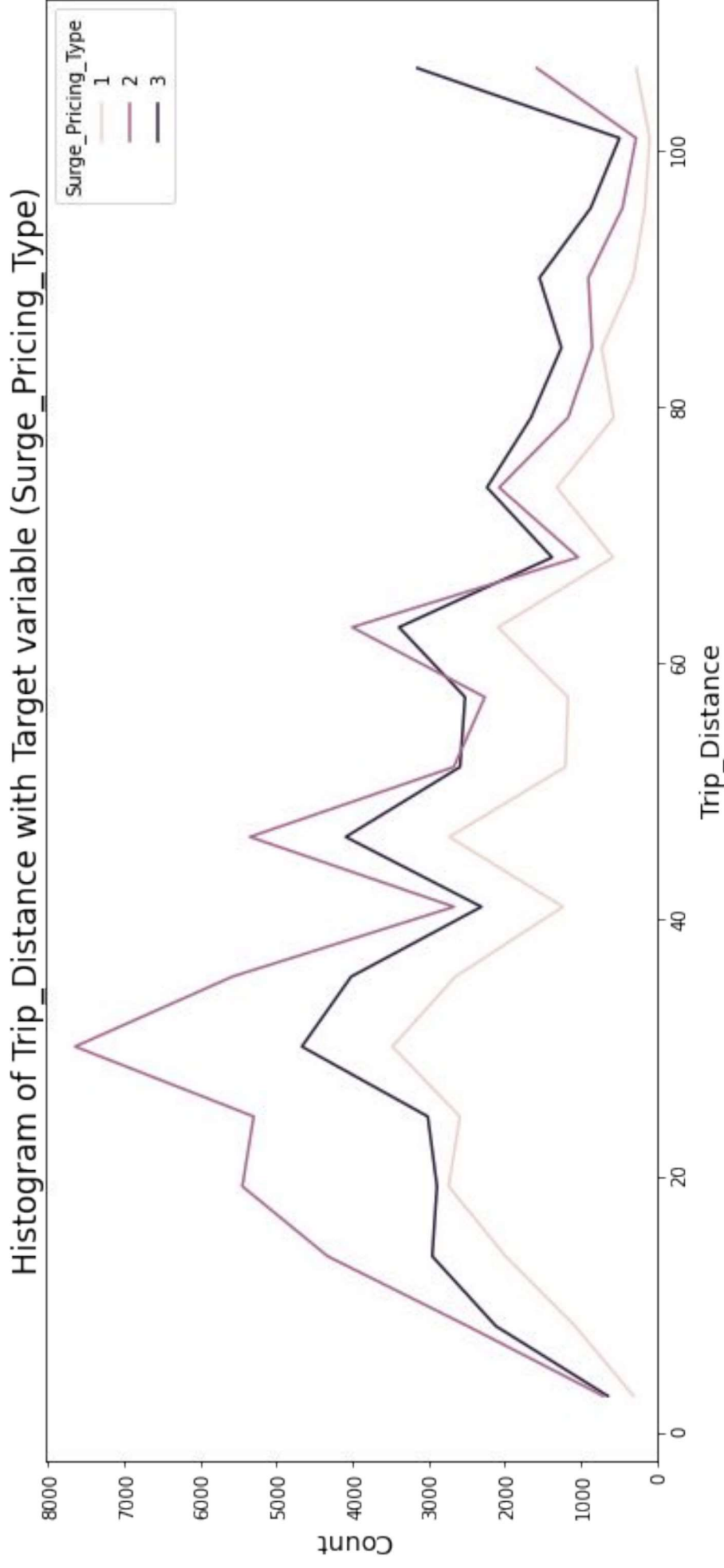## Data set name  Data Sigma Cabs

## Shape
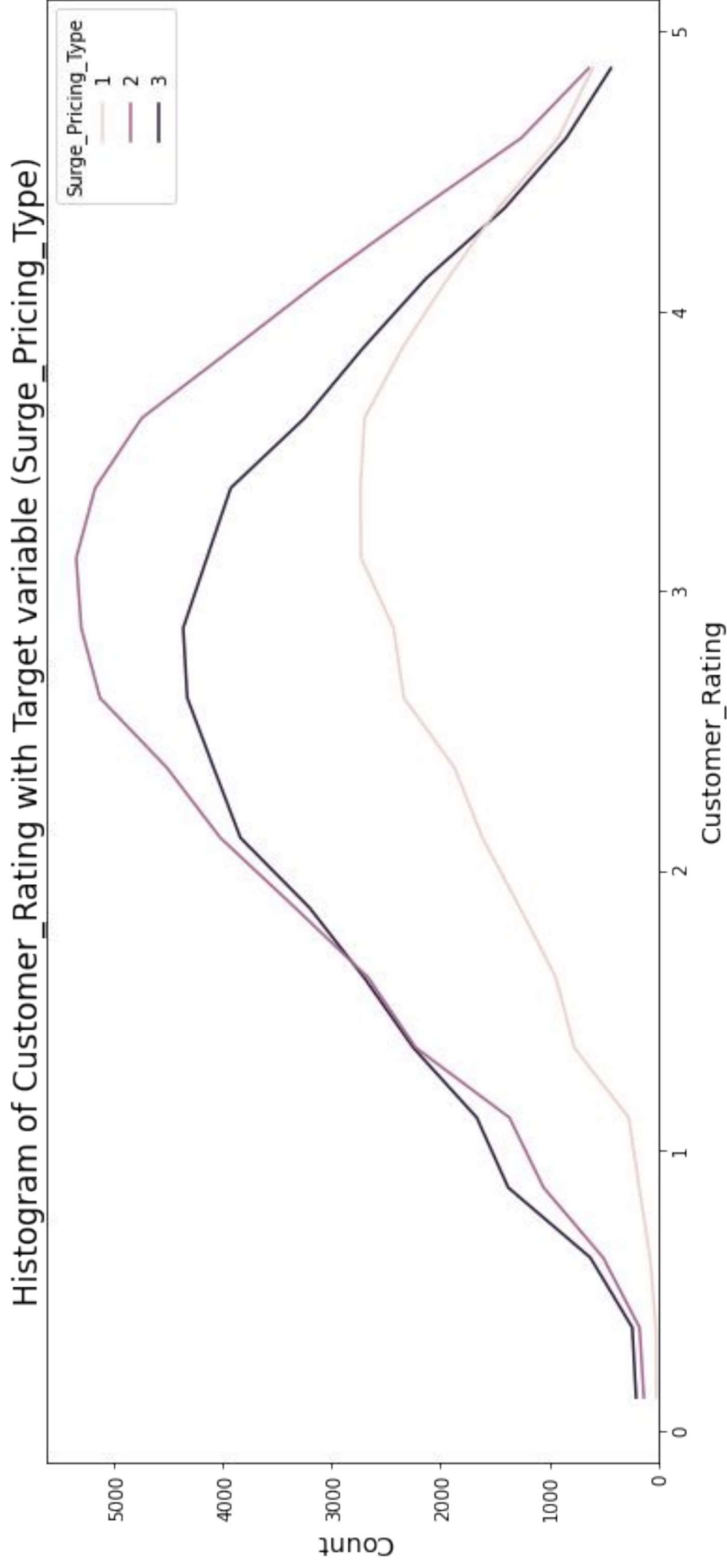
- **Rows -- 131,662**
- **Columns--14**

## Features

Trip_ID,Trip_Distance,Type_of_Cab,Customer_since_months,
Life_Style_Index, Confidence_Life_Style_Index,Destination_Type,
Customer_Rating,Cancellation_Last_1Month,Var1,Var2,Var3,Gender,
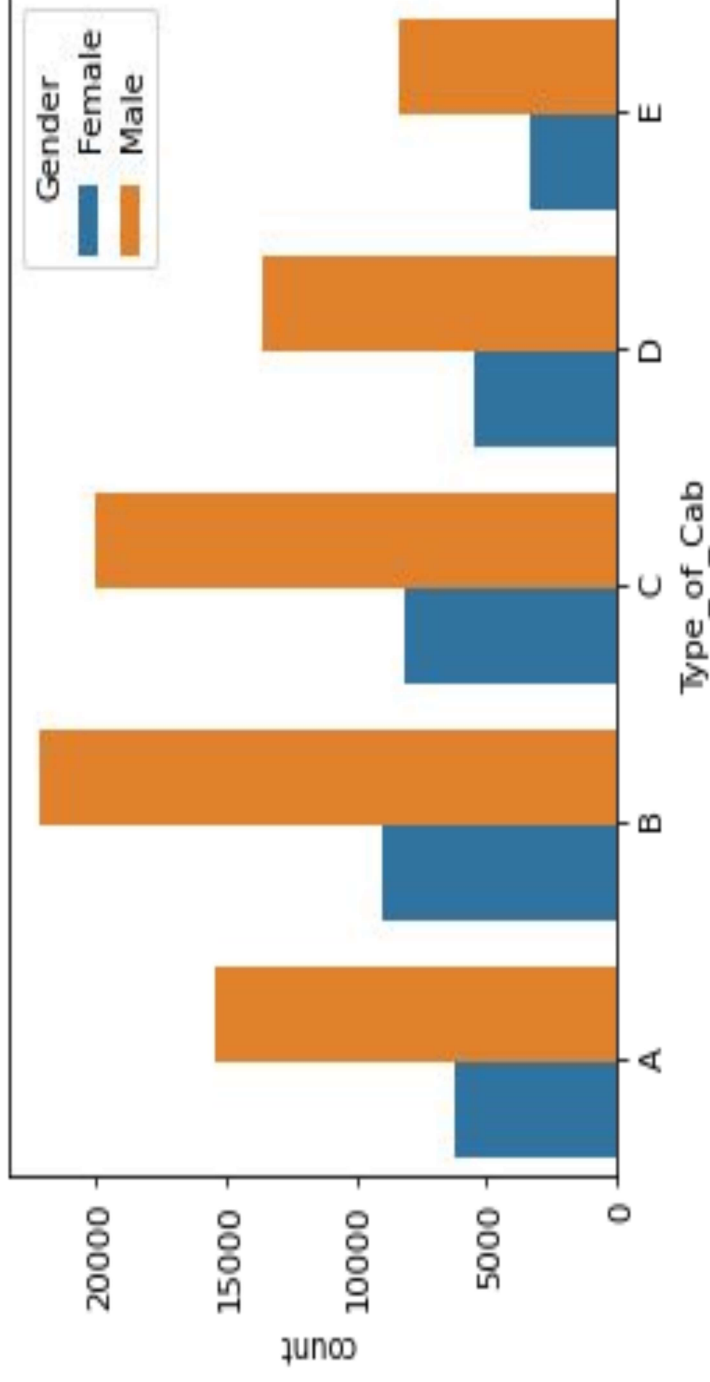Surge_Pricing_Type

# Comparing Trip Distance with Surge Pricing Type:

**AI**



Histogram of Trip_Distance with Target variable (Surge_Pricing_Type)

# Comparing Customer Rating with Surge Pricing Type:

Histogram of Customer_Rating with Target variable (Surge_Pricing_Type)

# Count of Type of Cab with Gender Filter

**AI**

# Features selection:

**Methods used:**

- **Extra Tree Classifier**
- **ANOVA**
- **Chi-Square**

# Extra Trees Classifier:

Comparison of different Features



Feature Importances

Feature Labels

ANOVA:

