



# Predicting Restaurant Rating Star

Presentation by **Hoda Shoghi**

# Project Objective

## ■ Goal:

Develop a model to predict restaurant star ratings.

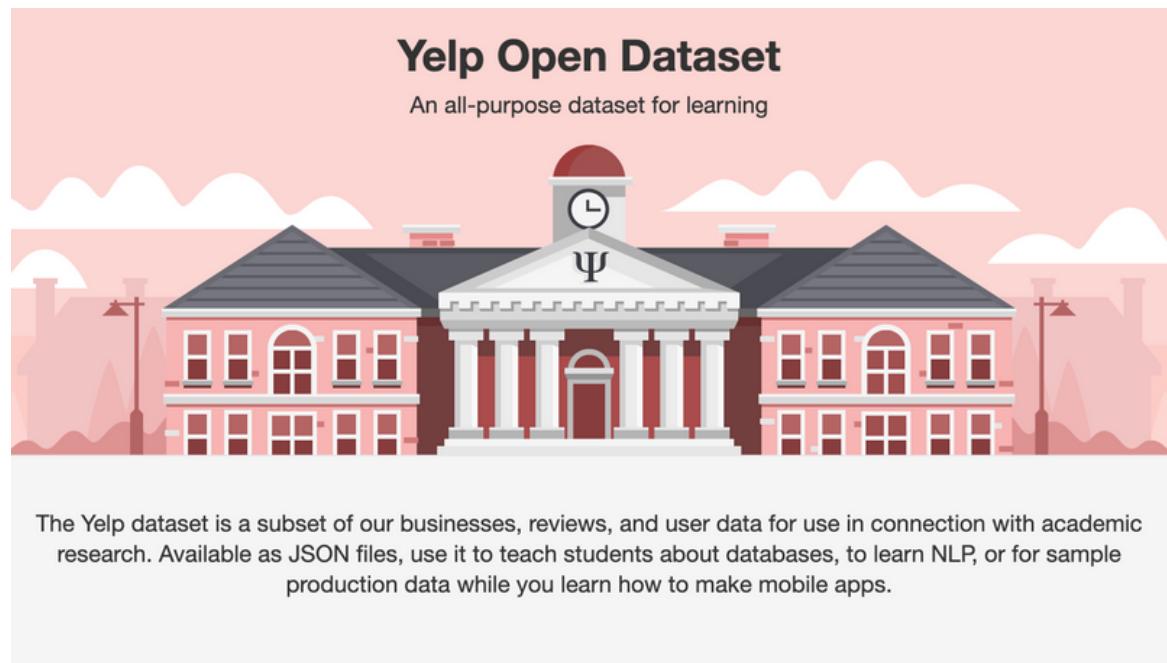
## ■ Why?:

To understand what influences ratings and help future restaurant owners succeed.

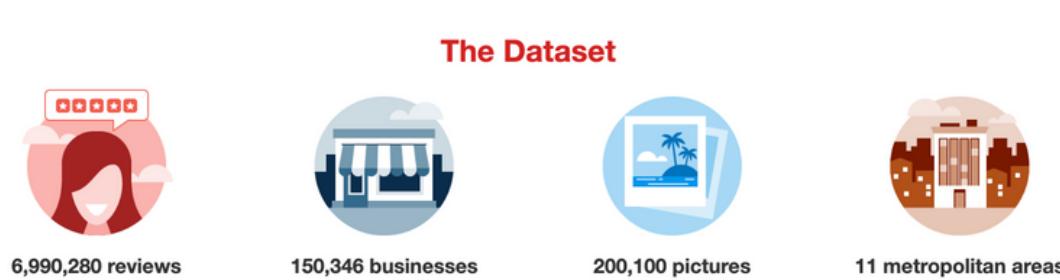
## ■ Purpose:

Offer valuable insights to potential restaurant owners and stakeholders about key factors that can influence a restaurant's success.

# Dataset Overview



The Yelp dataset is a subset of our businesses, reviews, and user data for use in connection with academic research. Available as JSON files, use it to teach students about databases, to learn NLP, or for sample production data while you learn how to make mobile apps.



**Source:**

Yelp's Public Dataset

**Size:**

150346 Businesses

**Characteristics**

Examples: Price Range, Cuisine Type, Location, Casual Dining, Ambience, etc.

**Target Response:**

Restaurant's Star Rating (1 to 5 stars)

# Pre-Processing

Handpick the best ingredients

Refine them

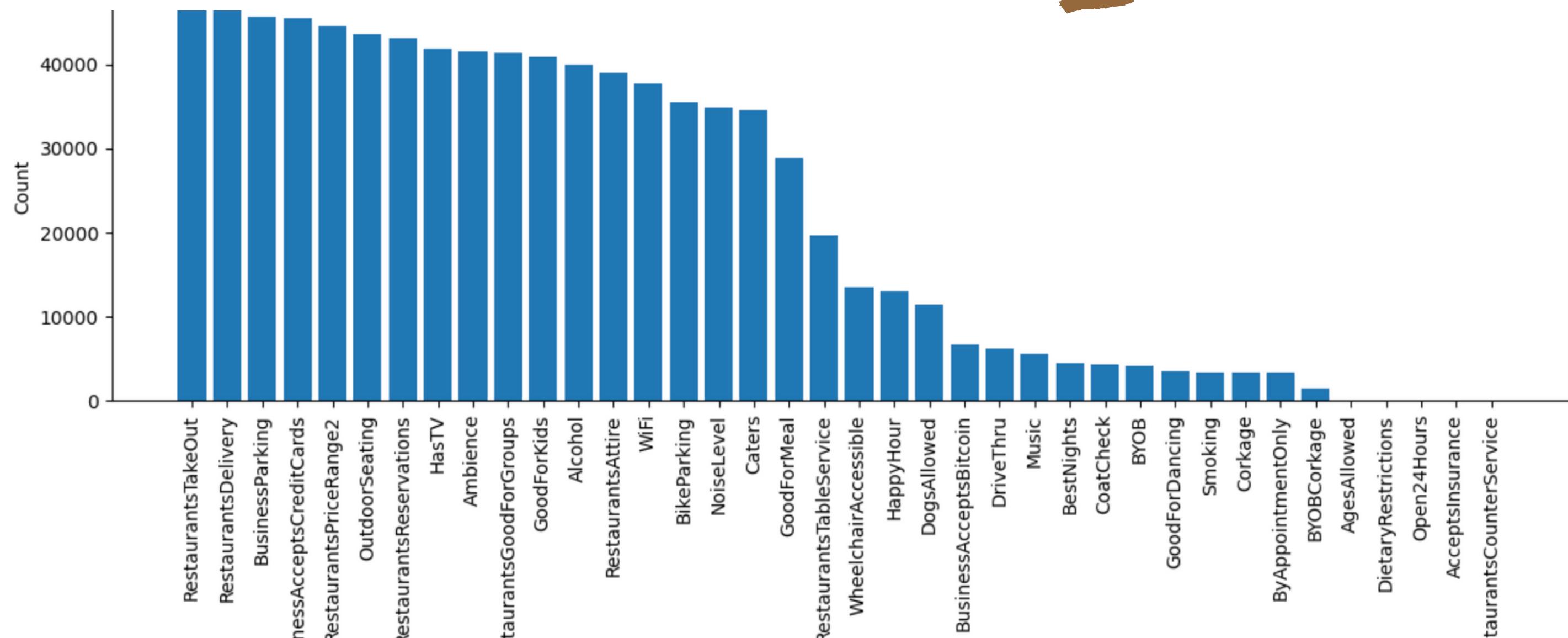
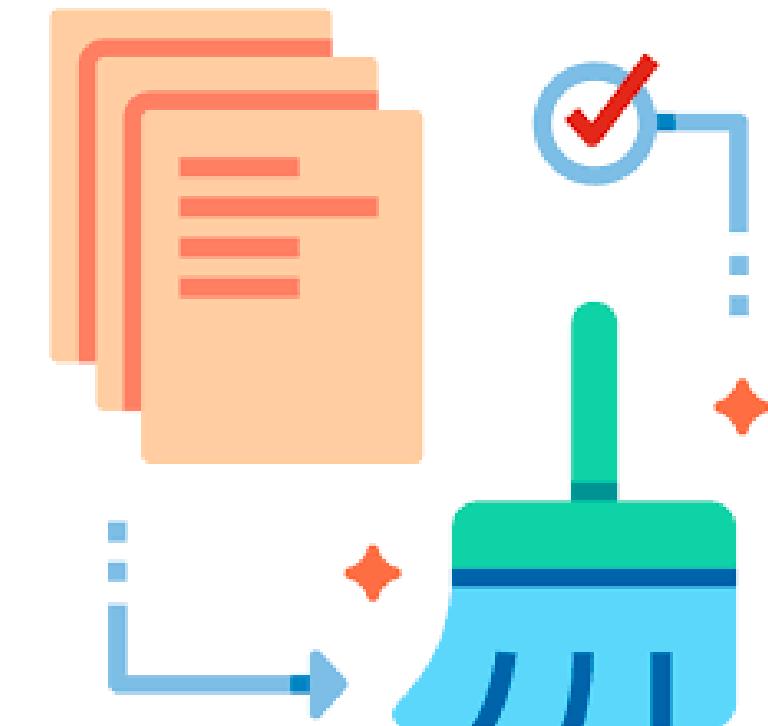
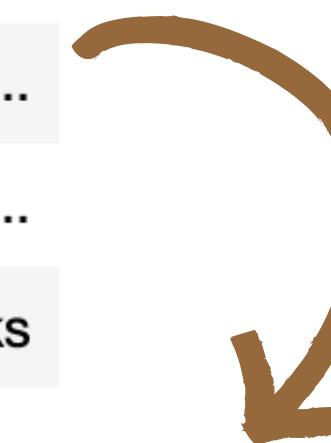
Mix them in just the right  
PROPORTIONS

Getting the foundation right!!

