

Yelp

Predicting Restaurant Success

Neel Vasa
Aditya Vaidya
Shreya Kamani
Manan Upadhyay
Mark Thomas

IOM 528 SPRING 2014

Executive Summary

In this day and age of social media, with its instant reviews and referrals, many businesses seek to find the best way to improve their online reviews in order to attract and keep clientele or they seek to establish a new business that is mindful of what drives good reviews. This is particularly true for restaurants and food establishments. Several studies have been done that have looked at the correlation between a restaurant's star rating on Yelp and their success related to the number of reservations and the revenue increase. We have obtained a dataset from Yelp, which lists businesses in and around Phoenix AZ. The dataset has each business's number of reviews, including their ratings, and other common attributes. Through this info we plan to mine the data for the aspects that make a particular business type successful.

We decided to narrow our current models on restaurant business types as this is what Yelp is primarily used for. We started by determining what successful exactly means given the data we have obtained. We decided to use the number of reviews and the number of star ratings, and then chose a minimum threshold that the business must obtain to be considered successful. Using this threshold we assigned a binary number to each business with a 1 denoting success. This allows us to run various prediction models, such as decision tree, logical regression, and neural networks in order to find what independent variables are most likely to lead to the success of the business.

The data set was quite large to begin with, so we grouped many of the attributes together, and aimed our focus on these grouped categories. We then narrowed the number of businesses we included to those with zip codes that contained more than 20 businesses in that particular zip code. So each of our business IDs had the following independent variables to help us determine the success outcome: Zip Code, Mexican, Asian, American, European, Latin American, Fast Food, Cafes, # Meals, Meal Options, # Features and Amenities.

Using these independent variables we ran several models in SAS to identify those variables that most accurately predict the success of a business on Yelp. We ultimately arrived at a model that used the most significant variable of success, and found that approximately 50% of success, as determined for our model, could be explained by the variables. The independent variables that had the greatest effect on the success were Fastfood, European Dishes, Latin American Dishes, Number of Meals and Number of Features and Amenities.

This information can assist businesses in improving their success through ratings and number of reviews. It can also provide valuable information to those who want to establish a new restaurant.

As businesses try to cater to more and more customers through social media, and as they develop strategies and methods for attracting those customers, social media can help to supply valuable info on customers' preferences. Through better mining of the data that is available through sites like Yelp and others, businesses can gain the info they need to market better to their customers.

Problem

We set out to discover if business attributes given by Yelp can help to predict whether a restaurant is successful. Business attributes that are collected by Yelp are things such as location (address), Food Type, Number of Meals, Amenities (TV, WiFi, etc.), Late Night Hours, Delivery, etc. We want to use this information to run models that help point to those attributes required to gain success.

We will run prediction analysis using decision trees, multiple regression, and neural networks using both JMP and SAS to find the model that has the best prediction rate. The main application of our result will be in helping restaurateurs decide the factors to consider for improving an existing restaurant or while setting up a new restaurant. On the basis of this idea we have defined the following hypotheses.

Hypothesis

Ho: Food type or categories would not determine the success of a restaurant

Ha: Food type or categories would determine the success of a restaurant

Ho: Having a large number of features and amenities would not increase the chances of a restaurant being successful

Ha: Having a large number of features and amenities would increase the chances of a restaurant being successful

Ho: Location of the restaurant does not play an important role

Ha: Location of the restaurant plays an important role

Data

Our project sources data from the Yelp Academic Dataset available at:

https://www.yelp.com/academic_dataset

Yelp is providing a generous sample of data from the greater Phoenix, AZ metropolitan area including:

15,585 businesses

111,561 business attributes

11,434 check-in sets

70,817 users

151,516 edge social graph

113,993 tips

335,022 reviews

Dataset Description: The dataset is a single compressed file, composed of one json-object per line. Every object contains a 'type' field, which tells whether it is a business, a user, or a review. The data consists of Business objects that contain basic information about local businesses. The 'business_id' field can be used with the Yelp API to fetch even more information for visualizations. The fields are as follows:

```
{ 'type': 'business',
  'business_id': (encrypted business id),
  'name': (business name),
  'neighborhoods': [(hood names)],
  'full_address': (localized address),
  'city': (city),
  'state': (state),
  'latitude': latitude,
  'longitude': longitude,
  'stars': (star rating, rounded to half-stars),
  'review_count': review count,
  'categories': [(localized category names)]
  'open': True / False (corresponds to closed, not business hours),
  'hours': { (day_of_week):
              { 'open': (HH:MM), 'close': (HH:MM) }, ... },
  'attributes': { (attribute_name): (attribute_value), ... }, }
```

These business objects also contain restaurants which of primary interest and these are extracted to enrich the raw dataset. The restaurant business objects contain various other fields specific to restaurants fields in json representation that are then converted to more amenable formats. The restaurant data fields also contain additional fields for the type of content served and also amenities available eg. Wi-Fi, TV screens etc.