# Chi-Square:



P-values for categorical features
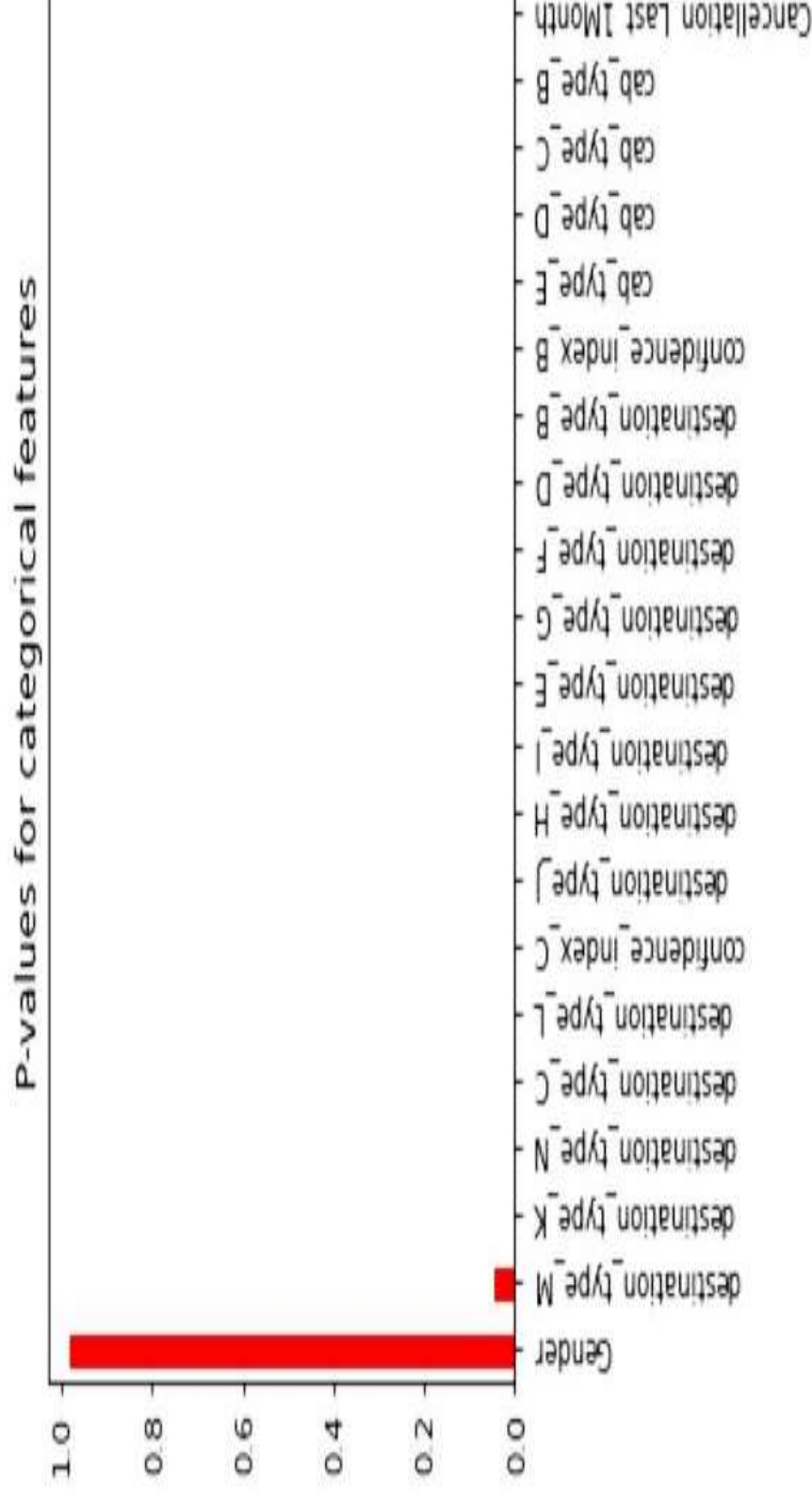