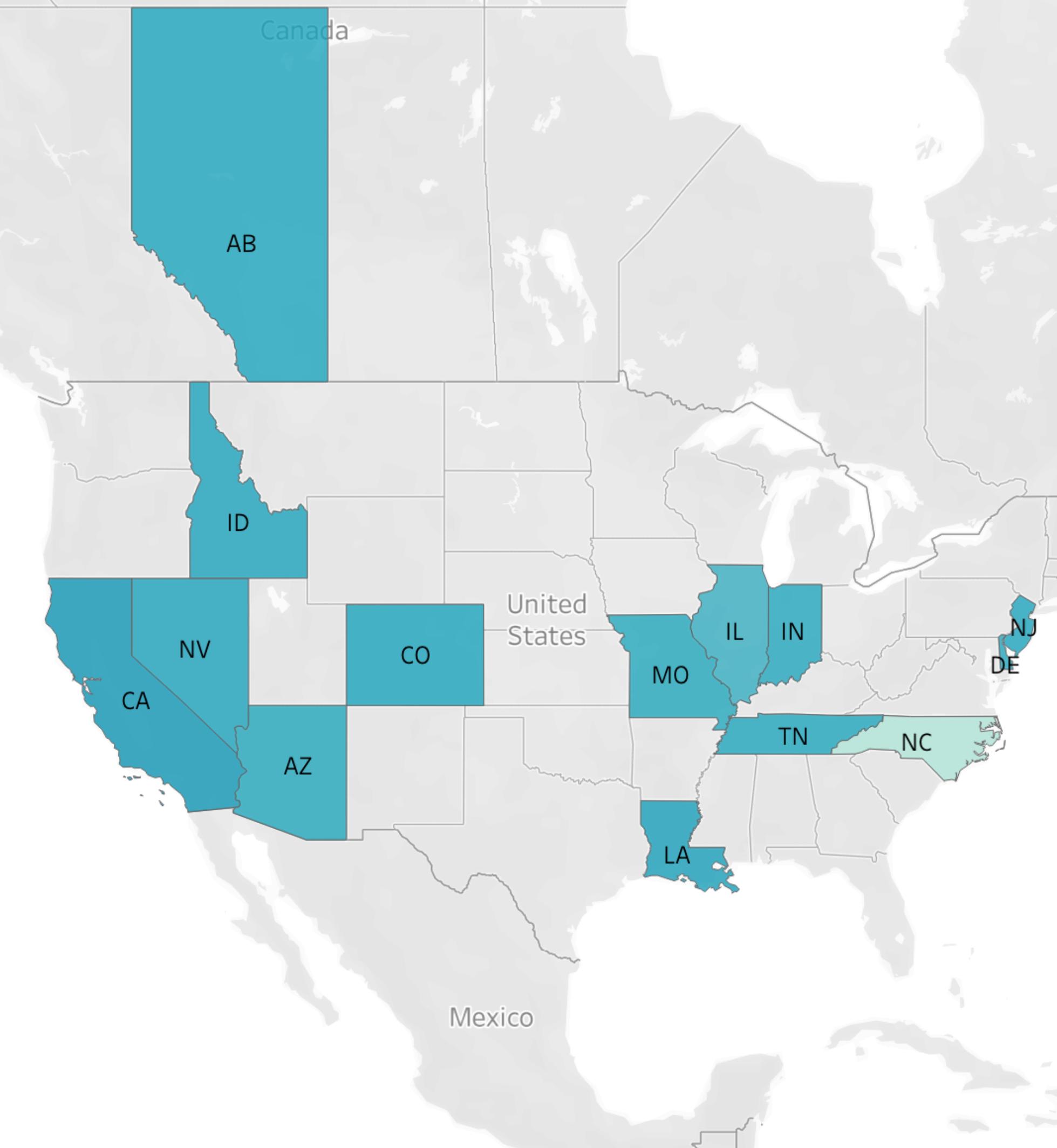
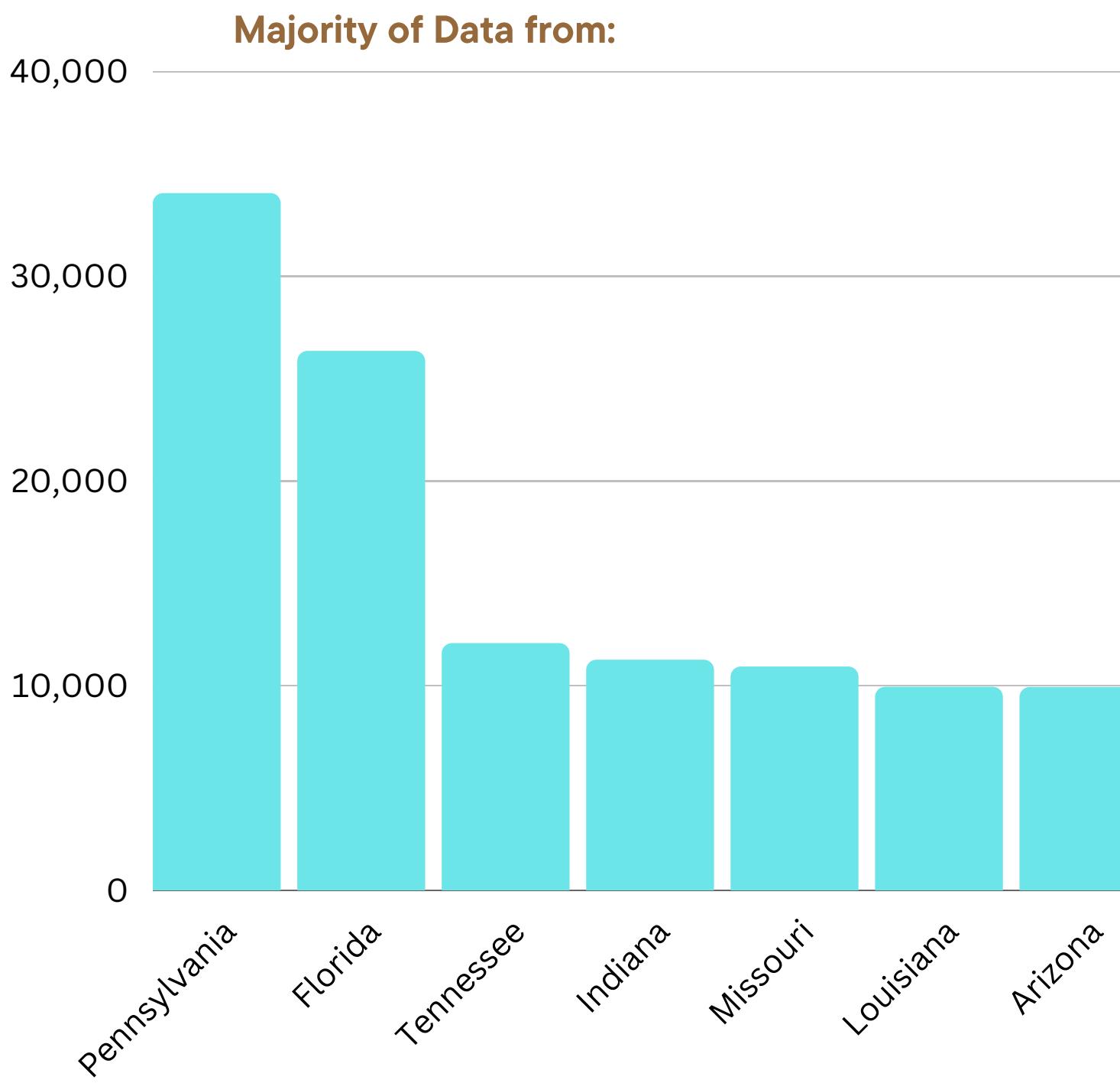


# Data Cleaning and Pre-processing

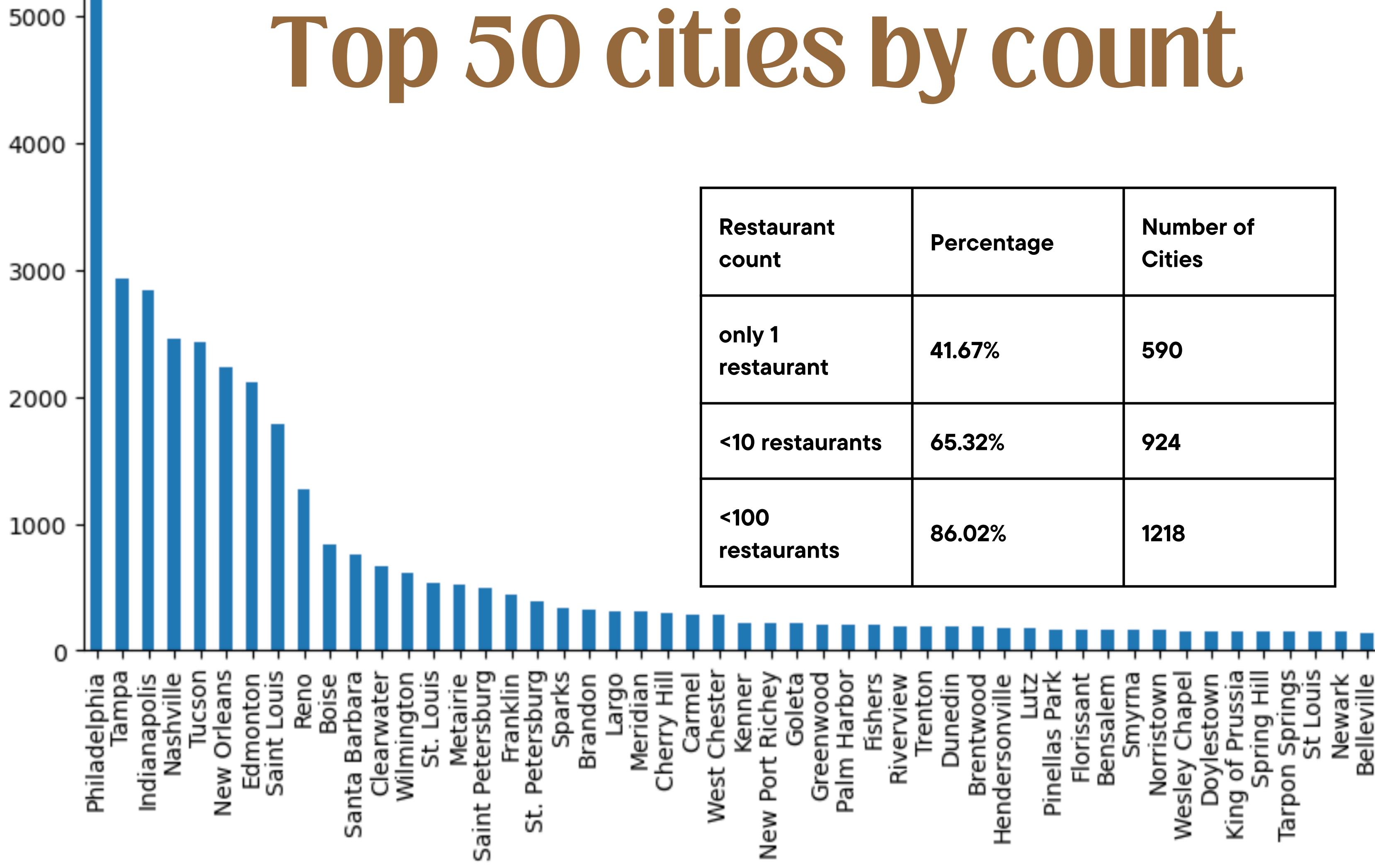
attributes	categories
{'RestaurantsDelivery': 'False', 'OutdoorSeati...}	restaurants, food, bubble tea, coffee & tea, b...
{'BusinessParking': 'None', 'BusinessAcceptsCr...}	burgers, fast food, sandwiches, food, ice crea...
{'Caters': 'True', 'Alcohol': 'u'full_bar'", '...	pubs, restaurants, italian, bars, american (tr...
{'RestaurantsAttire': "casual", 'Restaurants...}	ice cream & frozen yogurt, fast food, burgers,...
{'Alcohol': "none", 'OutdoorSeating': 'None'...}	vietnamese, food, restaurants, food trucks



