

**INFO7250 Engineering Big-Data Systems - Fall 2016**

# **Yelp: Predicting User Ratings for New Business**

Section - 01

**Team Members** – Kunal Singh Deora, Manasi Laddha, Swarna Satishkumar  
Dommeti, Bhavna Menghrajani, Sarthak Agarwal

**Professor** – Dino Konstantopoulos

December 12, 2016

## I. Table of Contents

<b>1. Introduction</b>	<b>3</b>
1.1 Idea	3
1.2 Challenge	3
1.3 Approach	3
<b>2. Dataset Description</b>	<b>4</b>
<b>3. Data Cleaning &amp; Feature Engineering</b>	<b>6</b>
3.1 Dataset Preparation	6
3.2 Feature Generation using Sentiment Analysis	9
<b>4. Machine Learning</b>	<b>12</b>
4.1.1 Data Splitting	12
4.1.2 Algorithm	12
4.1.3 Finding new business	14
4.1.4 Prediction	15
<b>5. Application</b>	<b>18</b>
<b>6. Challenges faced</b>	<b>20</b>
<b>7. Tools &amp; Technologies Used</b>	<b>21</b>
<b>8. Conclusion</b>	<b>22</b>
<b>9. Future scope</b>	<b>23</b>
<b>10. Work Distribution</b>	<b>24</b>
<b>11. References</b>	<b>25</b>

# Chapter 1

## 1. Introduction

### 1.1 Idea

The idea is to build a predictive system that predicts the degree of likeness of new food joints for users based on users past historical data on a global geographical scale.

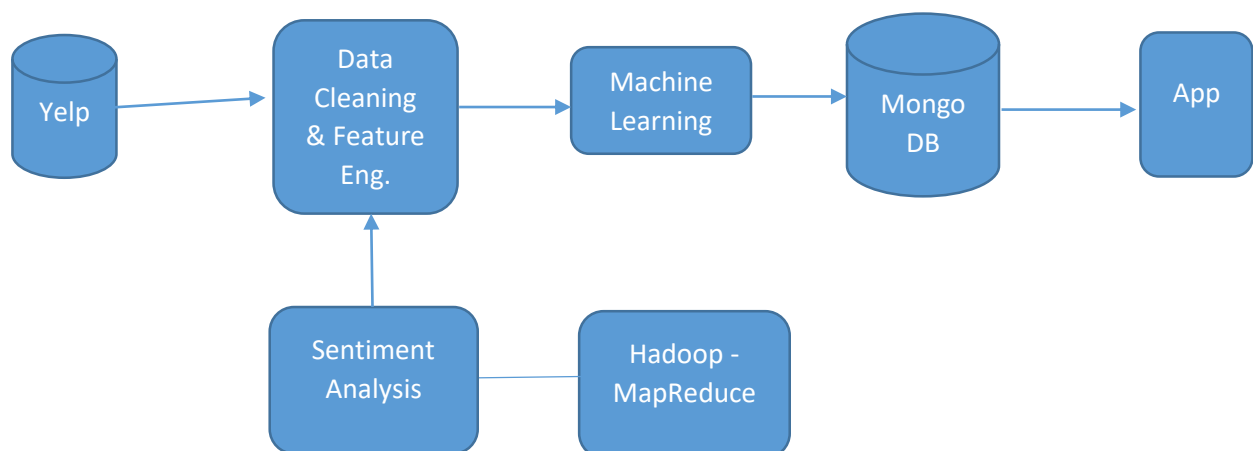
e.g. If user travels to California from Boston, the system will predict which restaurants he will like in his five-mile radius and how will he rate the restaurant based on his/her previous pattern on Yelp.

### 1.2 Challenge

Out of more than 100 features across five different data feeds, choose the appropriate features or create a new feature to get a high accuracy for each user.

### 1.3 Approach

A high level diagram/flow of the project:



# Chapter 2

## 2. Dataset Description

The Yelp dataset consists of five data feeds but we primarily work on the below mentioned two feeds.

- Business – Information of all the businesses in Yelp.
- Review – Reviews of all the business.

The data is provided by Yelp on their website which consists of 85,539 businesses and 2,685,066 reviews. The file sizes were 75MB and 2.2GB respectively.