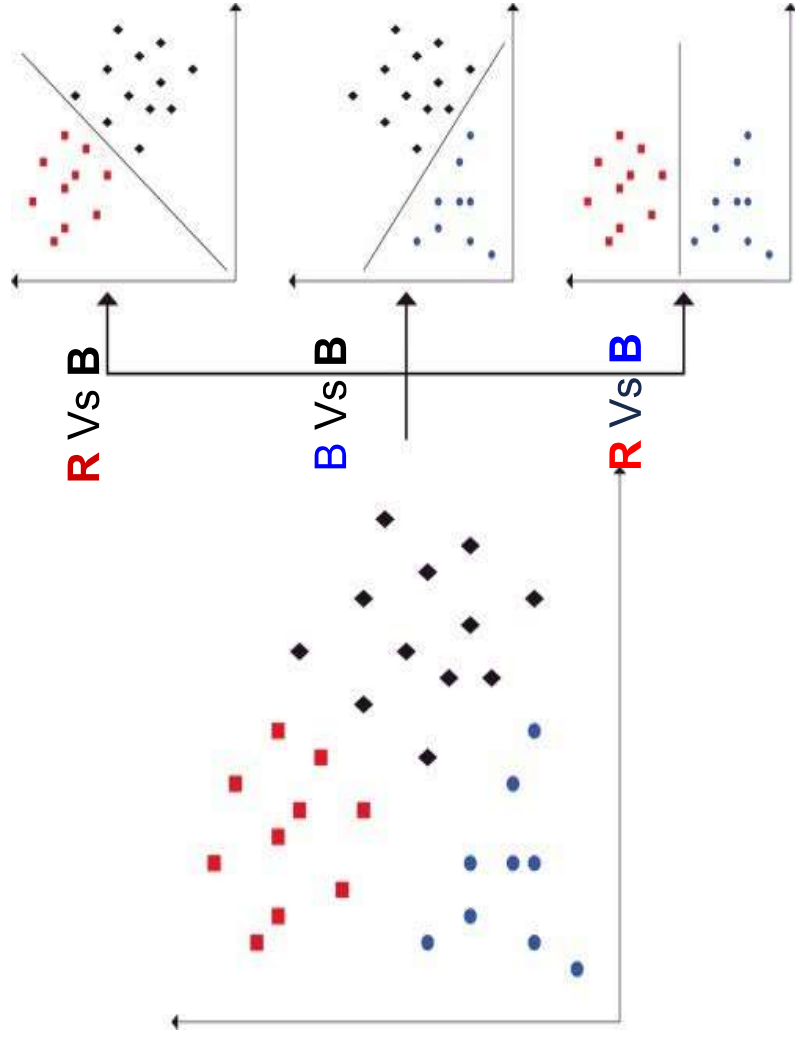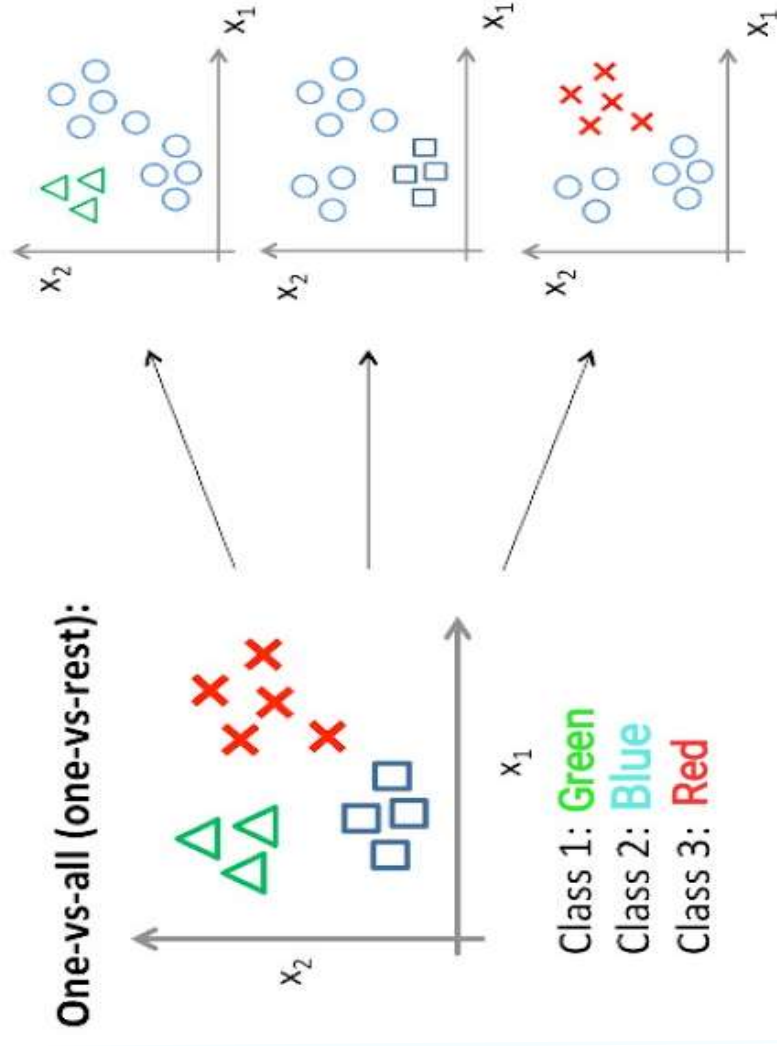
**AI**