# Data Geography

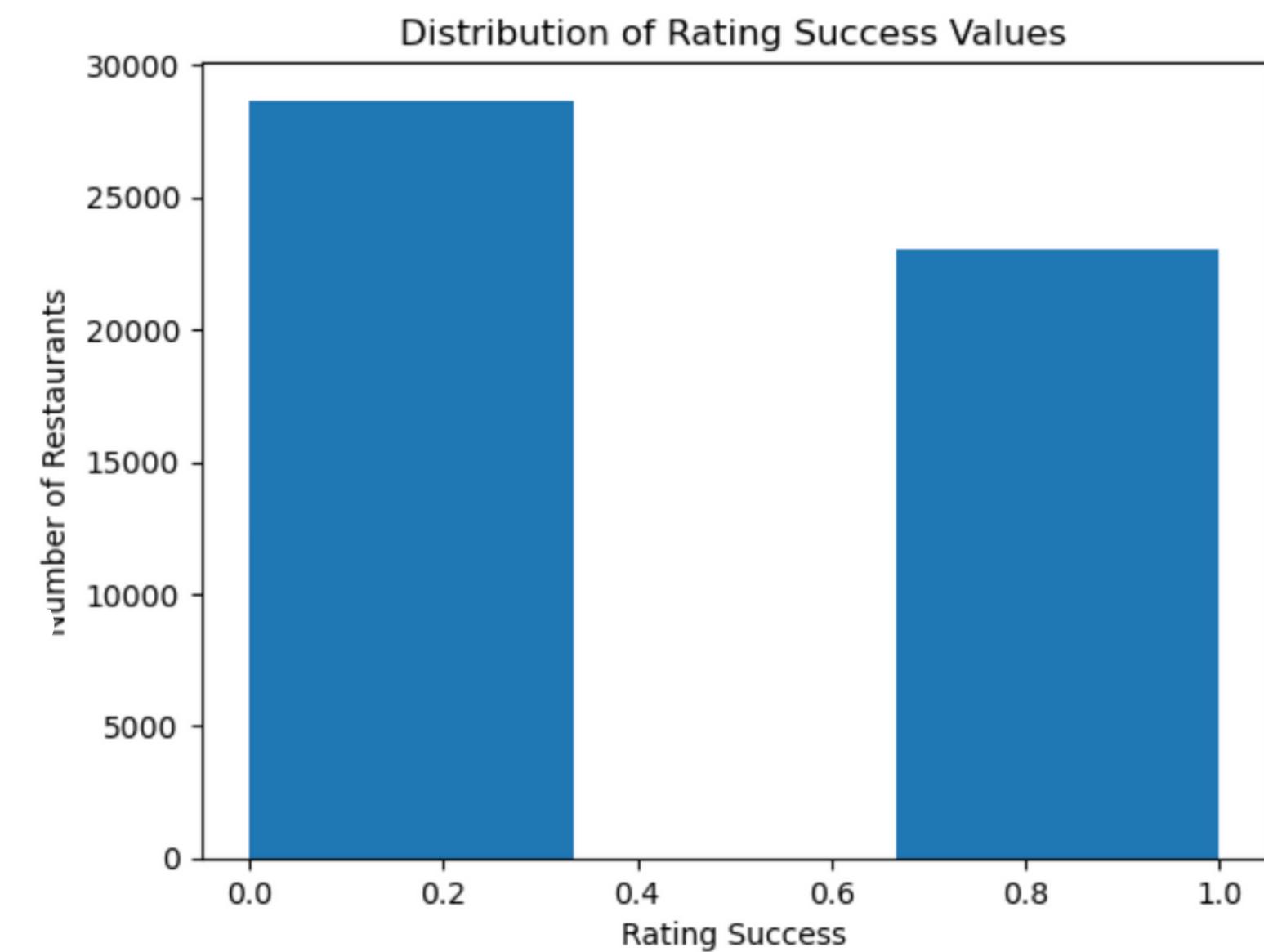
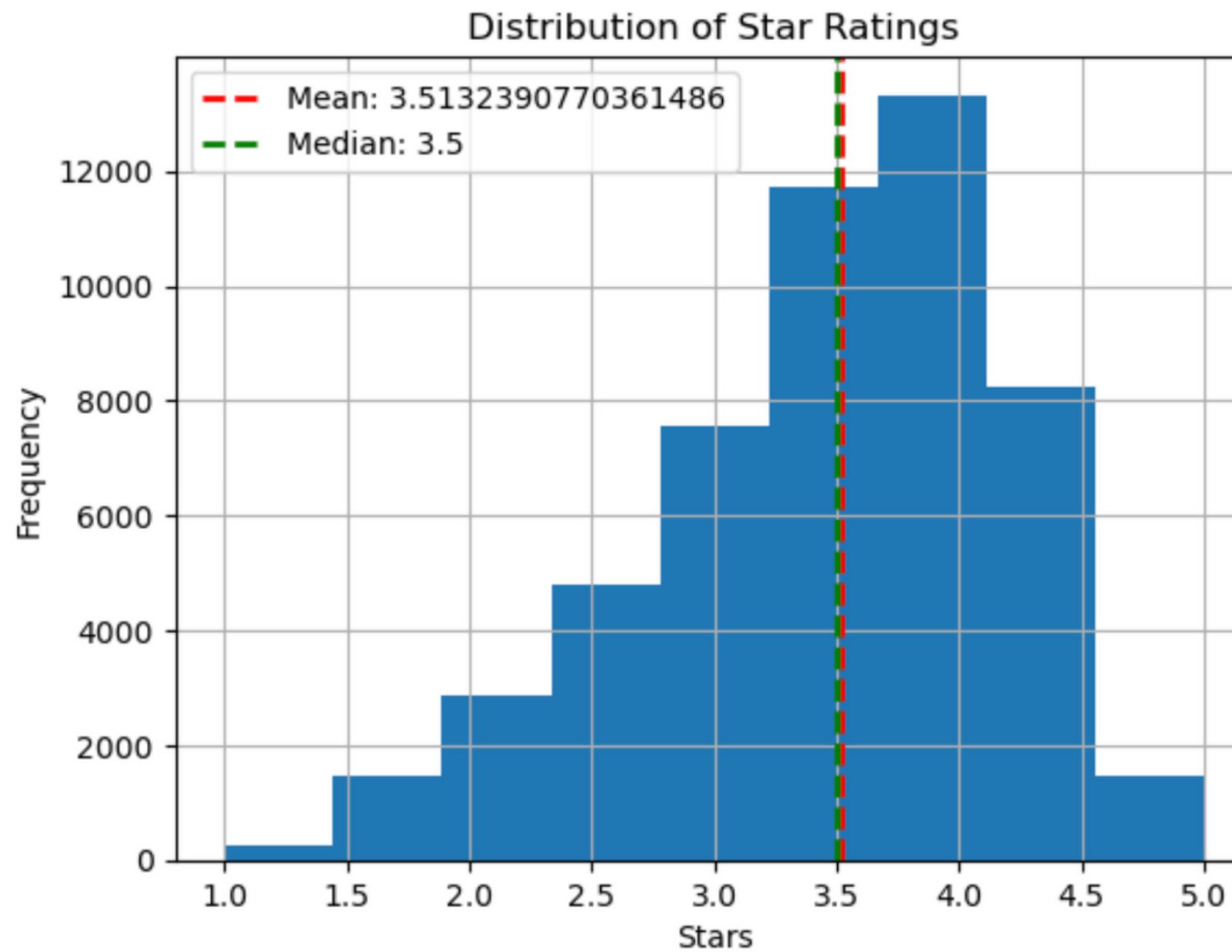


