

California Restaurant "Likes" Prediction Using Foursquare API and Machine Learning

By : Khalid.Meguenni

Capstone Project

IBM Data Science Professional Certificate

1. Introduction

California boasts an incredibly diverse collection of restaurants catering to different palettes and appetites. A large part of marketing for a modern restaurant (or any company) is social media, where the number of "likes" that the company can receive will dictate its brand and image to the general public.

For a new business owner (or existing company) to open a new restaurant in California, knowing ahead of time the potential social media image they can have would provide an excellent solution to the ever present business problem of uncertainty. In this case the uncertainty is regarding performance of social media presence.

We can mitigate this uncertainty through leveraging data gathered from FourSquare's API, specifically, we are able to scrape "likes" data of different restaurants directly from the API as well as their location and category of cuisine. The question we will try to address is, how accurately can we predict the amount of "likes" a new restaurant opening in this region can expect to have based on the type of cuisine it will serve and which city in California it will open in. (For the purposes of this analysis, we will contain the geographical scope of analysis to three heavily populated cities in California, namely San Francisco, Los Angeles, and San Diego).

Leveraging this data will solve the problem as it allows the new business owner (or existing company) to make preemptive business decisions regarding opening the restaurant in terms of whether it is feasible to open one in this region and expect good social media presence, what type of cuisine and which city of three would be the best. This project will analyze and model the data via machine learning through comparing both linear and logistic regressions to see which method will yield better predictive capabilities after training and testing.

2. Data

2.1 Data Scraping and Cleaning

In this section we will first retrieve the geographical coordinates of the three cities (San Francisco, Los Angeles, and San Diego). Then, we will leverage the FourSquare API to obtain URLs that lead to the raw data in JSON form. We will separately scrape the raw data in these URLs in order to retrieve the following columns: "name", "categories", "latitude", "longitude". and "id" for each city. We can also provide another column ("city") to indicate which city the restaurants are from.

It is important to note that the extracts are not of every restaurant in those cities but rather all of the restaurants within a 1000KM range of the geographical coordinates that geolocator was able to provide. However, the extraction from the FourSquare API actually obtains venue data so it will include venues other than restaurants such as concert halls, stores, libraries etc. As such, this means that the data will need to be further cleaned somewhat manually by removing all of the non-restaurant rows. Once this is complete, we have a shortened by cleaned list to pull "likes" data. The reason the cleaning takes precedence is mainly that pulling the "likes" data is the computing process which takes the longest time in this project so we want to make sure we are not pulling information that will end up being dropped anyways.

The "id" is an important column as it will allow us to further pull the "likes" from the API. We can retrieve the "likes" based on the restaurant "id" and then append it to the data frame. Once this is complete, we finally name the dataframe 'raw_dataset' as it is the most complete compiled form before needing any processing for analysis via machine learning.

2.2 Data Preparation

The data still needs some more processing before it is suitable for model training and testing. Mainly, the "categories" column contains too many different types of cuisines to allow a model to yield any meaningful results. However, the different types of natural cuisines have natural groupings based on conventionally accepted cultural groupings of cuisine. Broadly speaking, all of the different types of cuisine could be reclassified as European, Latin American, Asian, North American, drinking establishments (bars), or casual establishments such as coffee shops or ice cream parlours. We can implement manual classification as there really aren't that many different types of cuisines.

As this project will compare both linear and logistic regression, it makes sense to have "likes" as both a continuous and categorical (but ordinal) variable. In the case of turning into a categorical variable, we can bin the data based on percentiles and classify them into these ordinal percentile categories. I tried different ways of binning but in the end, splitting the sample into three different bins proved to yield the best classification results from a prediction standpoint.

As the last stage of data preparation, it is important to note that the regressors are categorical variables (3 different cities and 6 different categories of cuisines). Hence, they require dummy variable encoding for meaningful analysis. We can accomplish this via one-hot encoding.

3. Methodology

This project will utilize both linear and logistic regression machine learning methods to train and test the data. Namely, linear regression will be used in an attempt to predict the number of "likes" a new restaurant in this region will have. We will utilize the Sci-Kit Learn Package to run the model.

We can also utilize logistic regression as a classification method rather than direct prediction of the number of likes. Since the number of "likes" can be binned into different categories based on different percentile bins, it is also potentially possible to see which range of "likes" a new restaurant in this region will have.

Since the "likes" are binned into multiple (more than 2) categories, the type of logistic regression will be multinomial. Additionally, although the ranges are indeed discrete categories, they are also ordinary in nature. Therefore, the logistic regression will need to be specified as being both multinomial and ordinal. This can be done through the Sci-Kit Learn Package as well.

4. Results

4.1 Linear Regression Results

A linear regression model was trained on a random subsample of 80% of the sample and then tested on the other 20%. To see if this is a reasonable model, the residual sum of squares score and variance score were both calculated (56861.68 and 0.08 respectively). Given the low variance score, this is probably not a valid/good way of modelling the data. Therefore, we move on to logistic regression.

4.2 Logistic Regression Results

A multinomial ordinal logistic regression model was trained on a random subsample of 80% of the sample and then tested on the other 20%. To see if this is a reasonable model, its jaccard similarity score and log-loss were calculated (44% and 1.18 respectively). Although this is not a perfect prediction, a similarity of 44% between the training set and test set is a reasonable result. The classification report is also printed later on below.

Given the modestly accurate ability of this model, we can also run the model on the full dataset. The coefficients show that opening a restaurant in San Francisco, opening a bar, or serving cuisine that is American or Asian in nature, are associated negatively with "likes."

5. Discussion

The first thing to note is that given the data, logistic regression presents a better fit for the data over linear regression. Using logistic regression we were able to obtain a Jaccard Similarity Score of 44%, which although not perfect, is more reasonable than the low variance score obtained from the linear regression. As stated before, please note that for the purposes of this project, we are assuming that likes are a good proxy for how well a new restaurant will do in terms of brand, image and by extension how well the restaurant will perform business-wise. Whether or not these assumptions hold up in a real-life scenario is up for discussion, but this project does contain limitations in scope due to the amount of data that can be fetched from the FourSquare API.

As such, to obtain insights into this data, we can proceed with breaking down the results of the logistic regression model. The results showed that the precision score for classifying whether the new restaurant would fall into classes 1, 2, were 44% , 40%,. Therefore, the model is better at predicting if a restaurant will fall into the best or worst percentile of likes. This is good as we are mostly concerned with whether the restaurant will perform well or not so the high accuracy of predictions for the two extremum is a welcome feature. This allows us to fairly accurately predict the general performance of the business opportunity. Different binning methods for the classes were attempted, but the use of 3 bins by far yielded the best Jaccard Similarity Score.

Additionally, not only are we attempting to predict the general business performance but also pull insights to inform on business strategy. In this case strategy insight can be gleamed from the coefficient values from running the logistic regression on the full dataset. As such, we can see that opening a restaurant in San Francisco, opening a bar, or serving cuisine that is American or Asian in nature, are associated negatively with "likes." This suggests that the business opportunity should be opening a restaurant in either Los Angeles or San Diego, with a cuisine that is European, Latino, or casual in nature would be the best approach for maximizing likes.

6. Conclusion

In conclusion, after analyzing restaurant "likes" in California from 300 restaurants, we can conclude that the approach to best take in regard to maximizing business performance (as measured by "likes") is to open a restaurant that is either European, Latino, or casual and that opening the venue in either Los Angeles or San Diego rather than San Francisco would be the best approach. Additionally, the predictive capabilities of the logistic regression prediction model are most accurate for classifying whether a restaurant will fall in either the best or worst classes when the data is binned into 3 classes.