# Restaurant Success and Revenue Prediction Model Analysis

Marcus Lui

# What is a measurement of success?

Google Rating?

Social Media Followers?

Annual Revenue?

Years in Business?
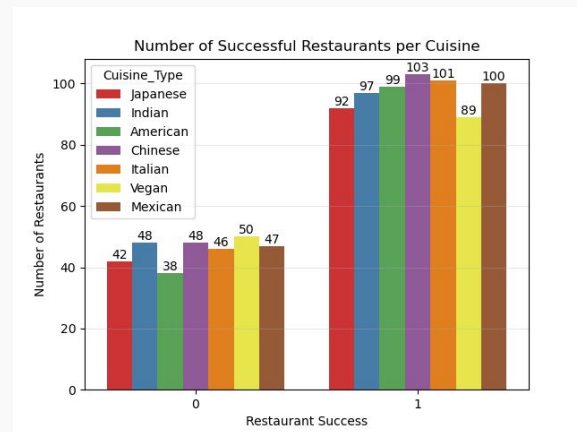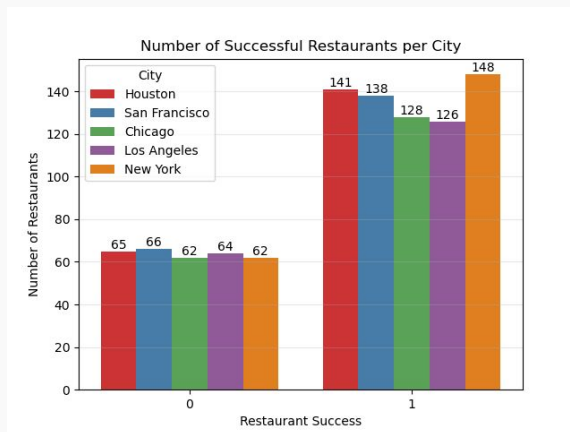
Average Meal Price?

Seating Capacity?

The goal of this project is to identify the best models for restaurant success classification and for revenue prediction

... and to provide information that may help restaurants improve their recipe for success
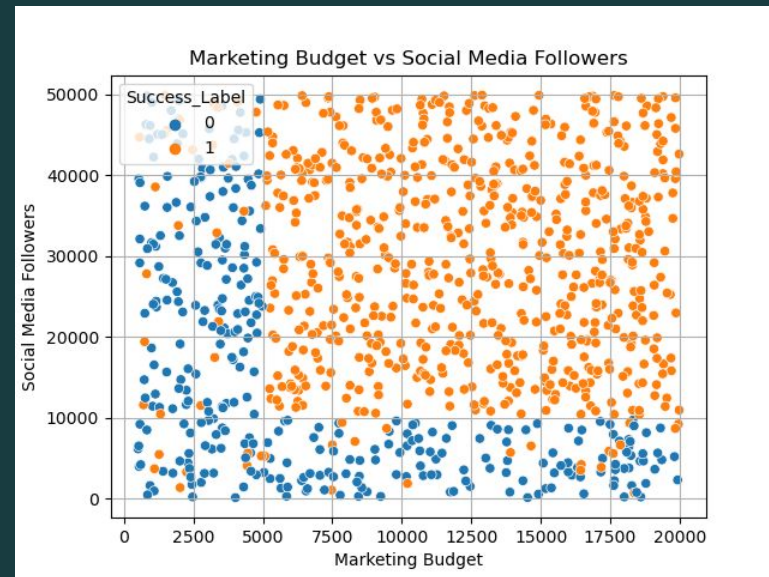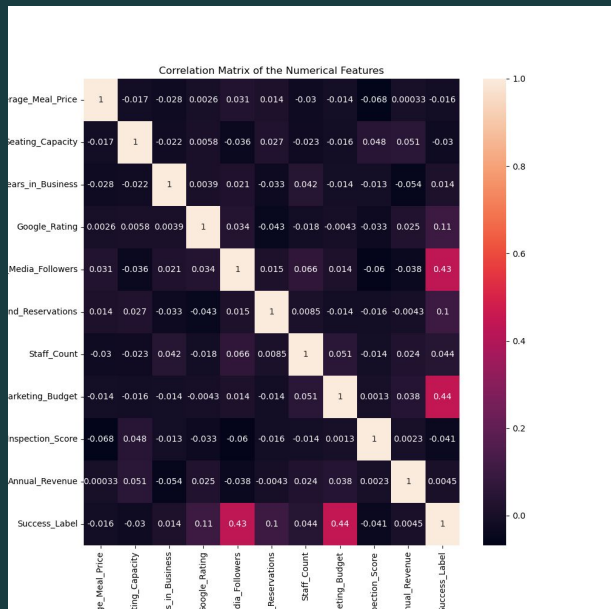
# The *Data*

- 4 Categorical Features
  - City
  - Cuisine Type
  - Delivery Service
  - Success Label
- 9 Numerical Features
  - Average Meal Price
  - Seating Capacity
  - Years in Business
  - Google Rating
  - Social Media Followers
  - Weekend Reservations
  - Staff Count
  - Marketing Budget
  - Health Inspection Score
  - Annual Revenue

# The First Look



Number of Successful Restaurants per City



Number of Successful Restaurants per Cuisine

# A Deeper Look...





The correlation heatmap (left) shows the features that shows some direct effect on restaurant success and is then confirmed via the scatterplot (above).