# Top 50 cities by count

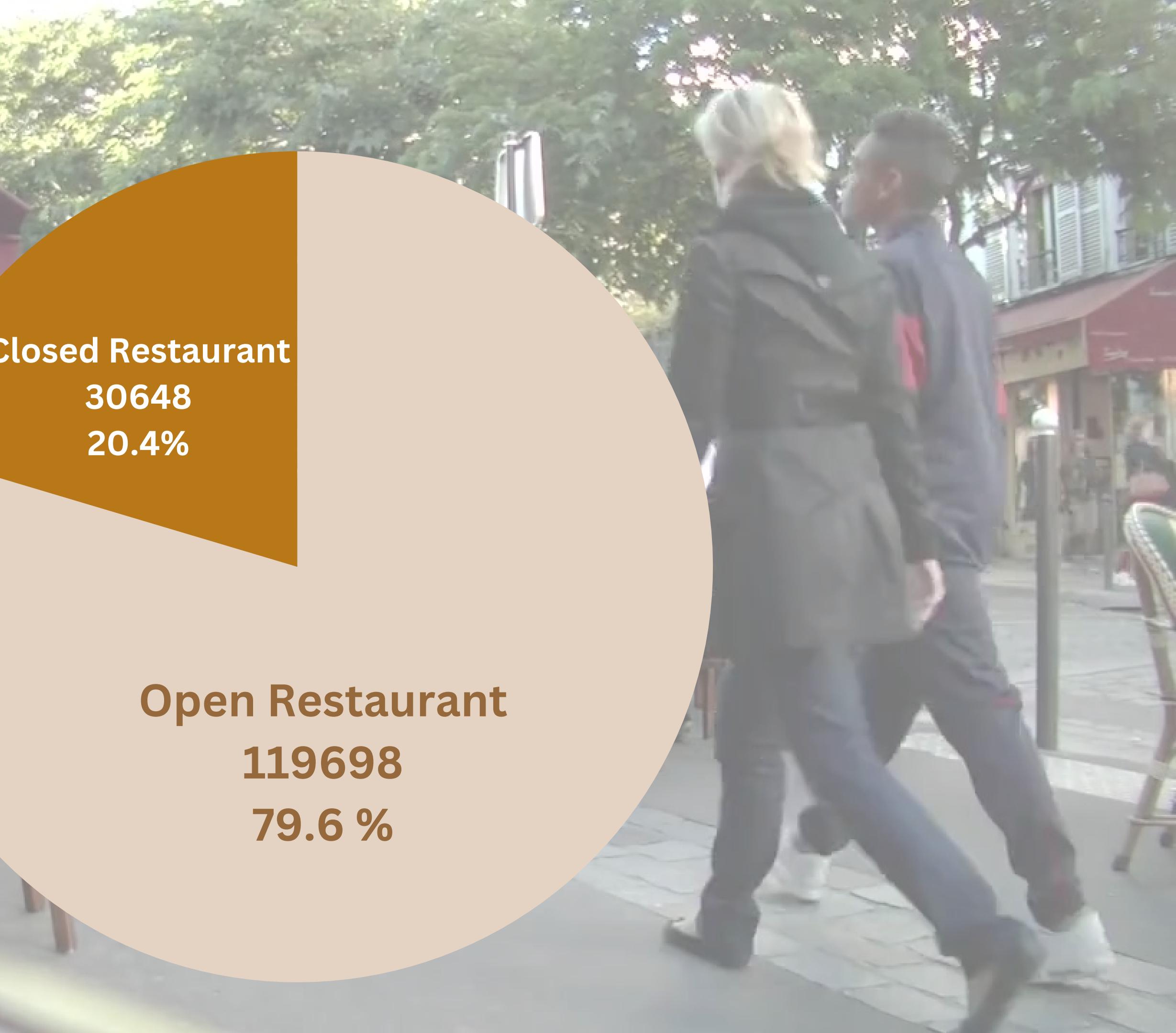


# Target Response

converted to  
binary target variable



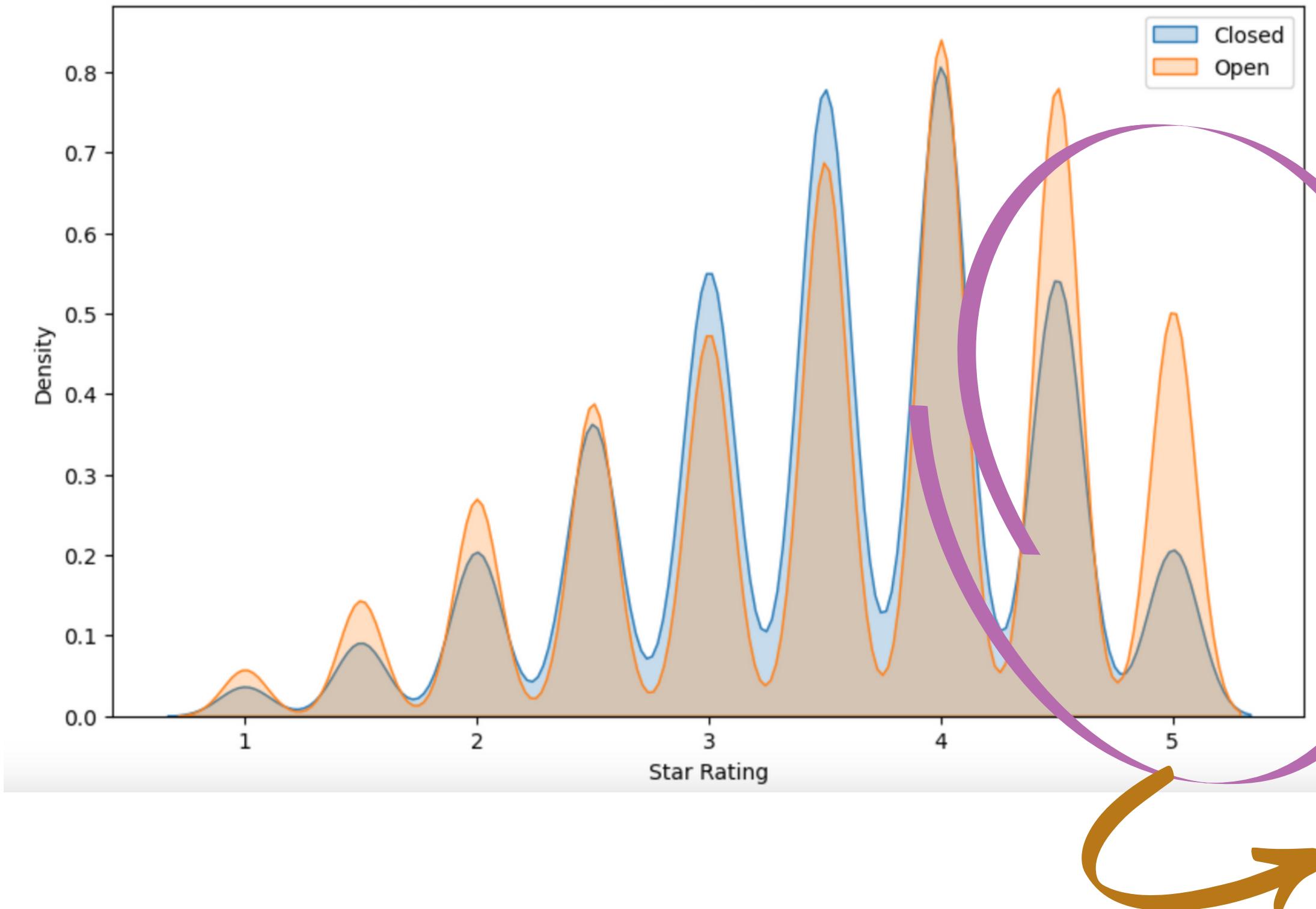
# Permanently Closed!?



# Average rating, Open Vs Cloesd



# Distribution of ratings for open and closed restaurants.



Rating stars	is_open = 0 (%)	is_open = 1 (%)
1.0	1.03	1.40
1.5	2.55	3.47
2.0	5.70	6.50
2.5	10.13	9.37
3.0	15.43	11.47
3.5	21.73	16.59
4.0	22.50	20.24
4.5	15.16	18.83
5.0	5.78	12.14