# Models used:

- Logistic Regression Classifier
- SVM Classifier
- Random Forest Classifier
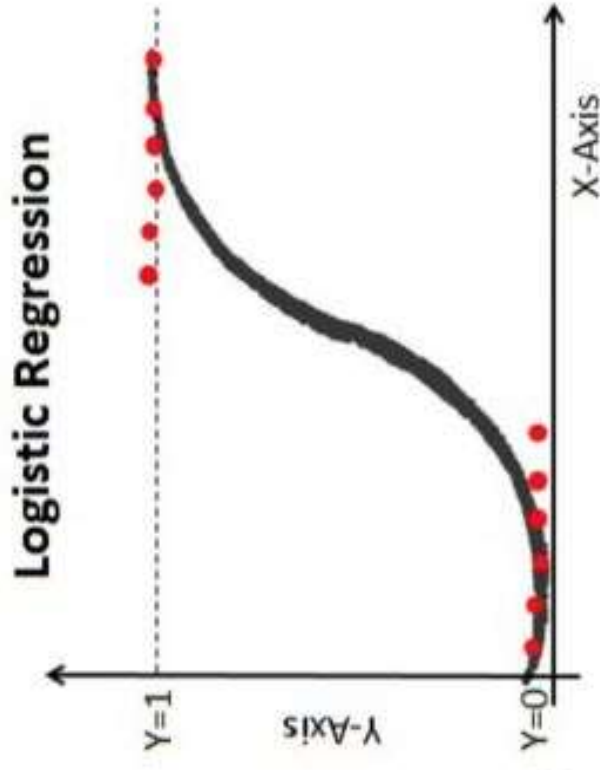- XGBoost Classifier