# Model Types

Classification Models
- Dummy
- Logistic Regression
- K-Nearest Neighbors
- Decision Tree
- Random Forest

Regression Models
- SVM
- Linear Regression
- Lasso
- Random Forest

# Classification Models

### Logistic Regression (LogReg)

A statistical model used for binary classification

### K-Nearest Neighbors (KNN)

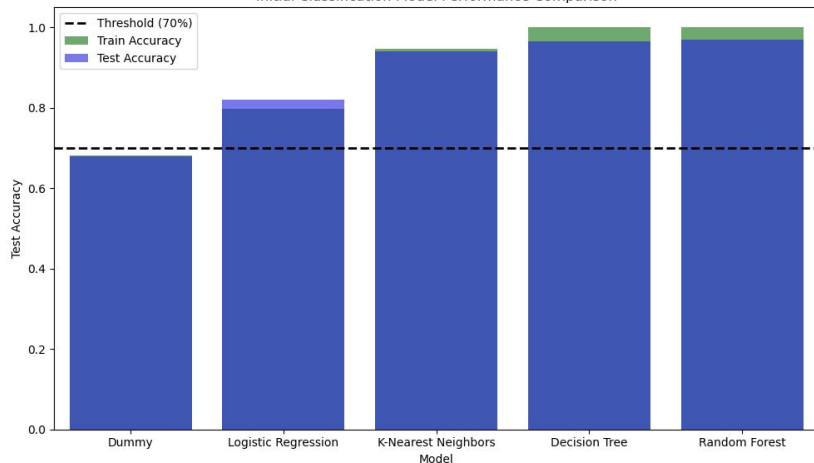A simple, non-parametric algorithm used for classification

### Decision Tree (DT)

An algorithm that recursively splits data into branches based on the features that provide the best separation until all nodes belong to a class

### Random Forest (RF)

An ensemble learning method that builds multiple decision trees and combines their outputs to improve accuracy and prevent overfitting
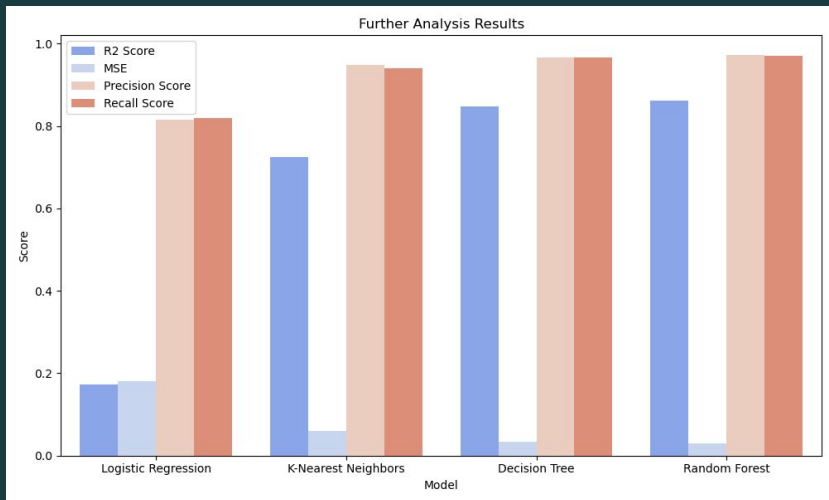
# Classification Model Analysis



Initial Classification Model Performance Comparison

LogReg and KNN models performed better than expected, recording a 82% and 94% test accuracy, respectively.

GridSearchCV was then applied to determine the best hyperparameters per each model. Results from the enhanced models were similar to the initial model performances.

Decision Tree and Random Forest models both recorded a 100% training accuracy, while both testing accuracies regressed by ~3%

Further Analysis Results

KNN, DT, RF all had a precision score and recall score above 90%. These two models also had a low MSE and a high R2 Score.

For this dataset and this task, DT and RF seem like good models to use for predicting success.

## Precision Score

It answers the question:

"Of all the positive predictions made by the model, how many were actually correct?

## MSE

It measures the average squared difference between the actual (true) values and the predicted values. Low MSE means higher prediction accuracy.