The data is provided as JSON files and their structure is as follows:

### **BUSINESS:**

```
{
  '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),
    ...
  },
}
```

**REVIEW:**

```
{
  'type': 'review',
  'business_id': (encrypted business id),
  'user_id': (encrypted user id),
  'stars': (star rating, rounded to half-stars),
  'text': (review text),
  'date': (date, formatted like '2012-03-14'),
  'votes': {(vote type): (count)},
}
```

# Chapter 3

## 3. Data Cleaning & Feature Engineering

### 3.1 Dataset Preparation

3.1.1 Convert JSON to CSV We used Python to convert JSON to CSV and the code is as follows

```
#!/usr/bin/env python3
# coding: utf-8 -*-

import argparse
import collections
import csv
import simplejson as json

def read_and_write_file(json_file_path, csv_file_path, column_names):
    """Read in the json dataset file and write it out to a csv file, given the column names."""
    with open(csv_file_path, 'wb+') as fout:
        csv_file = csv.writer(fout)
        csv_file.writerow(column_names)
        with open(json_file_path) as fin:
            for line in fin:
                line_contents = json.loads(line)
                line_contents = line_contents.replace("\n", "")
                csv_file.writerow(get_row(line_contents, column_names))

def get_superset_of_column_names_from_file(json_file_path):
    """Read in the json dataset file and return the superset of column names."""
    column_names = set()
    with open(json_file_path) as fin:
        for line in fin:
            line_contents = json.loads(line)
            column_names.update(
                set(get_column_names(line_contents).keys())
            )
    return column_names

def get_column_names(line_contents, parent_key=''):
    column_names = []
    for k, v in line_contents.items():
        column_name = "{0}.{1}".format(parent_key, k) if parent_key else k
        if isinstance(v, collections.MutableMapping):
            column_names.extend(
                get_column_names(v, column_name).items()
            )
        else:
            column_names.append((column_name, v))
    return dict(column_names)

def get_nested_value(d, key):
    if '.' not in key:
        if key not in d:
            return None
        return d[key]
    base_key, sub_key = key.split('.', 1)
    if base_key not in d:
        return None
    sub_dict = d[base_key]
    return get_nested_value(sub_dict, sub_key)

def get_row(line_contents, column_names):
    """Return a csv compatible row given column names and a dict."""
    row = []
    for column_name in column_names:
        line_value = get_nested_value(
            line_contents,
            column_name,
        )
        if isinstance(line_value, unicode):
            row.append('{0}'.format(line_value.encode('utf-8')))
        elif line_value is not None:
            row.append('{0}'.format(line_value))
        else:
            row.append('')
    return row

if __name__ == '__main__':
    """Convert a yelp dataset file from json to csv."""
    parser = argparse.ArgumentParser(
        description="Convert Yelp Dataset Challenge data from JSON format to CSV.",
    )
    parser.add_argument(
        'json_file',
        type=str,
        help='The json file to convert.',
    )
    args = parser.parse_args()

    json_file = args.json_file
    csv_file = '{name}.csv'.format(name=json_file.split('.')[0])
```

### 3.1.2 Used R to analyze the data and prepare it for Machine Learning input.

Following steps were performed:

- Filter the data related to food and similar categories.
- Filter the data for the users that have reviewed more than hundred businesses.
- Extract the relevant features out of 95 attributes in the Business dataset which will be used in our Machine learning algorithm.
- Rename the attributes to more meaningful names.
- Build and structure the data for our input to the Machine Learning Algorithm

#### Load the Business, Review and User CSVs into dataframe

```
#### BUSINESS
business <- read.csv("D:/GRAD_SCHOOL/Fall2016/Project_Yelp/DatasetsInCSV/DatasetsInCSV/yelp_academic_dataset_business.csv")
head(business)
colnames(business)
```

```
##### REVIEW
file.path <- "D:/GRAD_SCHOOL/Fall2016/Project_Yelp/DatasetsInCSV/DatasetsInCSV/yelp_academic_dataset_review.csv"
library(data.table)
review <- fread(file.path)
```

```
#### USER
user <- read.csv("D:/GRAD_SCHOOL/Fall2016/Project_Yelp/DatasetsInCSV/DatasetsInCSV/yelp_academic_dataset_user.csv")
head(user)
str(user)
```

#### Filter the data related to food and similar categories

```
business_foodData <- sqldf('select distinct * from business where categories like \'%Food%\'
                             or categories like \'%Restaurants%\'
                             or categories like \'%Lounges%\'
                             or categories like \'%Nightlife%\'
                             or categories like \'%Bars%\'')
str(business_foodData)
```