Data Usage in Project

The raw data format is JSON and is not supported by natively by JMP and SAS Miner. Also the signal to noise ratio in the raw data is too large and thus we clean the data using a JAVA code and convert it to a .CSV format that is supported by various applications.

Of the approximately 15000 business objects, about 5000 restaurant business objects were extracted yielding a Signal to Noise Ratio (SNR) = 0.33. This apprised us of the fact that enrich of the data is an extremely important part of the project and we further concentrated more efforts on enriching the data to improve the SNR. Also, these 5000 rows were partitioned into training and testing datasets. K-Fold partitioning was done to negate the effects of potential bias and randomization factors. This partitioning approach helped us achieve stable performance that was resilient and robust.

Procedure

We learned numerous ways in class with respect to analysis of large datasets and their usage in deriving business insights. Of them, we chose the MAGIC framework of analysis where we proceed through the steps of Modeling, Analysis, Grilling, and Improving Continuously (MAGIC) in-order to create an accurate model to address the problem that we tackle.

The following steps were followed in line with the MAGIC framework learnt in class:

1. Research Data sources
 - a. We thoroughly researched data sources at various websites for ways and means to get data and potential ways to enrich it so as it is pertinent to our project problem statement. The library session held in class about getting data was especially helpful as it opened up more avenues to direct our research and ultimately helped us reach Yelp and finalize it as a primary source of our data.
2. Enrich Data and Format Conversion
 - a. The data obtained from Yelp as explained above had a low SNR(signal to noise ratio). This adversely affected our preliminary models as we were not able to achieve the desirable accuracy in prediction results.
 - b. We thus, enriched the raw data, using various approaches and also converted the json format data into a more universal .csv format data.
 - c. The format conversion helped us in using various models and analytic tools on our data to compare on and we were able to achieve the best results in SAS Miner.
3. Model Building
 - a. This was the crux of the project and we considered various models from Decision Trees to Logistic Regression to Neural Networks.
 - b. Each approach had its own advantages and fallbacks for eg where neural networks provided a tremendous accuracy in training data set, the failed miserably on the testing dataset.
4. Improving Continuously
 - a. We implemented various models constantly tweaking and changing the parameters in order to achieve the desired accuracy.
 - b. We also followed a novel way of K-Fold partitioning in the data set to obtain better results and robust performance.
5. Final Model: A final model was thus reached which provided us with a satisfying accuracy, robust performance and was consistent with our business sense.

Extract, Transform, and Load Data

The original data set was a 300 MB JSON file containing a large number of extraneous attributes. We refined this in a first pass of the extraction to restrict only business objects.

In the second pass, the business objects were further refined using various JSON parsing libraries to remove various attributes, like location and category data and then these were stored in a data format (.csv) which can be opened using JMP and SAS Miner. This data set still had numerous independent variables, and therefore grouping some of these variables needed to be considered to improve our model outputs.

In a third refinement of the data, the attributes were counted and accumulated, replacing individual binary 0/1 values of attributes with a count of the number of attributes and amenities the restaurant provides. In this pass over of the data set, outliers in zip codes were eliminated as well, retaining only zip codes with 20 or more restaurants. Data attributes accumulated include meals served, open times (late-night etc.), parking, Wi-Fi etc. Additionally, food categories were grouped into broader categories, for example various categories like Pizza, Burgers, and Sandwiches were grouped under Fast Food.

