

### **Model description:**

I used the principles of model-based recommendation systems, XGBRegressor(a regressor based on the decision tree) to train the model.

There were total of 5 files that were given:

- a. train\_review.json
- b. user.json – user metadata
- c. business.json – business metadata, including locations, attributes, and categories
- d. user\_avg.json – containing the average stars for the users in the train dataset
- e. business\_avg.json – containing the average stars for the businesses in the train dataset

### **Attribute selection:**

I used these files to extract features required for the classification. From the user.json file I extracted features such as compliment\_photos, fans, funny and important among this was the friends list that I could use to analyze the ratings made by a user's friends on a business. The mean and mode of the friend rating is being used to train the model.

From the business.json file, I extracted attributes such as is\_open and category of business, postal code and number of days the business is open. I stored this information in the business\_dict data structure.

Similarly I extracted the information from the user\_avg.json and business\_avg.json files.

For every (user\_id, business\_id) pairs, I extracted the feature attribute information, tabulated them in lists and collected the actual star rating given by the user on a business. I also calculated the average rating by all the friends of a user.

### **Hyperparameters selection:**

To find the optimal hyper parameters that will give the best results, I ran GridSearchCV on a set of parameters to find the best parameter values for the XGBRegressor. Some parameters that were searched was learning\_rate, max\_depth, gamma, min\_child\_weight and colsample\_bytree. These parameters must be optimal to avoid overfitting.

I divided my dataset into 2 parts one having the values of user-average ratings and the other one which is null. I generated 2 model files.

### **Prediction:**

For every (user\_id, business\_id) pairs, I extracted the feature attribute information, tabulated them in lists. If the user was a new user or business was a new business, no features were added. It was then passed to the xgb predict function which provided the results. This information was written to the output file.