#### Filter the data for the users that have reviewed more than hundred businesses

```
xx <- sqldf('select r.user_id, count(r.business_id)
             from review_user_business r, business_foodData b
             where r.business_id = b.business_id
             group by r.user_id
             having count(r.business_id) > 100')
head(xx)

xx2 <- sqldf('select a.user_id, b.business_id from
             xx a,
             business_foodData b,
             review_user_business r
             where a.user_id = r.user_id
             and b.business_id = r.business_id')
head(xx2)
write.csv(xx2, 'D:/GRAD_SCHOOL/Fall2016/Project_Yelp/DatasetsInCSV/DatasetsInCSV2/Our_Users_GreaterThan100_Businesses.csv')
```

Rename the attributes to more meaningful names.

```
nrow(business_foodData_relevantFeatures)

colnames(business_foodData)[32] <- "attributes_Price_Range"
colnames(business_foodData)[40] <- "attributes_Accepts_Credit_Cards"
colnames(business_foodData)[44] <- "attributes_Take_out"
colnames(business_foodData)[55] <- "attributes_Delivery"
colnames(business_foodData)[60] <- "attributes_Wheelchair_Accessible"
colnames(business_foodData)[42] <- "attributes_Good_For_Lunch"
colnames(business_foodData)[68] <- "attributes_Good_For_dinner"
colnames(business_foodData)[20] <- "attributes_Good_For_brunch"
colnames(business_foodData)[33] <- "attributes_Good_For_breakfast"
colnames(business_foodData)[51] <- "attributes_Takes_Reservations"

nrow(business_foodData_relevantFeatures)
```

Build and structure the data for our input to the Machine Learning Algorithm.

We structured the data based on the input that our ML algorithm needs. The input to the algorithm will be something like this:

Users\_Id, Business\_id, ----Relevant Features---, Label

```
AllUsersAllBusinessAllFeatures <-
  sqldf('select r.user_id, r.business_id, a.latitude, a.longitude, a.review_count,
    a.attributes_Price_Range, a.attributes_Accepts_Credit_Cards, a.attributes_Take_out,
    a.attributes_Delivery, a.attributes_Wheelchair_Accessible, a.attributes_Good_For_Lunch,
    a.attributes_Good_For_dinner, a.attributes_Good_For_brunch, a.attributes_Good_For_breakfast,
    a.attributes_Takes_Reservations, round(a.stars) as business_stars, a.attributes_Parking_lot,
    a.attributes_Parking_street, a.attributes_Parking_garage, a.attributes_Parking_valet,
    a.attributes_Parking_validated, a.attributes_Good_For_Groups, a.attributes_Good_For_kids,
    a.attributes_Ambience_casual, a.attributes_Ambience_romantic, a.attributes_Ambience_upscale,
    a.attributes_Ambience_classy,
    r.review_avg_stars
  from business_foodData a, review_final_round r
  where a.business_id = r.business_id')

colnames(AllUsersAllBusinessAllFeatures)
nrow(review_final_round)
```



## 3.2 Feature Generation using Sentiment Analysis

We generated a new feature using sentiment analysis of Reviews in Yelp dataset. This feature is a sentiment score of a business based on the sentiments of all the reviews that the business has reviewed. According to us this feature is a true representation of what user thinks about the business.

We wrote the code in R and tried performing sentiment analysis, but since the review file is approximately 2.2 GB and we have to process each word of each review, R was going out of memory. Hence we implemented it using Hadoop Map Reduce and ran the code in Distributed mode.

Our MapReduce algorithm is as follows:

1. The input to our Mapper is the reviews file in JSON format and used Google's GSON library to read the file and assign it to the class.
2. In the Mapper we calculate the sentiment score of all the reviews and output the Business\_id as key and a combination of Review\_id and sentiment score as value.

```
public class SentimentMapper extends Mapper<LongWritable, Text, Text, Text> {  
    @Override  
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {  
        Gson gson = new Gson();  
        // used gson object so that we can assign the headers of the json into a class.  
        if (!value.toString().contains("{\"text\":\"\"}")) {  
            return; // Not a review we're interested in.  
        }  
        // takes the values from the json and stores in the Class so that we can use the reviews as a String in our program.  
        Review review = gson.fromJson(value.toString(), Review.class);  
        DoubleWritable sentimentValue = new DoubleWritable(review.calculateSentimentValue());  
        Text businessId = new Text(review.business_id);  
        // making the value as review_id+ sentimentValue (which is calculated by using lexicon.ttf file which stores the sentiment score of that particular word.  
        Text value_text = new Text(review.review_id + ":" + sentimentValue);  
        context.write(businessId, value_text);  
    }  
}
```