# Machine Learning Models

**Logistic Regression**

70 %

**Random Forest**

70 %

**Random Forest with  
Feature reduction**

69 %

**Gradient Boosting**

70 %

**Gradient Boosting in Pipeline**

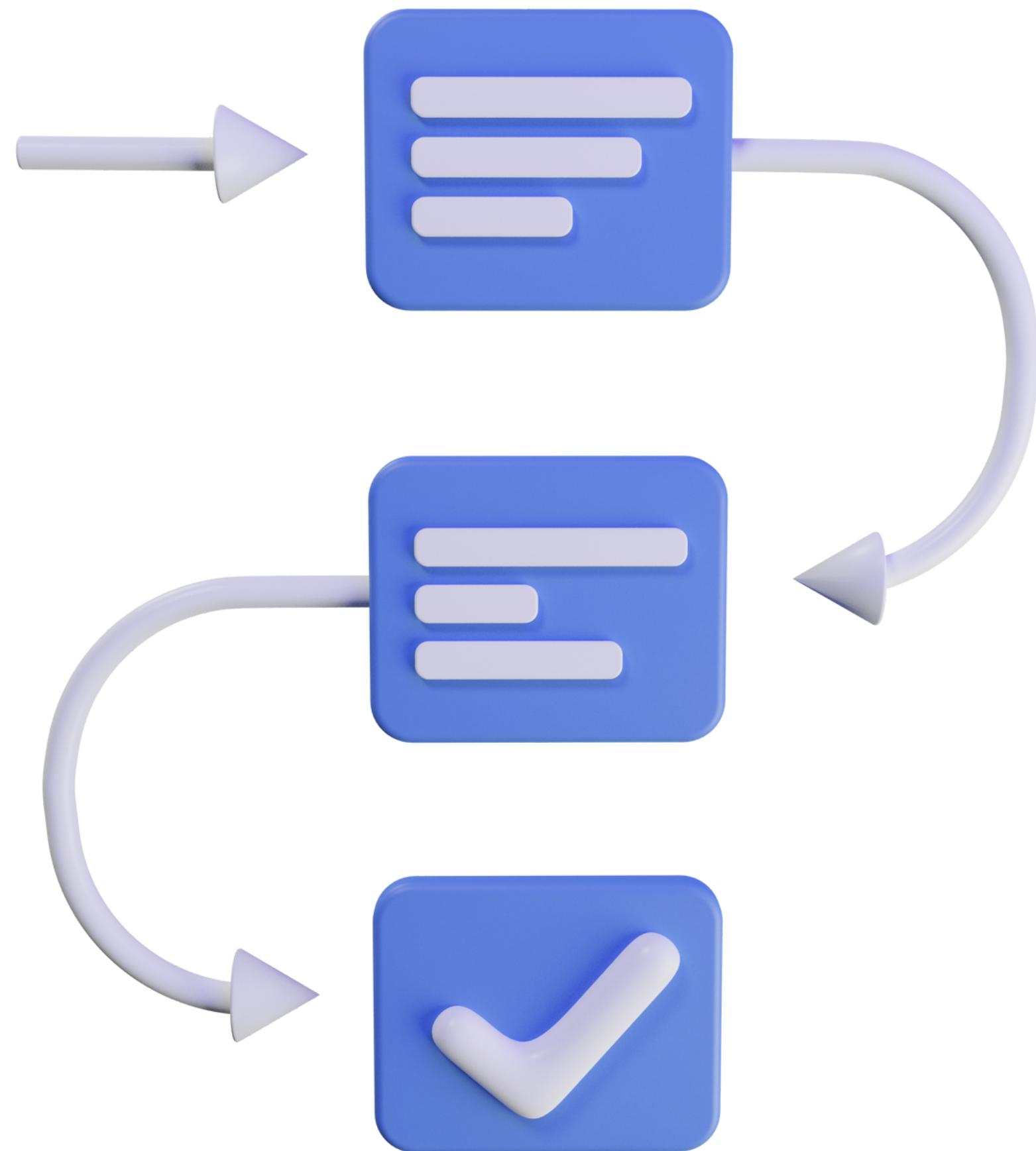
71%

**XGboost**

71%

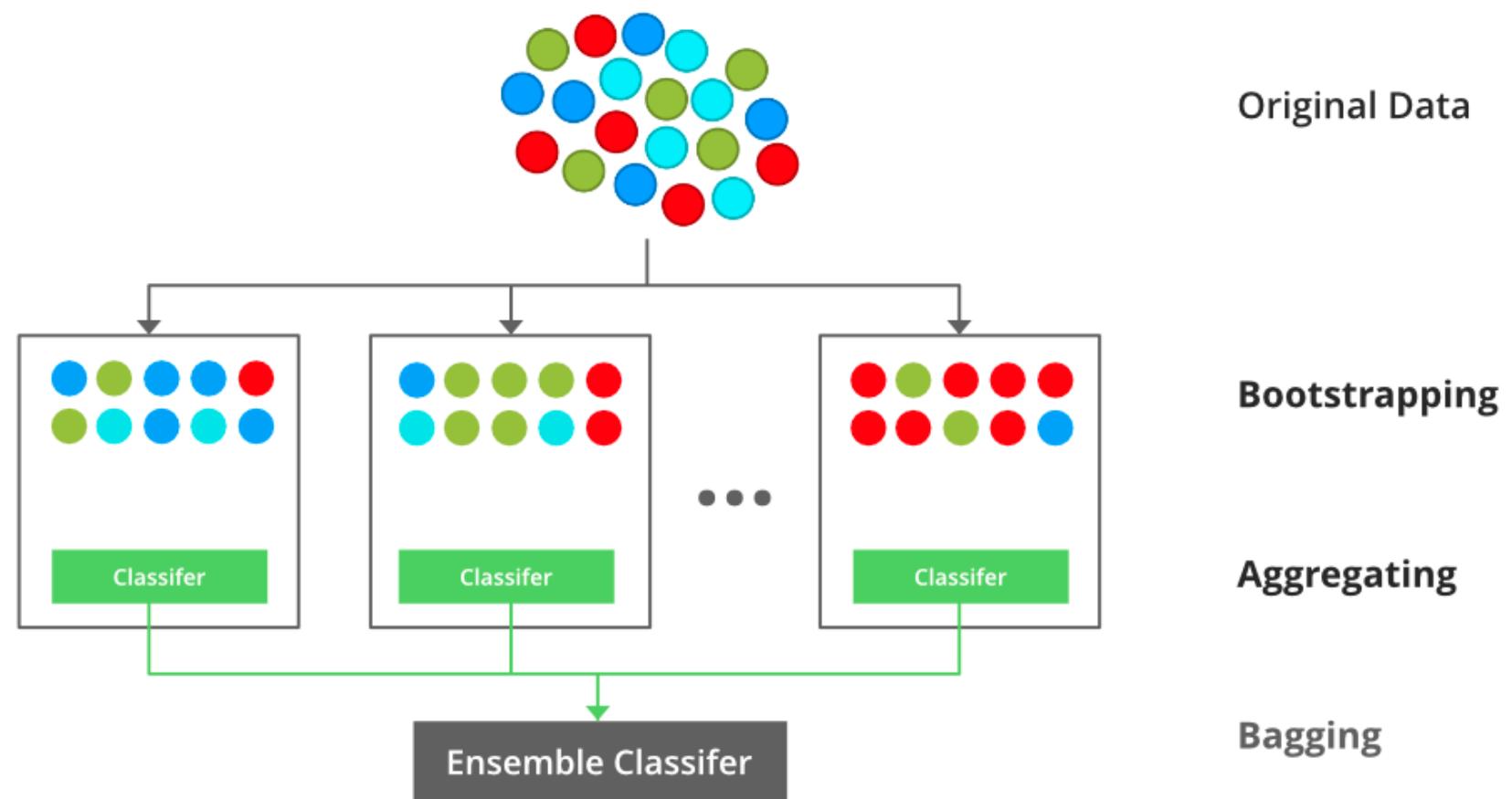
**XGboost in Pipeline**

71%



# XGBoost

XGBoost is the winner!



## Classification Report

Metric	Class 0	Class 1
Precision	0.73	0.69
Recall	0.77	0.64
F1-Score	0.75	0.66
Accuracy		0.71