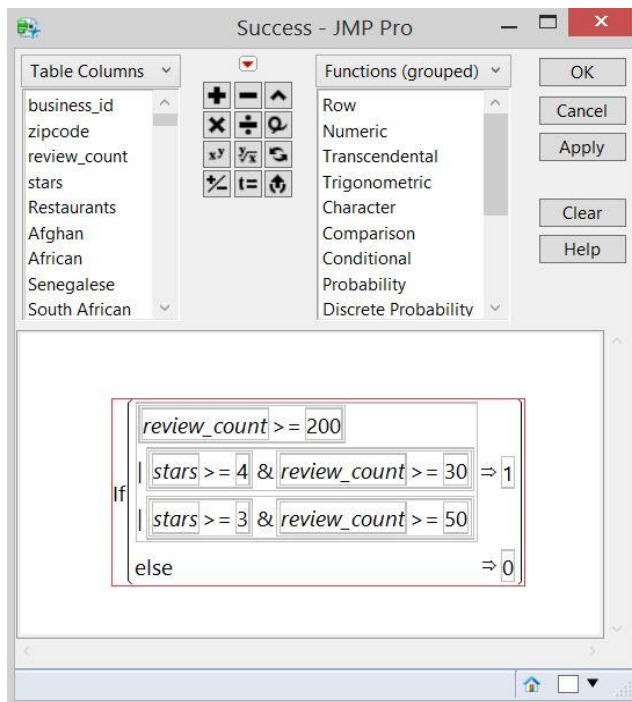
Appendix 1 shows a sample tuple of the refined data set.

Key Summary Stats

Appendix 2 shows the distribution info between several of the independent variables from the dataset before we condensed them down into smaller groups. These were the most significant variables from the larger dataset. They were review count and stars, which are quantitative, and then Mexican, Caters, Price Range, Really Fast Food, and Alcohol, which are all qualitative. The greatest number of stars for the restaurant businesses was 3, which had 2538 instances; the second most number was 4 with 2015 instances. The review counts had a greater distribution, ranging from 0 to over 1000. The mean for reviews was at 44, with a standard deviation of 76, and the median was 16. So a majority of the restaurants had less than 100 reviews, with a small number receiving over 100, including some outliers that received over 500.

Appendix 3 shows the distribution from the refined data set. These are some of the key variables which include review count and stars, as well as categories Mexican, Really Fast Food, American Dishes, and the # Meals, and # Features and Amenities. Similar to the larger data set, the mean for review counts is 47 and the median is 18, so the majority of reviews are under 100 with a few outliers over 500. The stars also are largest at 3 with 2328, and second is 4 with 1857 at that level. Another variable we created is Features and Amenities, which has a fairly normal distribution with a mean around 4. We also found that there are some businesses that have some missing data, like number of meals per day listed is 0. So we need to consider the best way to handle this, for example defaulting all those with number of meals zero to a one.

Below is the calculation used for measuring success according to the data we obtained.



Models

Our next best model was run using Decision Tree. We used the optimal values and got to 7 splits. It separated first by the number of Features and Amenities and then split by a number of different zip codes as well as number of meals and then the type of food.

Appendix 4 shows the full tree. Appendix 5 shows the R Square for the model.

The R square for the testing and validation were very far off though, so it was over fit for training data. We decided to try Neural Network to test the model and see if we could improve on the R Square and the significance of the variables.

Our best model was run using Neural Network. We tried several variables and came up with the following parameters to give us the best R Square model for both Training and Validation.

Parameters:

Zipcode
Fastfood
European Dishes
Latin American Dishes
Meals
Features and Amenities

Rsq: 0.58 (Training) 0.51 (Testing)

I had first Fit Y by X for all our variables versus success. These 6 variables were among the ones that had maximum significance on success. The reason why these are statistically significant and make business sense:

1. Zip code: Geographical areas play a significant role in contributing towards success. Downtown areas and ASU campus areas in particular have a high success rate.
2. Fastfood: Quick bites. People in need of a quick snack or a cheap meal frequent fast food joints.
3. European Dishes: More often than not, fine dining restaurants where food quality and service is monitored and scrutinized strictly. Clients usually are big spenders, rich families or a venue for business meals.
4. Latin American: Frequency of Latin American restaurants is low compared to Mexican/Fast Food/Asian (less than a quarter of these places). But it is a popular cuisine because of lot of Latin Americans in Phoenix.
5. Meals: Restaurants serving 1 or 2 meals (lunch and dinner) are more successful than those serving 3 or more meals (like brunch). Maybe because of less popularity of brunch?
6. Features and Amenities: More the features and amenities (WiFi, Parking, Home Delivery, Alcohol etc) more likelihood of being a success.

Appendix 6 and Appendix 7 show the neural network.

Improvements

The project has a few shortcomings which can be overcome by utilizing various other data sources. The accuracy of Yelp data might be questionable, as competitors might deliberately give bad reviews or the owners themselves might start with very good reviews of their own business. To eliminate this, a few techniques could be used, such as grouping of reviews by time period (month-wise or quarterly).

Another improvement could be the addition of other regions' data. Including 5-6 cities could give a robust model, as this data might suffer from over fitting to the Phoenix area.