3. In the Reducer, we split the value and extracts the sentiment score and calculate the average sentiment score for each business.
4. The output of the reducer is the Business\_id and its sentiment score.

```

public class RatingReducer extends Reducer<Text, Text, Text, DoubleWritable> {

    @Override
    public void reduce(Text key, Iterable<Text> values, Context context)
        throws IOException, InterruptedException {
        double total = 0;
        double length = 0;

        List<Double> dd = new ArrayList<Double>();

        // key=business_id and values= review_id + sentiment_score
        // split the value on : so that we can get the sentiment score of the against that business id

        for (Text t : values ) {
            String[] arrVal = t.toString().split(":");
            String sentiment_score = arrVal[1];
            Double sentiment_score_double = Double.parseDouble(sentiment_score);
            dd.add(sentiment_score_double);
        }

        // store the sentiment value in list of double so that we can aggregate while its on the loop.
        for (Double sentiment : dd) {
            total += sentiment;
            length++;
        }

        //calculate the average sentiment score of each business

        double average = total / length;

        if(key.toString().equalsIgnoreCase("--0ZoBTQWQra1FxD4rBWmg"))
            System.out.println("Found it again");

        //emit bbusiness_id, sentiment_score
        context.write(key, new DoubleWritable(average));
    }
}

```

Since the reviews data was of 1.9 GB, we ran the map reduce job on fully distributed system. It took around 10 minutes to complete the job. Below are the screenshot of the job-tracker sandbox.

JOBTRACKER SAND – Fully Distributed Mode

**Cluster Summary (Heap Size is 123 MB/889 MB)**

Maps	Reduces	Total Submissions	Nodes	Map Task Capacity	Reduce Task Capacity	Avg. Tasks/Node	Blacklisted Nodes
2	1	4	1	2	2	4.00	0

**Scheduling Information**

Queue Name	Scheduling Information
default	N/A

**Filter (Jobid, Priority, User, Name)**  
 Example: 'user:smith 3200' will filter by 'smith' only in the user field and '3200' in all fields

**Running Jobs**

Jobid	Priority	User	Name	Map % Complete	Map Total	Maps Completed	Reduce % Complete	Reduce Total	Reduces Completed	Job Scheduling Information
job_201612101541_0007	NORMAL	manasi-desktop/manasiladdha	Yelp Sentiment Analysis	11.42% <div style="width: 11.42%;"></div>	35	4	3.80% <div style="width: 3.80%;"></div>	1	0	NA

## JOBTRACKER SAND – Fully Distributed Mode

### Cluster Summary (Heap Size is 123 MB/889 MB)

Maps	Reduces	Total Submissions	Nodes	Map Task Capacity	Reduce Task Capacity	Avg. Tasks/Node	Blacklisted Nodes
0	0	4	<a href="#">1</a>	2	2	4.00	<a href="#">0</a>

### Scheduling Information

Queue Name	Scheduling Information
<a href="#">default</a>	N/A

Filter (Jobid, Priority, User, Name)

Example: 'user:smith 3200' will filter by 'smith' only in the user field and '3200' in all fields

### Running Jobs

[none](#)

### Completed Jobs

Jobid	Priority	User	Name	Map % Complete	Map Total	Maps Completed	Reduce % Complete	Reduce Total	Reduces Completed	Job Scheduling Information
<a href="#">job_201612101541_0007</a>	NORMAL	manasi-desktop\manasiladdha	Yelp Sentiment Analysis	<div><div>100.00%</div></div>	35	35	<div><div>100.00%</div></div>	1	1	NA

# Chapter 4

## 4. Machine Learning

### 4.1.1 Data Splitting

Using python, scikit-learn library, we split the data into training and test data using *train\_test\_split* function of *sklearn.model\_selection* library. We set the function to shuffle the data and then split it into 80 percent training and 20 percent test data.

### 4.1.2 Algorithm

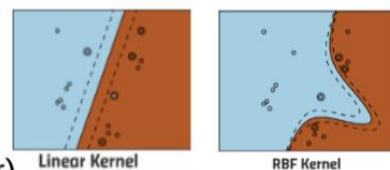
We implemented the two models using scikit libraries: support vector machines (SVM) and random forest classifiers. These algorithms were implemented with the defaults from scikit, which can be found in their user guide. However, we investigated the types of kernels for SVM and the regularization parameters for SVM.