# Key Drivers of Restaurant Ratings

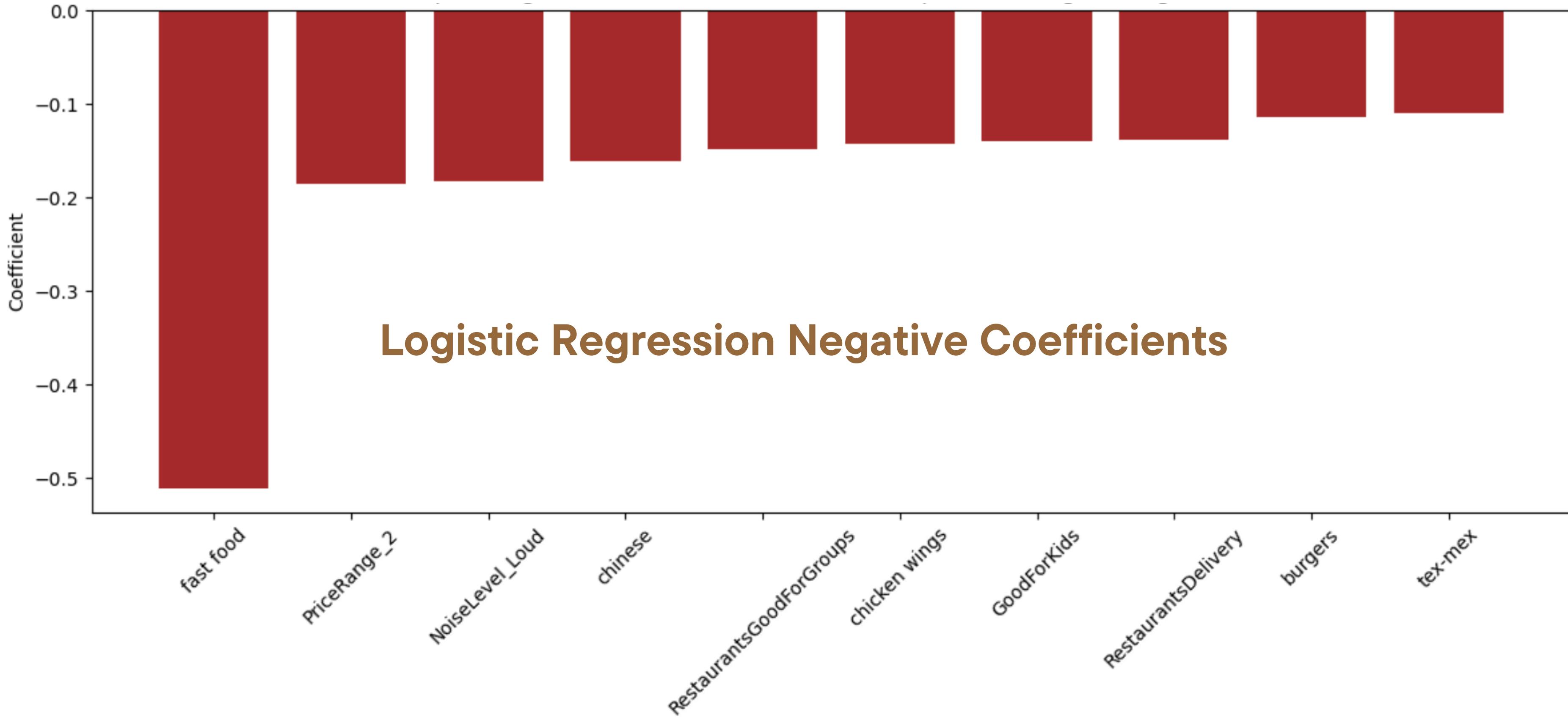


SHAP

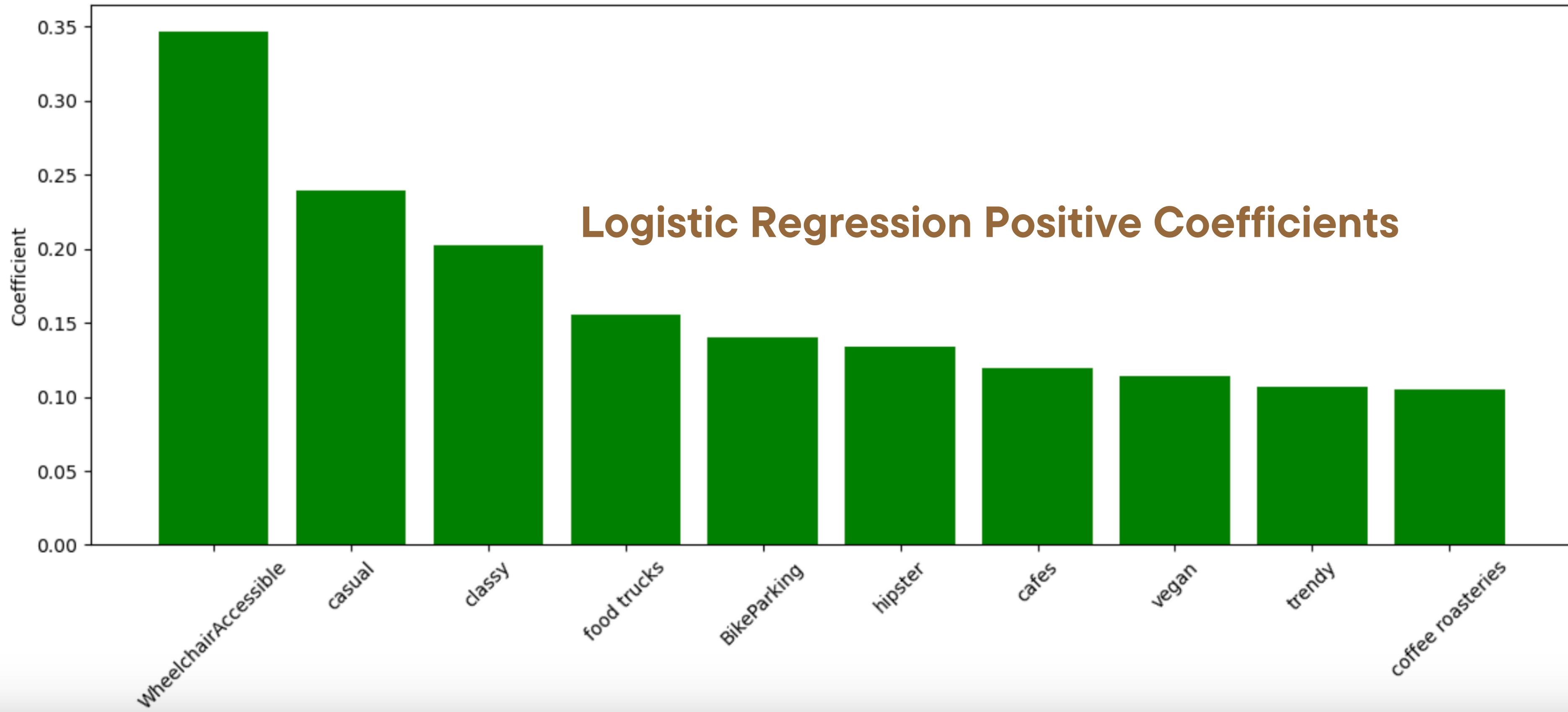
- **Logistic Regression Coefficients**
- **Random Forest Feature Importance**
- **SHAP Analysis with XGBoost**