# One vs One and One vs Rest:



One-vs-all (one-vs-rest):

Class 1: Green
Class 2: Blue
Class 3: Red

R Vs B

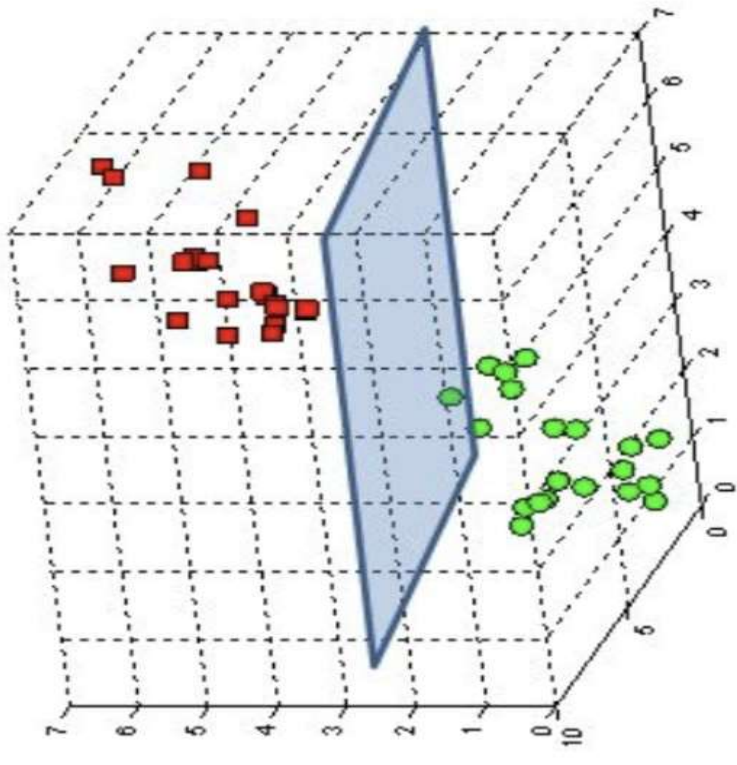B Vs B

R Vs B

One vs Rest

One vs One

# Logistic Regression:

- One vs Rest approach ("ovr")

- Hyperparameter Tuning( Bayesian Optimisation)-C:0.001, solver:"lbfgs",penalty=l2

- Metric Scores- Accuracy=72%, Precision=72%, Recall=70% & f1_score=71%



Logistic Regression

Y=1

Y-Axis

Y=0

X-Axis

# Support Vector Machine:

- One vs One approach ("ovo")