Lastly, including demographic data from census or other information sources could lead to a better indication of success/failure of a restaurant. For example, a restaurant in a prosperous neighborhood might have several 4/5 rated competitors, and would need a higher sustained rating to be deemed a success

Conclusion

The Yelp Phoenix data set proved to be a very rich and powerful data source for predicting restaurant success as defined in the project. Undertaking an extensive and rigorous ETL process helped significantly increase the accuracy rates for prediction in the models as well. The findings were clear and concise: A restaurant is highly likely to be a success if it caters to certain popular categories and provides a large number of services, regardless of its location. This particular finding is in contrast to the preliminary belief of location being the most crucial factor to success, and can probably be attributed to the fact that a better location also equals greater competition.

Addition of demographic data, additional cities' data and improvement in Yelp data as described in the Improvements section above has the potential to lead to a robust and highly accurate model for future use.

Bibliography:

The following resources were used for the development of this project:

http://www.yelp.com/dataset_challenge/

<http://www.jmp.com/support/downloads/documentation.shtml>

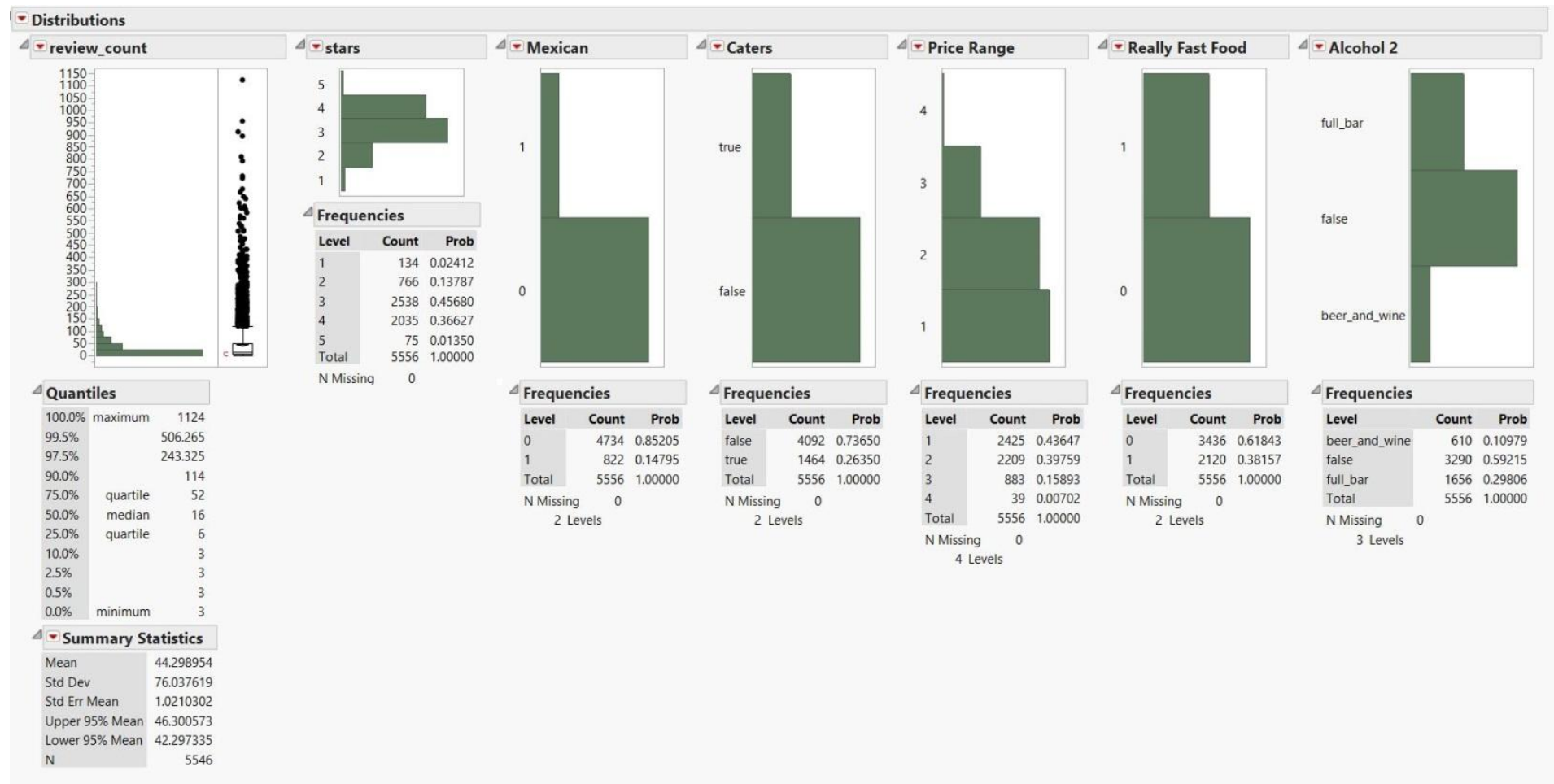
<http://support.sas.com/documentation/onlinedoc/miner/>