# What Diners Might Not Prefer



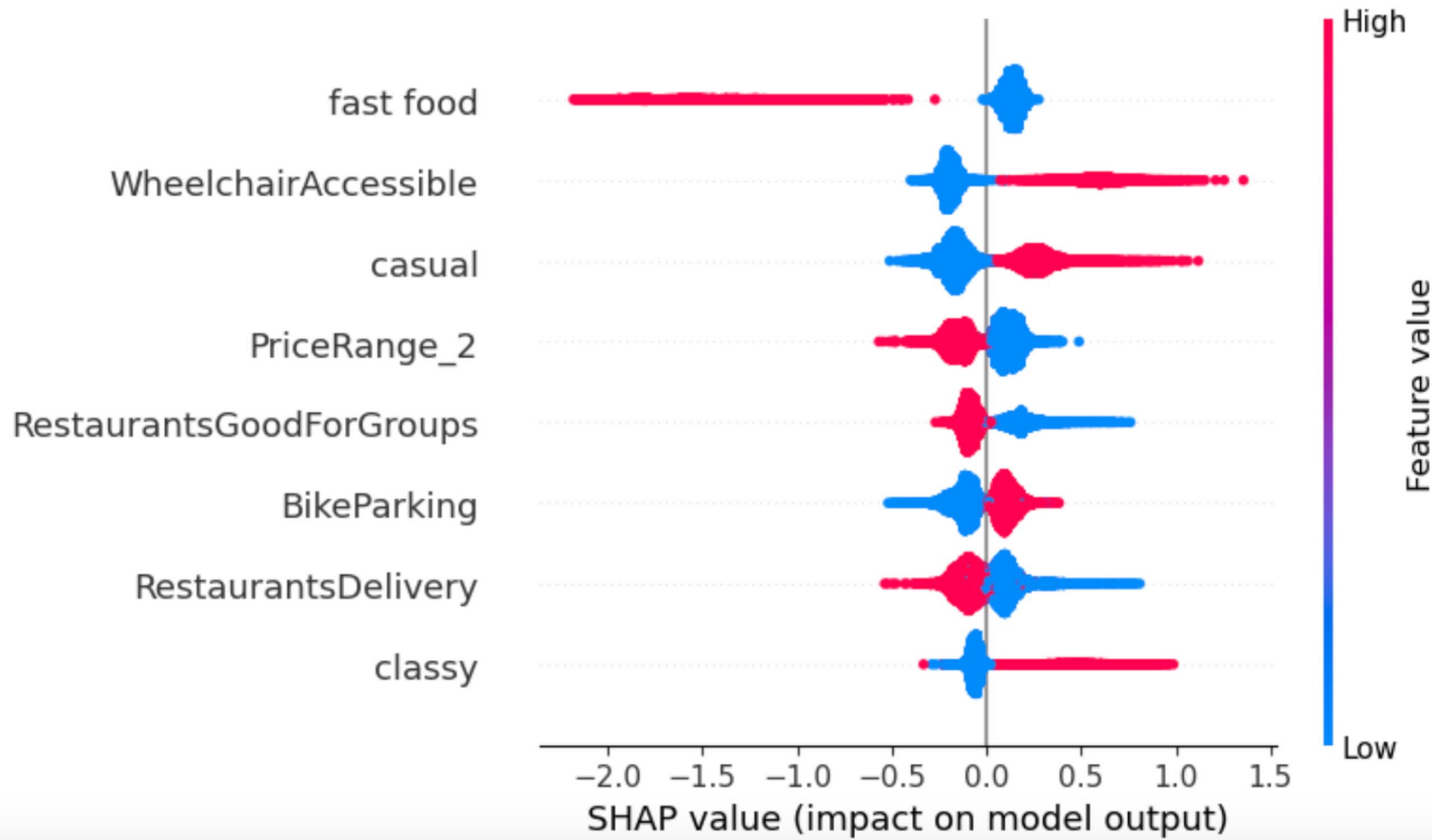
# Factors That Boost Restaurant Ratings





# SHAP Analysis with XGBoost

SHAP

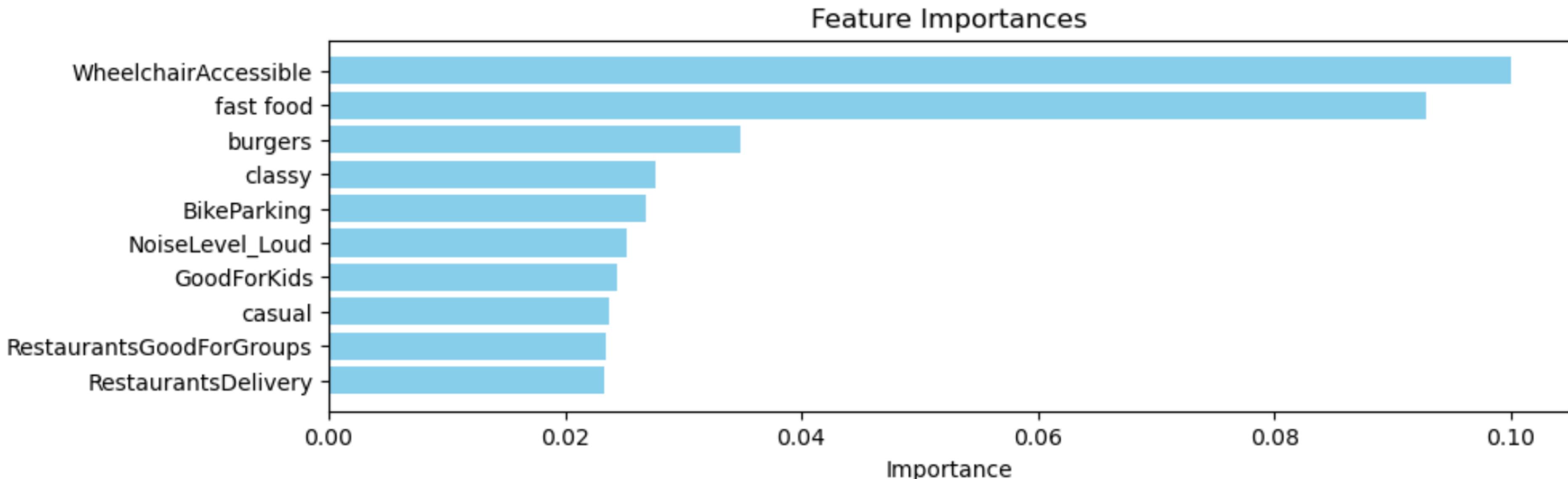




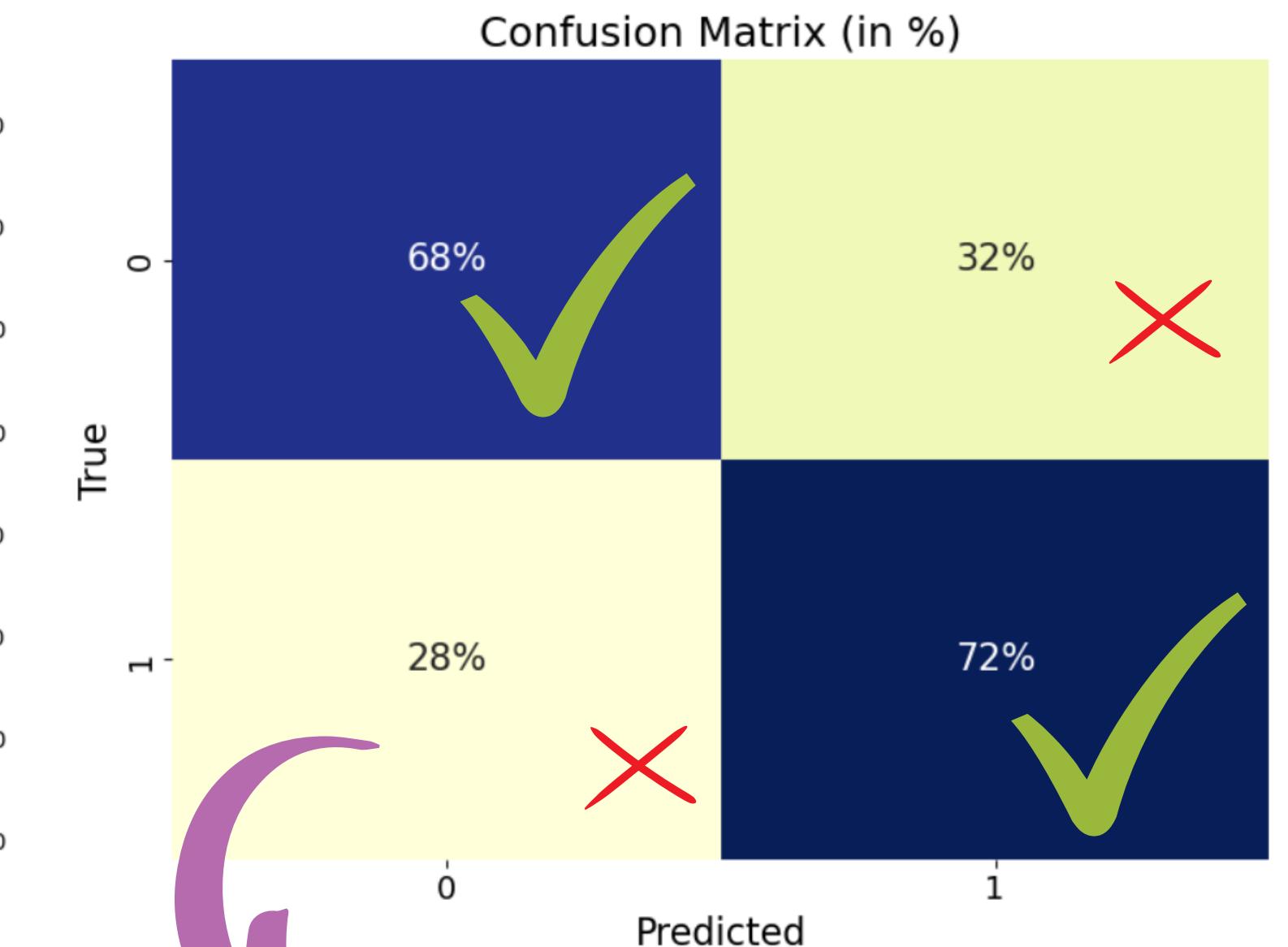
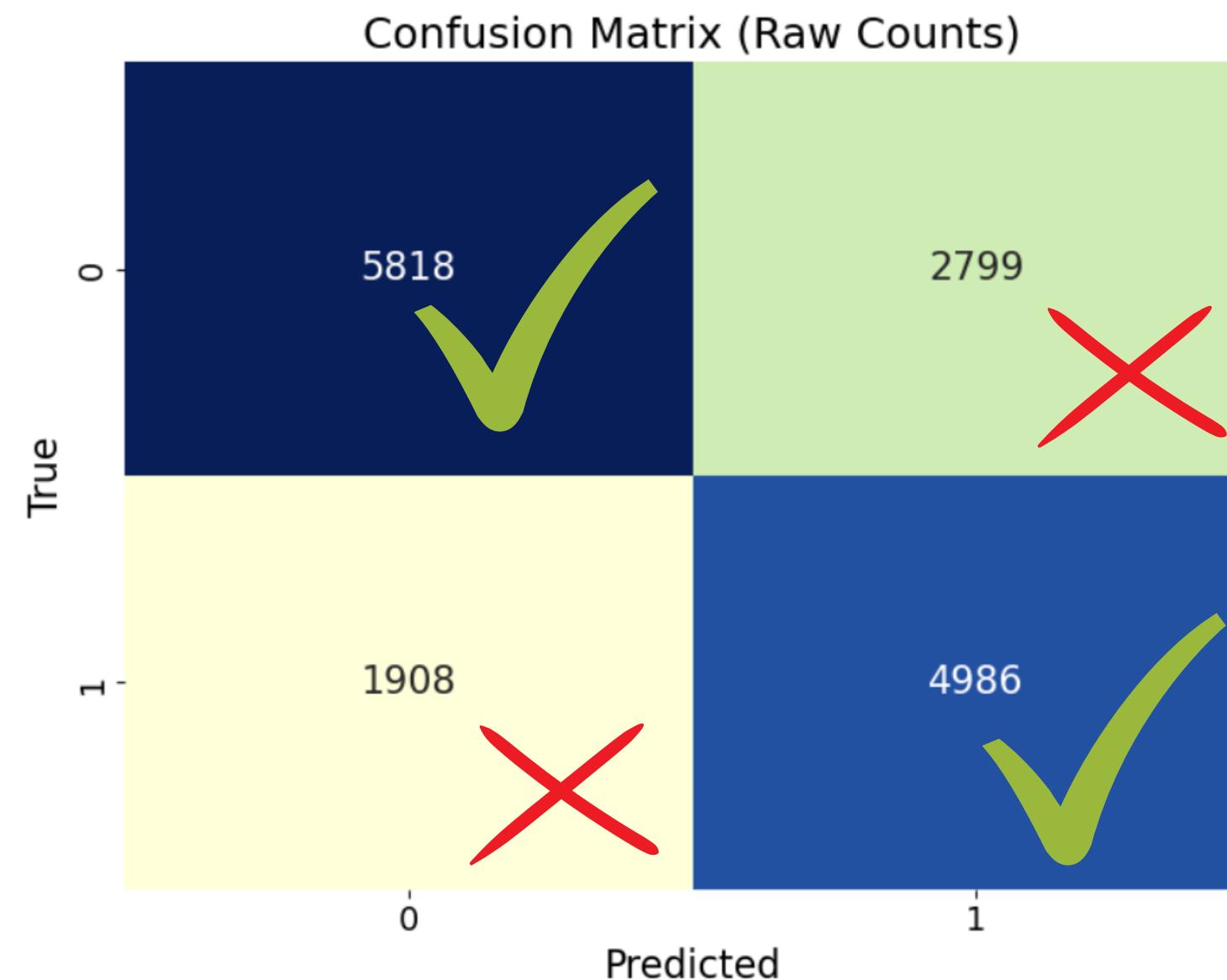
# Random Forest