- Parameters - C:1, degree =3,

- Kernel - Poly Kernel is giving us the best results. Accuracy i.e 72%, Precision=73%, Recall=70% & f1_score=70%
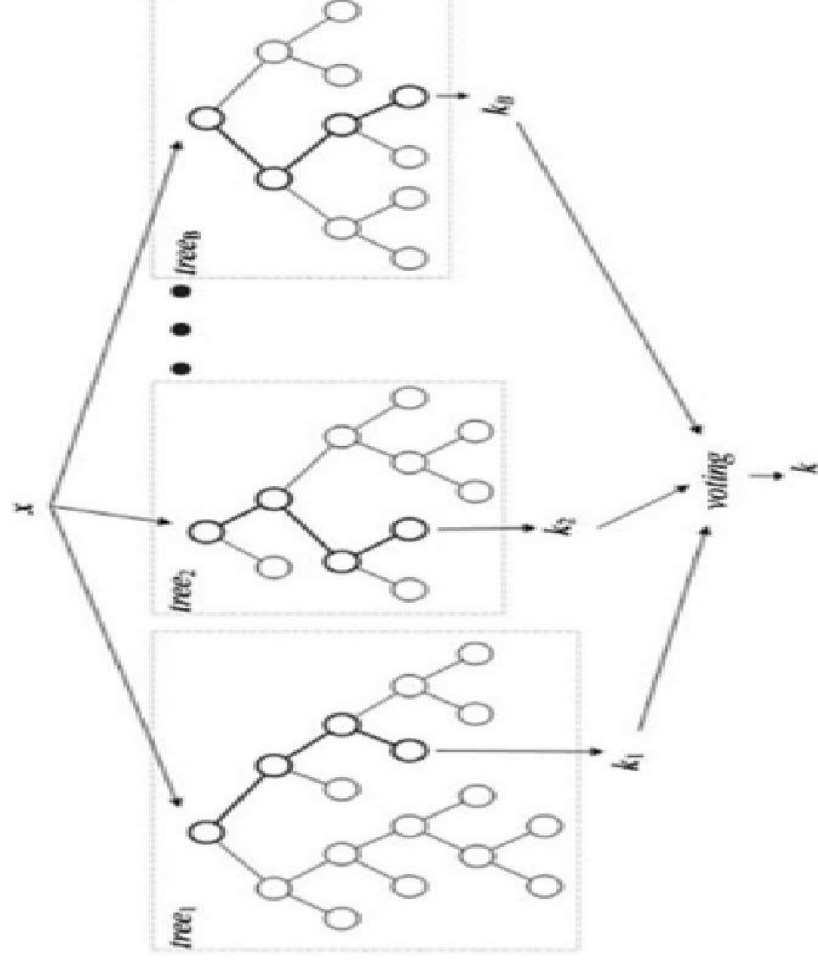
# Random Forest Classifier:



- Hyper parameter Tuning(Bayesian Search)-
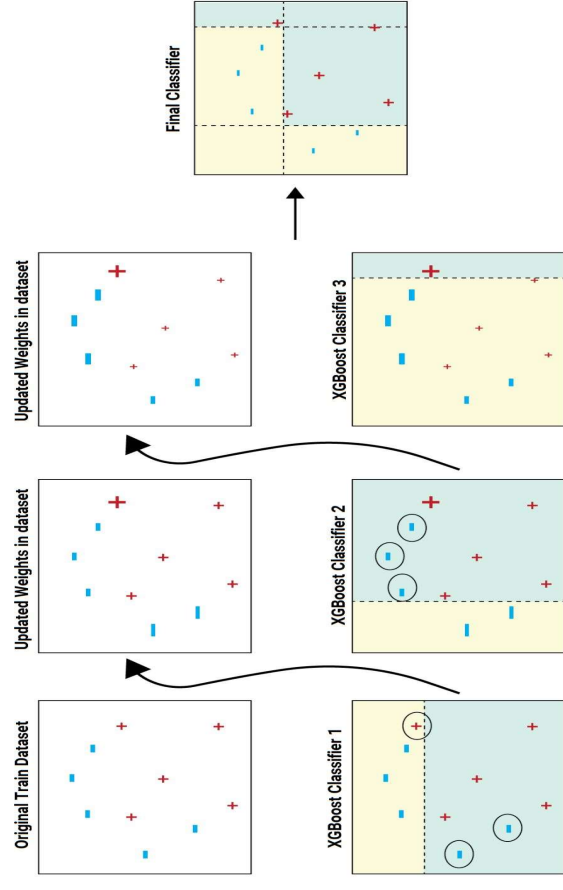  ('max_depth', 8),
  ('min_samples_leaf', 10),
  ('min_samples_split', 50),
  ('n_estimators', 100)
- accuracy= 72%,
  precision=73%,
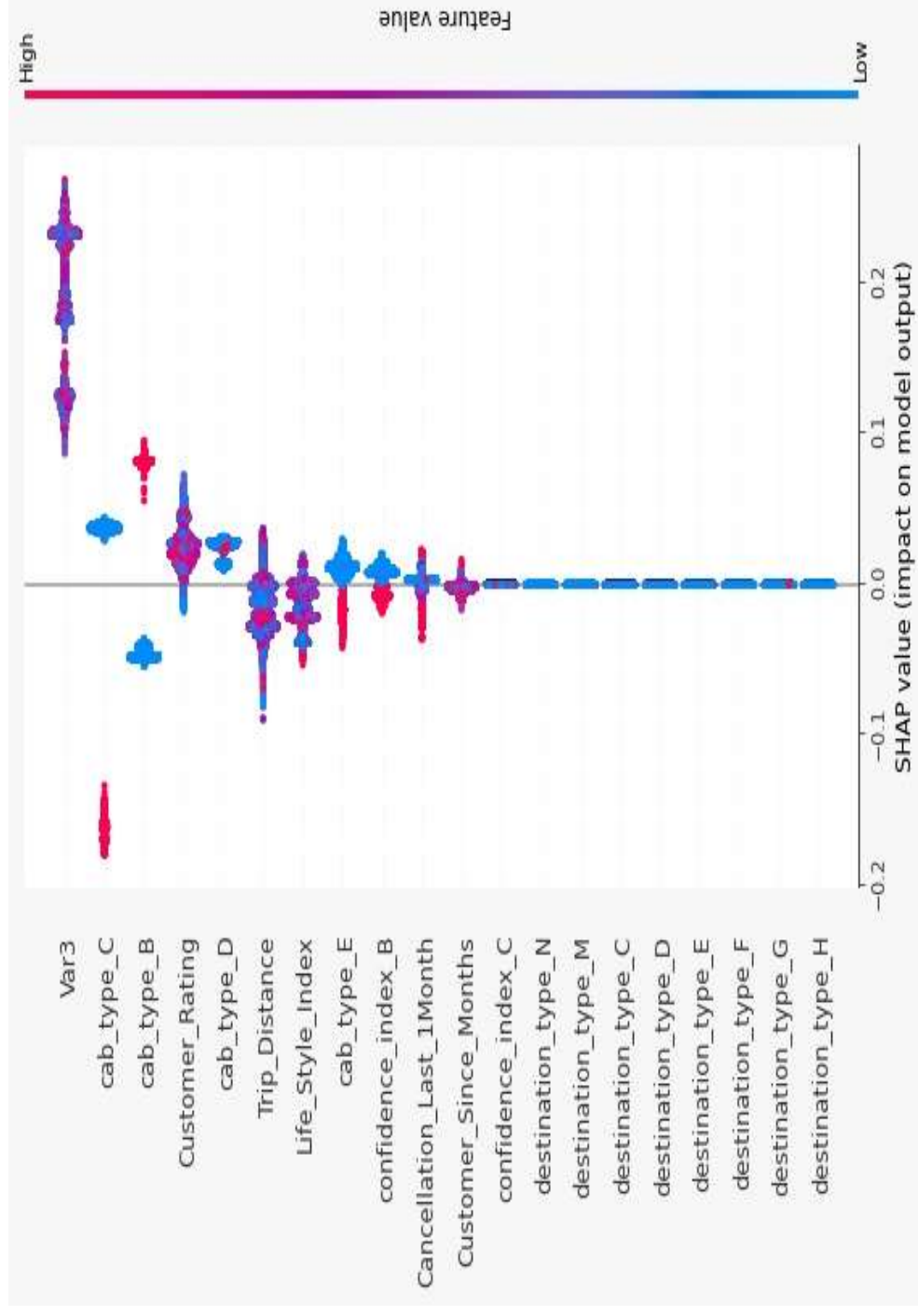  recall=70%,f1_score=71%

# XGBoost Classifier:



- Hyperparameters-gamma=0, learning_rate=0.1, max_depth=15, n_estimators=100, objective='multi:softprob'

- Metric Scores- accuracy=72%, precision=73%, recall=70%,f1_score=71%

# SHAP Values:

**AI**

# Which model did we choose and why?

- We choose logistic regression as it's evaluation scores is very similar to other complicated models but it is computationally cheaper and more interpretable.
- Accuracy : 72%
- Recall : 72%
- Precision : 72%
- This is the most consistent performing model with same scores for all metrics.

**AI**

# Challenges

- Lots of NaN values in the dataset.
- Some features like Var1,Var2,Var3 are not clearly explained.
- Choosing the right encoding technique for categorical features.
- Choosing the right features for modelling.
- Faced issues while running the models as the dataset is large.
- Choosing the right models as there is not much difference in accuracy.

**AI**

# Conclusion

- We build a predictive model which can help Sigma Cabs in predicting Surge Pricing Types proactively.
- This will helps in matching the right cab with the right customer quickly and efficiently
- They can increase their customer base and profit by providing better services.

# Q & A