#### a. Random Forest:

While implementing Random forest we used the scikit implementation and the results were over fitted. i.e. the difference between the training and testing accuracy was huge.

#### b. SVC (Support Vector Classification):

- SVC work well for large numbers of training samples spread across large numbers of features
- Supports Multi-Class classification
- Supports Model for different kernels
- SVC parameters
  - $C = 1$  (Penalty Parameter)
    - $C$  increases, overfitting increases
  - `random_state = 0` (Shuffling Paramter)
    - used to shuffle the data before model fitting
  - `kernel = 'RBF'` (data distribution)
    - RBF is Radial Basis Function for non-linear data points
  - `SVC.score`
    - Returns the mean accuracy on the given test data and labels.



We experimented with SVM using scikit's SVM with linear, polynomial and Gaussian kernels. We also performed a parameter sweep and the results were as follows:

Training Accuracy	Testing Accuracy	Parameters
0.623397036	0.536425359	c = 1, kernel = "rbf"
0.714505147	0.433934534	c =2, kernel = "rbf"
0.600952494	0.427557664	c = 1, kernel = "linear"
0.775047514	0.413936467	c =1, kernel = "rbf", no preprocessing
0.617386536	0.426587355	c = 1, kernel = "rbf", randomState = 0
0.577037933	0.423913976	Linear SVC, pentaly = l2
0.545430948	0.413995582	Linear SVC, pentaly = l2, loss = "hinge"

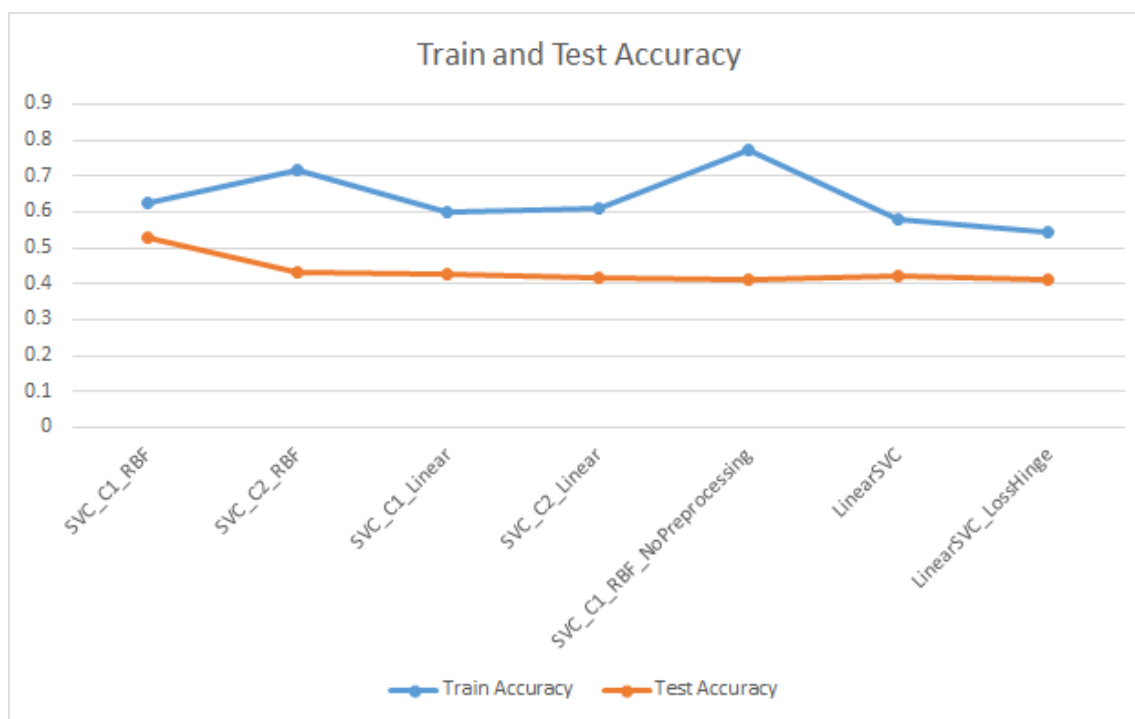


Figure 1: SVM results with different parameters

Based on the above results, we settled on C-SVC with RBF kernel to perform our prediction.