IOM 528 class notes and minicases

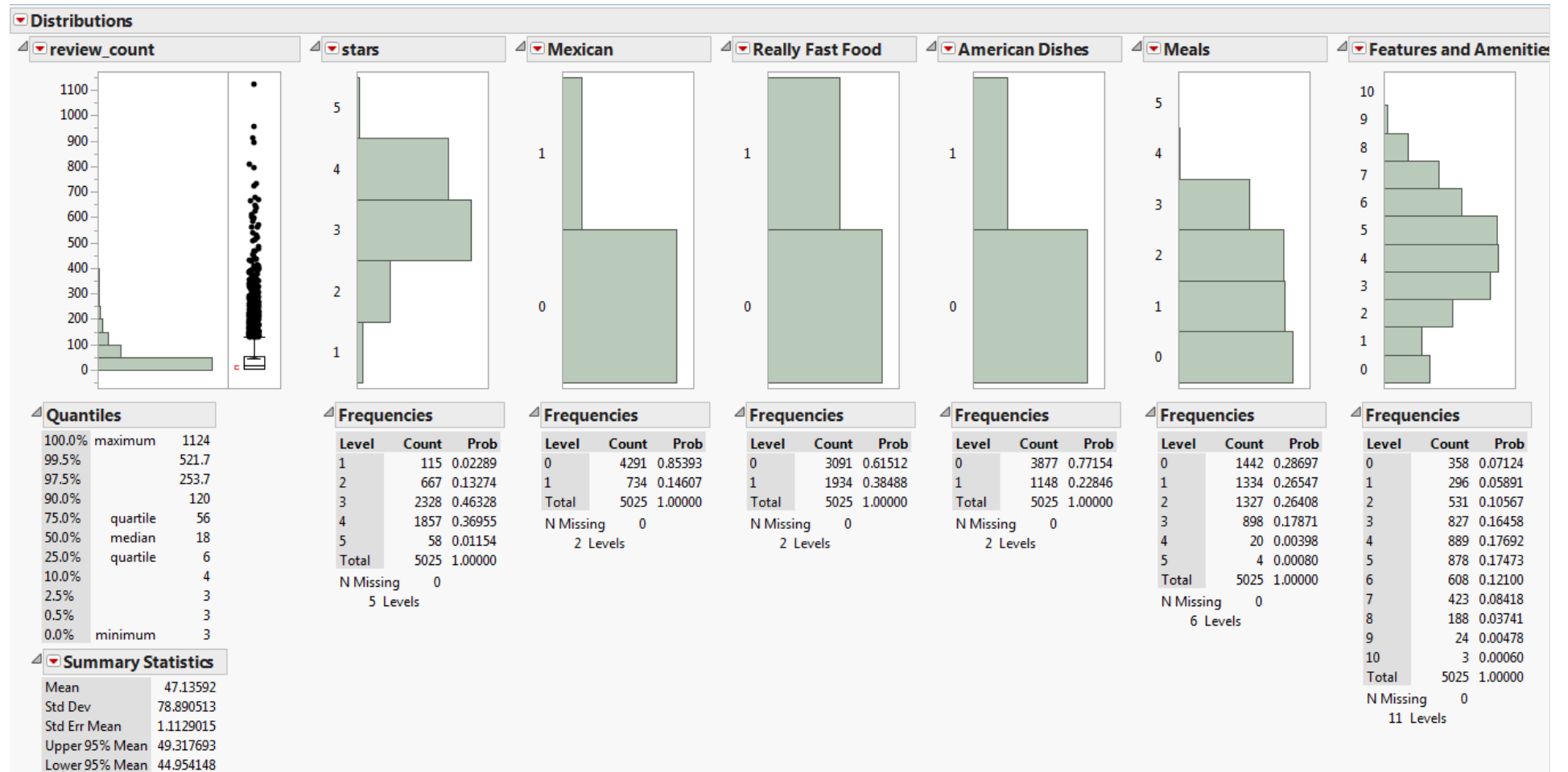
Appendix 1

business_id	zipcode	review_count	stars	Mexican	Really Fast Food	Asian Dishes	American Dishes	European Dishes	Latin American Dishes	Cafes	Meals	Meal Options	Features and Amenities	Success
OYlq2UoKPQ O_ElqhC3a8	85003	125	3	0	0	0	0	0	0	0	3	1	6	1