## R2 Score

It is a measurement of how well the models is at explaining variability

# Regression Models

### Linear Regression (LinReg)

A fundamental machine learning algorithm used to model the relationship between a dependent variable and one or more independent variables

### Lasso

A Linear Regression technique that should help with feature selection and thus reduce overfitting
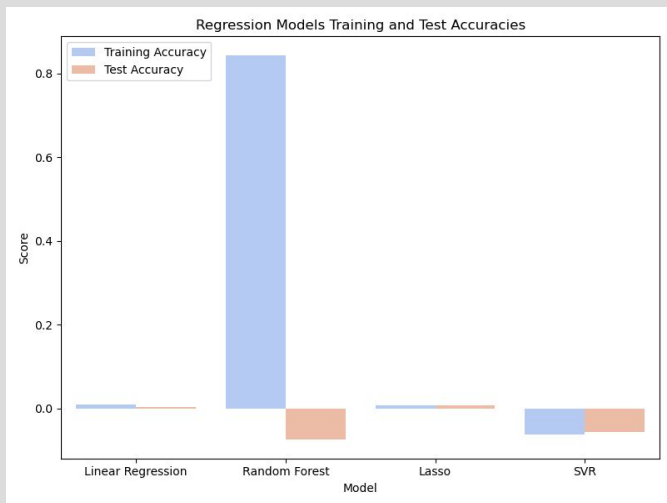
### Support Vector Regression (SVR)

Predicts continuous values by finding a function that fits the data within a certain margin of tolerance, ignoring outliers

### Random Forest (RF)

An ensemble machine learning method that uses multiple decision trees to predict a continuous output
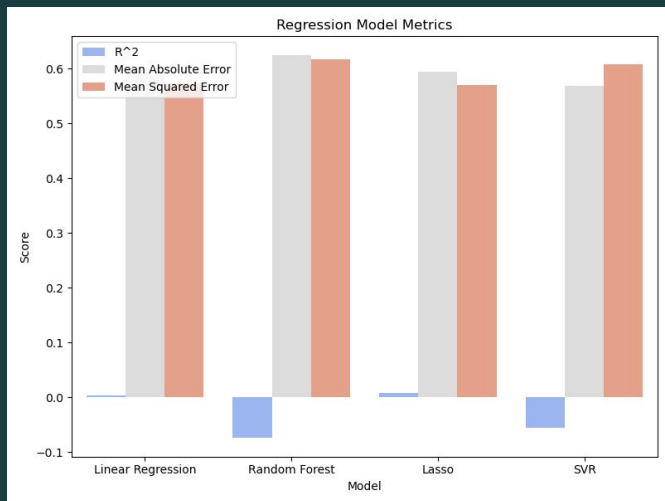
# Predictive Model Analysis



Regression Models Training and Test Accuracies

Unlike the results of the Classification models, the Linear Regression (LinReg), Random Forest (RFR), Lasso, and SVR models all had a low testing accuracy.

GridSearchCV was then applied to determine the best hyperparameters per each model. Results from the enhanced models were no different to the initial regression model performances.

Such drastic results suggests the data is too scattered, and an unbiased model that ensures high accuracy and high precision cannot be created

Regression Model Metrics

LinReg, RFR, Lasso, and SVR all seemed to have a high MAE and MSE score while also having a low R2 score.

For this dataset and this task, there doesn't seem to be any models that have a high accuracy to use for predicting annual revenue.

## MAE

It is the average of the absolute differences between the predicted and actual values. Lower MAE means higher prediction accuracy.

## MSE

It measures the average squared difference between the actual (true) values and the predicted values. Low MSE means higher prediction accuracy.

## R2 Score

It is a measurement of how well the models is at explaining variability

# Model Results Summary

With DT and RF, you can accurately classify a restaurant as successful or not successful with over 95% accuracy.

The results of the DT and RF models have little variance and low error and thus is proven to be both precise and accurate

RFR had the best training accuracy, but recorded a surprisingly low test accuracy, which suggests that the models may be over biased.

All regression models tested performed poorly on this data, implying that the features available do not capture the complexity of what drives annual revenue

# Key Takeaways



Google Rating vs Annual Revenue
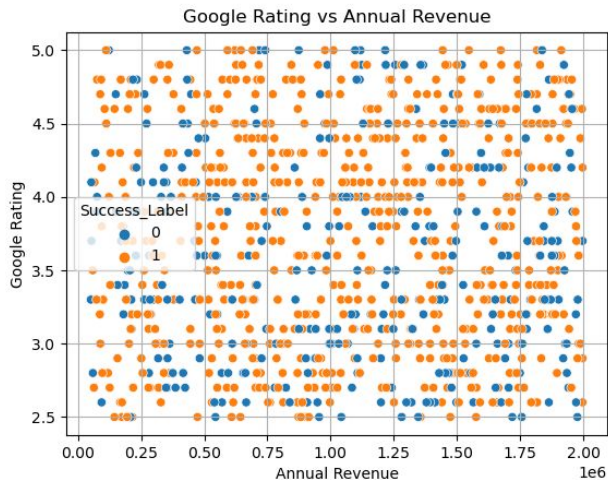
## Use Social Media and Boost Marketing

Restaurants that prioritized a higher marketing budget and a greater online presence were recognized as successful.

## Revenue Prediction Requires More Quality Data

The dataset used had both categorical data and numerical data, but a limited amount. Good data is better than noise

## New Features and Metrics

Features that could improve the dataset include customer loyalty data, seasonal trends, or competition density

# Thank you

https://github.com/MarcusLui9/Restaurant_Success_and_Revenue_Prediction
https://www.linkedin.com/in/marcus-chi-kin-lui/