# Restaurant Location Analysis in San Francisco using Machine Learning Techniques

### 1. Executive Summary

This project utilizes machine learning techniques to identify the ideal restaurant location in San Francisco for aspiring restaurateurs and real estate developers. By analyzing the web scraped data from Yelp and Uber Eats, I evaluated existing restaurant performance, demographics, and consumer preferences, ensuring thorough data cleaning and preprocessing for consistency. The endeavor extended to web scraping for detailed mapping of San Francisco's neighborhoods and their zip codes, enhancing our dataset to cover broader areas within the Bay Area and refining our geographic understanding to pinpoint factors influencing restaurant success.

Through feature engineering, I grouped varied cuisine types into broader categories, followed by training a random forest model on the dataset, now segmented into training and test sets. This approach allowed to understand the impact of factors such as cuisine categories, neighborhood locations, and popularity metrics on the model's predictive accuracy. The comprehensive use of data analytics and machine learning equips clients with essential insights, empowering them to make strategic decisions that greatly enhance the potential for success in San Francisco's competitive restaurant landscape.

### 2. Background, Context, and Domain Knowledge

This project is strategically positioned at the intersection of real estate and the restaurant industry,

designed to assist both real estate developers and restaurateurs aiming to penetrate San Francisco's competitive dining scene. This initiative reflects a fusion of data-driven analytics and industry acumen, aimed at identifying the optimal location for a new restaurant venture. By focusing on San Francisco's culinary environment, the project underscores the importance of selecting a site that not only aligns with the desired culinary theme but also integrates with the intended physical storefront, enhancing the establishment's prospects for success.

The methodology behind this venture involves a comprehensive analysis of the local restaurant ecosystem, with a spotlight on consumer behavior, competitive dynamics, and neighborhoods with high levels of customer engagement. The analysis leverages data from a variety of sources, including Yelp and Uber Eats for insights into market trends and restaurant performance, complemented by detailed web scraping of the mappings of neighborhood demographics through zip codes.

The utilization of advanced machine learning techniques, particularly the Random Forest algorithm facilitates the precise modeling of factors crucial to the success of a restaurant venture, enhancing the richness of our dataset and providing our clients with clear, actionable insights. In doing so, it is aimed to equip real estate developers and restaurateurs with the knowledge necessary to navigate San Francisco's competitive restaurant landscape successfully, fostering the foundation for a prosperous business venture.

### 3. Traditional Approaches and Strategic Alignment in Restaurant Market Entry

Traditionally, the restaurant industry has approached the challenge of market entry through a

combination of market research, personal intuition, and trial-and-error. This conventional strategy involves analyzing basic demographic information, scouting potential locations based on foot traffic and visibility, and relying heavily on the restaurateur's personal experience and understanding of the culinary landscape. The decision-making process is often influenced by factors such as rental costs, perceived demand for certain cuisines, and the presence of complementary businesses that might drive customer traffic.

Historically, this approach has yielded mixed results, with success dependent on the restaurateur's or the developer's acumen and often subject to unforeseen market dynamics. The traditional model emphasizes location as a critical factor for success, aligning with the adage that the location matters the most. While this strategy can lead to successful outcomes, it also poses significant risks due to its reliance on subjective judgment and limited quantitative analysis. Moreover, it may overlook deeper market insights such as the impact of local competition, changing consumer preferences, and demographic shifts that are critical in today's rapidly evolving culinary scene.

The integration of data analytics and machine learning has transformed traditional market entry strategies. By utilizing diverse data sources, including social media, online reviews, and economic indicators, businesses can align their decisions more closely with their business model and strategic objectives. This approach gives a deeper insight into consumer behavior, the competitive landscape, and market trends, offering a competitive edge in choosing locations that appeal to the desired demographic and address market gaps. Moving from intuition-based to evidence-driven strategies through advanced analytics will improve the chances of success in the competitive restaurant sector.

**4. Analyses**

**4.1 Logistic Regression**

In the logistic regression model, I aimed to predict the popularity category of restaurants based on various features such as rating, price, category group, and neighbourhood. The model achieved an accuracy of 79%, indicating that it correctly classified approximately 79% of the instances in the test dataset.

The classification report provides further insights into the model's performance across different classes. For the "low" popularity category, the precision, recall, and F1-score were high, indicating that the model effectively identified instances of this class with high accuracy and captured a high proportion of true positives. However, for the "high" popularity category, the precision, recall, and F1-score were low, indicating that the model struggled to accurately identify instances of this class.