## Random Forest Feature Importance



# Evaluating Our Model's Predictions

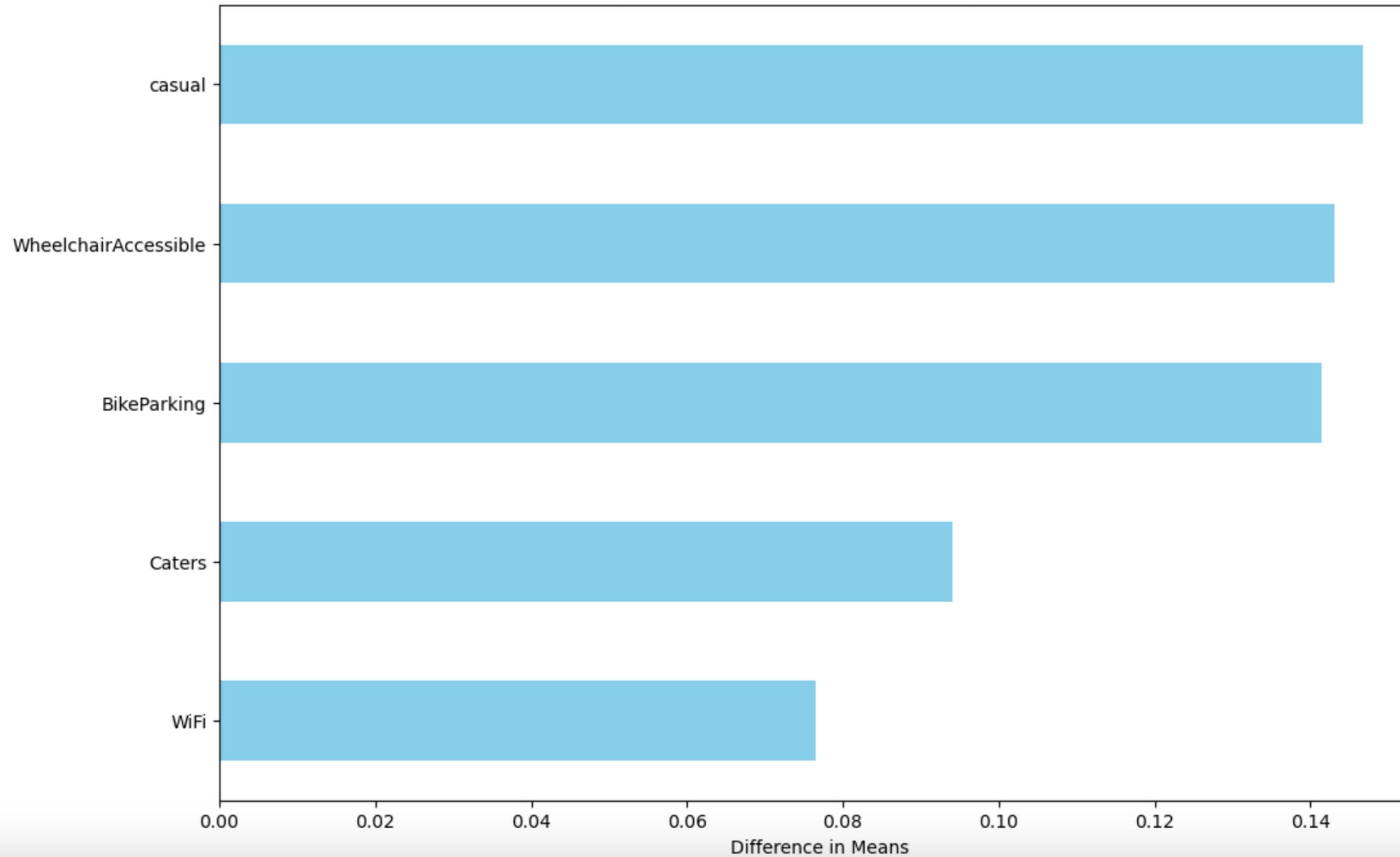


WHERE PREDICTIONS AND  
REALITY DON'T MATCH

# Model bias

## Feature Mean Differences Analysis

Top Features Influence on False Negatives



# Forging Ahead



## Potential Improvements

### Recommendations:

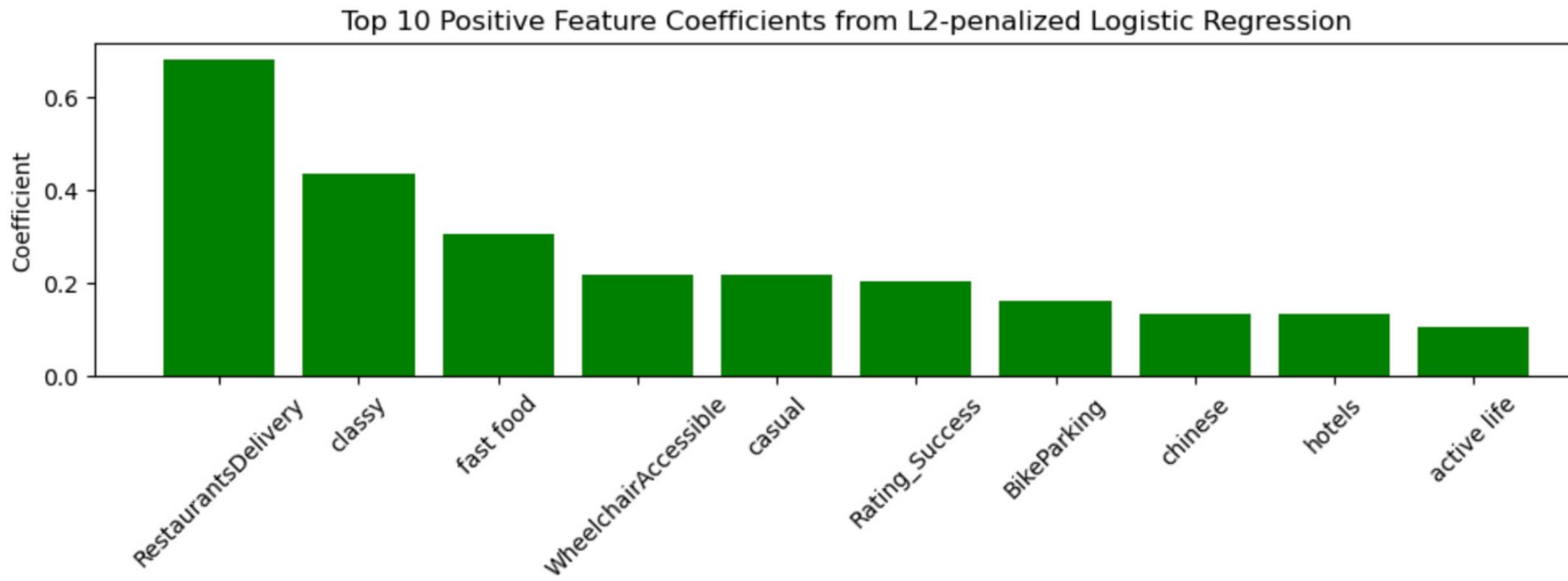
- a. **Collect More Data:** Especially on under-represented features.
- b. **Resampling:** Achieve balanced representation for wheelchair accessible venues.
- c. **Feature Weight Adjustment:** Refine importance of standout features for better prediction.

# Thank You

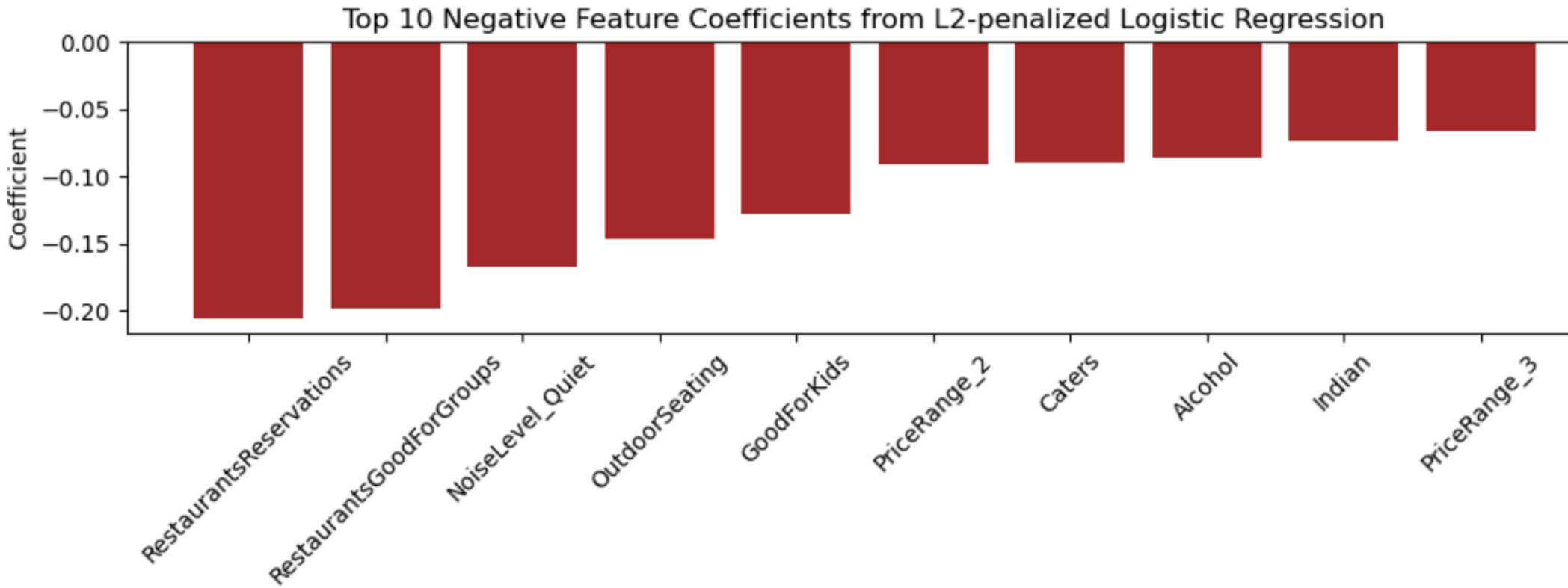


# Appendix

# Determinants of Closure



**Logistic Regression**    **Train Accuracy: 74.83%**  
**Model Performance:**    **Test Accuracy: 74.80%**



<b>Model</b>	<b>Parameters</b>	<b>Accuracy</b>	<b>Precision Low Rating</b>	<b>Recall Low Rating</b>	<b>F1-score Low Rating</b>	<b>Precision High Rating</b>	<b>Recall High Rating</b>	<b>F1-score High Rating</b>
<b>Logistic Regression</b>	<code>penalty='l2', solver='saga', C=100 ,test_size=0.30, random_state=42</code>	<b>0.70</b>	<b>0.71</b>	<b>0.77</b>	<b>0.74</b>	<b>0.68</b>	<b>0.60</b>	<b>0.64</b>
<b>Random Forest</b>	<code>bootstrap=False, max_depth=17, min_samples_leaf=2, min_samples_split=10, n_estimators=100</code>	<b>0.70</b>	<b>0.70</b>	<b>0.80</b>	<b>0.75</b>	<b>0.70</b>	<b>0.58</b>	<b>0.63</b>
<b>Random Forest Refined Model</b>	<b>feature Importance , Refined Model with Feature reduction</b>	<b>0.69</b>	<b>0.71</b>	<b>0.76</b>	<b>0.73</b>	<b>0.67</b>	<b>0.60</b>	<b>0.63</b>
<b>Gradient Boosting</b>	<code>n_estimators=250, random_state=42</code>	<b>0.70</b>	<b>0.71</b>	<b>0.78</b>	<b>0.74</b>	<b>0.68</b>	<b>0.61</b>	<b>0.64</b>
<b>PipeLine</b>	<b>feature extraction (PCA)</b>	<b>0.71</b>	<b>0.70</b>	<b>0.76</b>	<b>0.73</b>	<b>0.70</b>	<b>0.60</b>	<b>0.63</b>
<b>XGBoost</b>	<b>default parameters</b>	<b>0.71</b>	<b>0.73</b>	<b>0.77</b>	<b>0.75</b>	<b>0.69</b>	<b>0.64</b>	<b>0.67</b>
<b>PipeLine</b>	<code>classifier_colsample_bytree': 0.8, 'classifier_gamma': 0, 'classifier_learning_rate': 0.1, 'classifier_max_depth': 5, 'classifier_n_estimators': 250, 'classifier_subsample': 0.8</code>	<b>0.71</b>	<b>0.73</b>	<b>0.77</b>	<b>0.75</b>	<b>0.69</b>	<b>0.64</b>	<b>0.66</b>