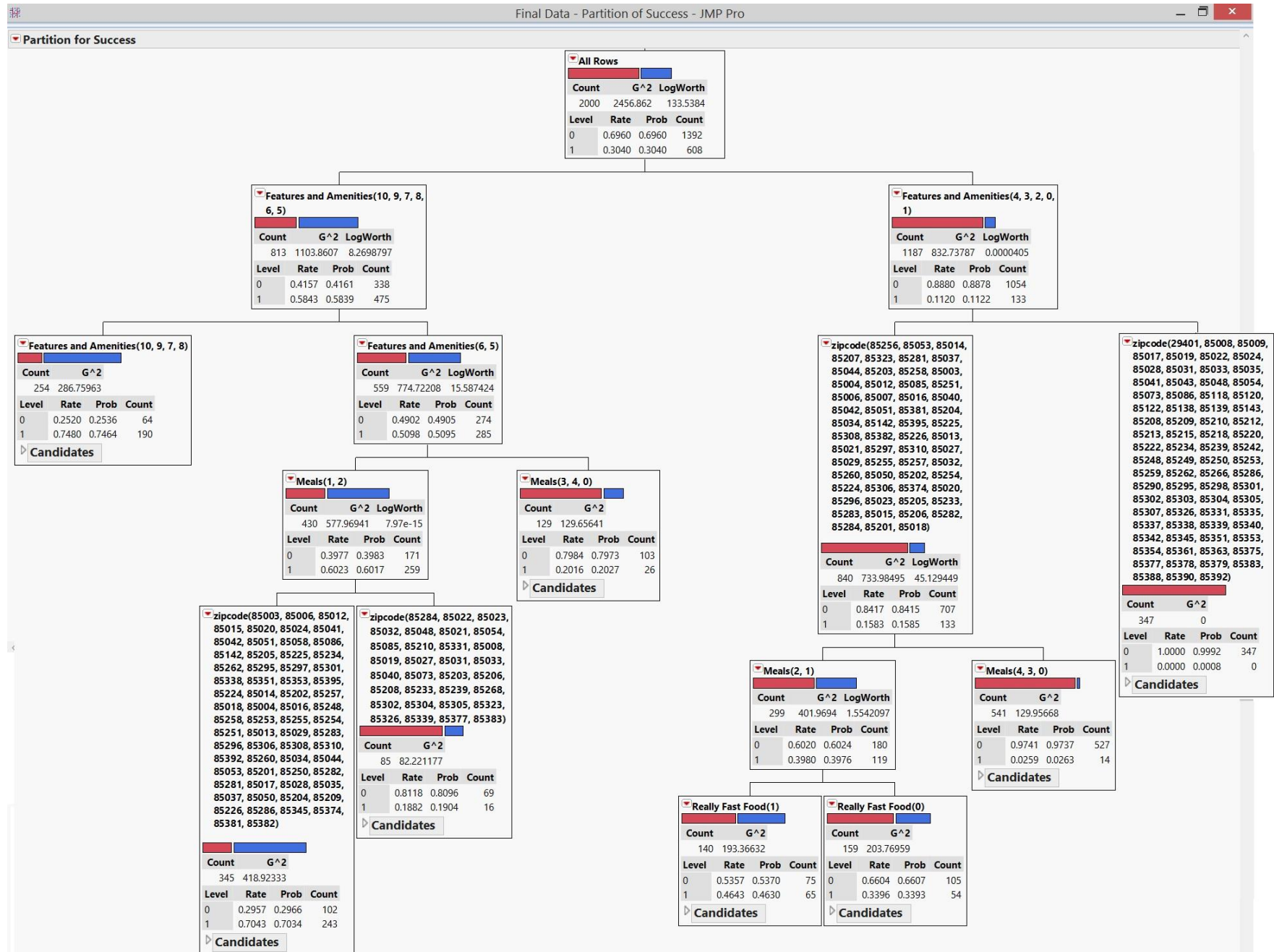
Appendix 2



Appendix 3



Appendix 4



Appendix 5

	RSquare	N	Number of Splits
Training	0.412	2000	7
Validation	0.280	3546	

Appendix 6