Overall, while the logistic regression model demonstrated reasonable performance, there is room for improvement, particularly in accurately predicting instances of the "high" popularity category. Further refinement of the model and exploration of additional features may help enhance its predictive capability.

- Accuracy: 79%
- Precision (weighted avg): 73%
- Recall (weighted avg): 79%

- F1-score (weighted avg): 75%

## 4.2 Decision Tree Classification

In the decision tree model, the objective was to predict the popularity category of restaurants based on features such as rating, price, category group, and neighbourhood. The model achieved an accuracy of 75%, indicating that it correctly classified approximately 75% of the instances in the test dataset.

The classification report provides a breakdown of the model's performance across different popularity categories. For the "low" popularity category, the precision, recall, and F1-score were high, suggesting that the model effectively identified instances of this class with high accuracy and captured a high proportion of true positives. However, for the "high" popularity category, the precision, recall, and F1-score were relatively low, indicating challenges in accurately identifying instances of this class.

Overall, the decision tree model demonstrated reasonable performance but exhibited limitations, particularly in accurately predicting instances of the "high" popularity category. Further optimization of the model parameters and feature selection may be necessary to improve its predictive accuracy.

- Accuracy: 75%
- Precision (weighted avg): 76%
- Recall (weighted avg): 75%

- F1-score (weighted avg): 75%

### 4.3 Random Forest Classification

In the Random Forest model, I aimed to predict the popularity category of restaurants using a more complex ensemble learning technique. This model incorporated decision trees to make predictions based on various input features such as rating, price, category group, and neighbourhood. The Random Forest model achieved an accuracy of 74%, indicating that it correctly classified approximately 74% of the instances in the test dataset.

The classification report provides a detailed breakdown of the model's performance across different popularity categories. For the "low" popularity category, the precision, recall, and F1-score were relatively high, indicating that the model effectively identified instances of this class with high accuracy and captured a high proportion of true positives. However, for the "high" popularity category, the precision, recall, and F1-score were low, indicating that the model struggled to accurately identify instances of this class.

Compared to the logistic regression model, the Random Forest model exhibited similar performance overall, with both models achieving comparable accuracy scores. However, like the logistic regression model, there is room for improvement in accurately predicting instances of the "high" popularity category. Further refinement of the model, parameter tuning, and feature engineering may help enhance its predictive capability.

- Accuracy: 74%

- Precision (weighted avg): 71%

- Recall (weighted avg): 74%

- F1-score (weighted avg): 73%

### 4.4 Random Forest Classification with Hyperparameter Tuning

In the Random Forest model with hyperparameter tuning, I aimed to predict the popularity category of restaurants using an ensemble of decision trees. The model achieved an accuracy of 80%, indicating that it correctly classified approximately 80% of the instances in the test dataset.

The classification report provides a detailed assessment of the model's performance across different popularity categories. For the "low" popularity category, the precision, recall, and F1-score were high, indicating that the model effectively identified instances of this class with high accuracy and captured a high proportion of true positives. Similarly, for the "new" popularity category, the precision, recall, and F1-score were also high, demonstrating the model's ability to accurately classify instances in this category.

However, for the "high" popularity category, the precision, recall, and F1-score were relatively low compared to the other categories, suggesting that the model struggled to accurately identify instances of this class. This indicates a potential area for improvement in the model's predictive performance.

Overall, the Random Forest model with hyperparameter tuning showed promising results,

achieving a higher accuracy compared to the logistic regression model. However, further refinement and exploration may be necessary to improve its performance, particularly in accurately predicting instances of the "high" popularity category.

- Accuracy: 80%

- Precision (weighted avg): 76%

- Recall (weighted avg): 80%

- F1-score (weighted avg): 78%

## 5. Recommendations and Business Value
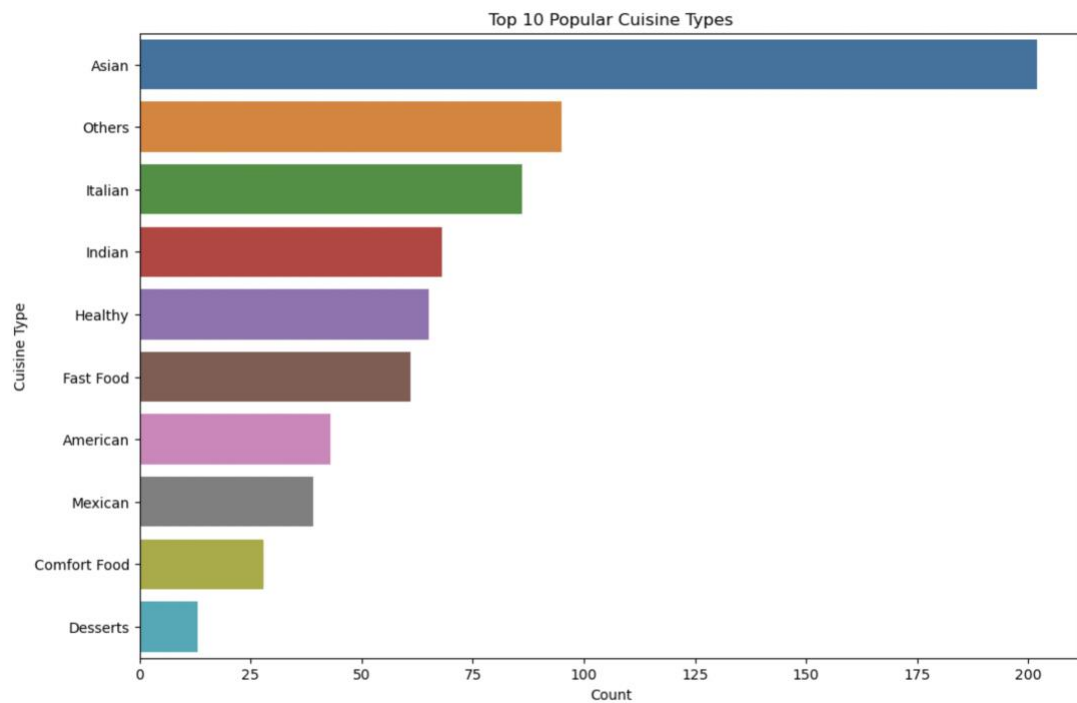
Recommendations:

The Random Forest classifier with hyperparameter tuning has demonstrated high accuracy in predicting restaurant popularity in San Francisco based on price, cuisine category group, and neighborhood. With a commendable accuracy of 80%, I would suggest the following strategic recommendations for stakeholders:

1) Utilize the model's predictions to make informed investment decisions. The model can help pinpoint restaurant profiles that are likely to be popular thus serving as a guide for potential restaurant owners or investors in choosing where to allocate their resources.

2) Direct marketing strategies toward enhancing restaurant features that align with the characteristics associated with high popularity scores. Adjustments in price levels or emphasizing specific cuisine categories (such as Asian cuisine in SF), can help in targeting the right customer demographics.

3) Leverage insights from the model for expansion planning. Restaurants looking to open

new locations should consider neighborhoods and restaurant features that align with higher popularity ratings to maximize the success rate.

Business Value

1) The data-driven approach provides a strategic advantage by identifying high-potential restaurant attributes, allowing businesses to outperform competitors.

2) The insights from the machine learning model can aid in identifying and mitigating potential risks, thereby making more secure investment choices.

3) In this analysis, I learnt that Asian cuisine is the most popular followed by Indian, Italian, Healthy and Others. Incorporating dishes from these cuisines could prove to be beneficial for a restaurant business.



4) I also got to know how the restaurants are distributed across different neighborhoods.

This can help a restaurant business come up strategies to choose the right location.

5) I also found out which cuisines have the most presence in neighborhoods. This information will help us understand the customer preferences of each region and can be particularly helpful in choosing the restaurant cuisine based on the area.

**6. Summary and Conclusions**

In conclusion, this data-driven investigation into the San Francisco restaurant scene has yielded a robust predictive model with substantial accuracy. The use of a Random Forest classifier with hyperparameter tuning has proven effective in assessing the popularity of restaurants, drawing on key factors such as price point, cuisine type, and neighborhood. This model not only underscores the potential of machine learning in strategic business applications but also offers a compelling tool for stakeholders looking to understand and capitalize on the dynamics of the restaurant industry. With the analytical efforts culminating into a reliable model, this project paves the way for informed decision-making that could enhance the success rate of new and existing restaurants in the competitive San Francisco market.