With C-SVC algorithm, our machine learning model achieves a decent accuracy of 62% and we are able to predict the rating of restaurants which have not been reviewed by a specific user nicely.

Radial Basis Function kernel network scale well to large numbers of features in the input space. The pre-processing module further provides a utility class StandardScaler that implements the Transformer API to compute the mean and standard deviation on the training set so as to be later reapply the same transformation on the testing set. The 'c' parameter trades off misclassifying of training examples against simplicity of the decision surface. When gamma is very small, the model is too constrained and cannot capture the complexity or "shape" of the data. The region of influence of any selected support vector would include the whole training set. The resulting model will behave

similarly to a linear model with a set of hyperplanes that separate the centers of high density of any pair of two classes.

### Python Code for Modelling

```
print('-----Model Making Started-----')
results = defaultdict(list)
users_models = defaultdict()
for user, business_dict in our_user_our_business_dict.items():
    final_features = []
    final_labels = []
    for business, features_list in business_dict.items():
        final_features.append(features_list[:-1]) #everything except the last item
        final_labels.append(features_list[-1]) #Last Label

    final_features = numpy.array(final_features)
    final_labels = numpy.array(final_labels)

    X_train, X_test, y_train, y_test = train_test_split(final_features, final_labels,
                                                         test_size=0.2, random_state=0)

    scaler = preprocessing.StandardScaler().fit(X_train)
    X_train_transformed = scaler.transform(X_train)
    clf = SVC(C=1, gamma=0.2).fit(X_train_transformed, y_train)
    print(clf.get_params())
    X_test_transformed = scaler.transform(X_test)
    train_accuracy = clf.score(X_train_transformed, y_train)
    test_accuracy = clf.score(X_test_transformed, y_test)

    if train_accuracy >=0.0 and test_accuracy>=0.0:
        users_models[user] = clf
        results[user].append(X_train.shape)
        results[user].append(X_test.shape)
        results[user].append(train_accuracy)
        results[user].append(test_accuracy)
```

#### 4.1.3 Finding new business

After creating the model for all the users based on the features of their already reviewed businesses, our python code finds out the businesses that the user has not yet reviewed and on which we will perform the prediction.

There were 982 Users and approximately 36991 were non predicted businesses for each users. So the final file with all the features required for the prediction resulted in a 3.47 GB file on which the prediction was done.

## Python Code for new businesses identification

```
print("----- Diff the businesses for users ----- ")
all_business = set(all_business)
all_users_new_businesses = defaultdict(list)
for u, business_set in all_users_old_businesses.items():
    diff_business = all_business.symmetric_difference(business_set)
    all_users_new_businesses[u] = list(diff_business)

print("-----GENERATING NEW BUSINESSES WITH FEATURES FOR USERS -----")
all_users_new_businesses_features = defaultdict(dict)
final_results = []
for u, business_list in all_users_new_businesses.items():
    for b in business_list:
        final_results.append([u, b] + all_business_features[b])
        all_users_new_businesses_features[u].update({ b : all_business_features[b]})

test_user = list(all_users_new_businesses_features.keys())[0]
print(all_users_new_businesses_features[test_user])

# write results to a csv
resultFile = open("FilteredUsers_NewBusinesses_OurFeatures.csv", 'w')
wr = csv.writer(resultFile, delimiter=',', lineterminator='\n')
wr.writerow(['user_id', 'business_id', 'attributes_Price_Range', 'attributes_Accepts_Credit_C',
            'attributes_Take_out', 'attributes_Delivery', 'attributes_Wheelchair_Accessible',
            'attributes_Good_For_lunch', 'attributes_Good_For_dinner', 'attributes_Good_For_l',
            'attributes_Good_For_breakfast', 'attributes_Takes_Reservations',
            'latitude', 'longitude', 'business_stars', 'review_count', 'sentimental_rating'])
wr.writerows(final_results)
```

### 4.1.4 Prediction

To predict, we use the SVC-model generated for each user to predict the User Rating for these non-predicted businesses. We use an AWS EC2 instance with 32GB, 1TB configuration and it took approx. 52 minutes 29seconds to complete the prediction.

For example,

user_id	records_trained_on	records_tested_on	train_accuracy	test_accuracy	actual_avg_stars	predicted_avg_stars
1	80	21	0.7375	0.666666667	2.9	3
2	116	30	0.784482759	0.766666667	4.04	4
3	96	24	0.791666667	0.75	4.75	5

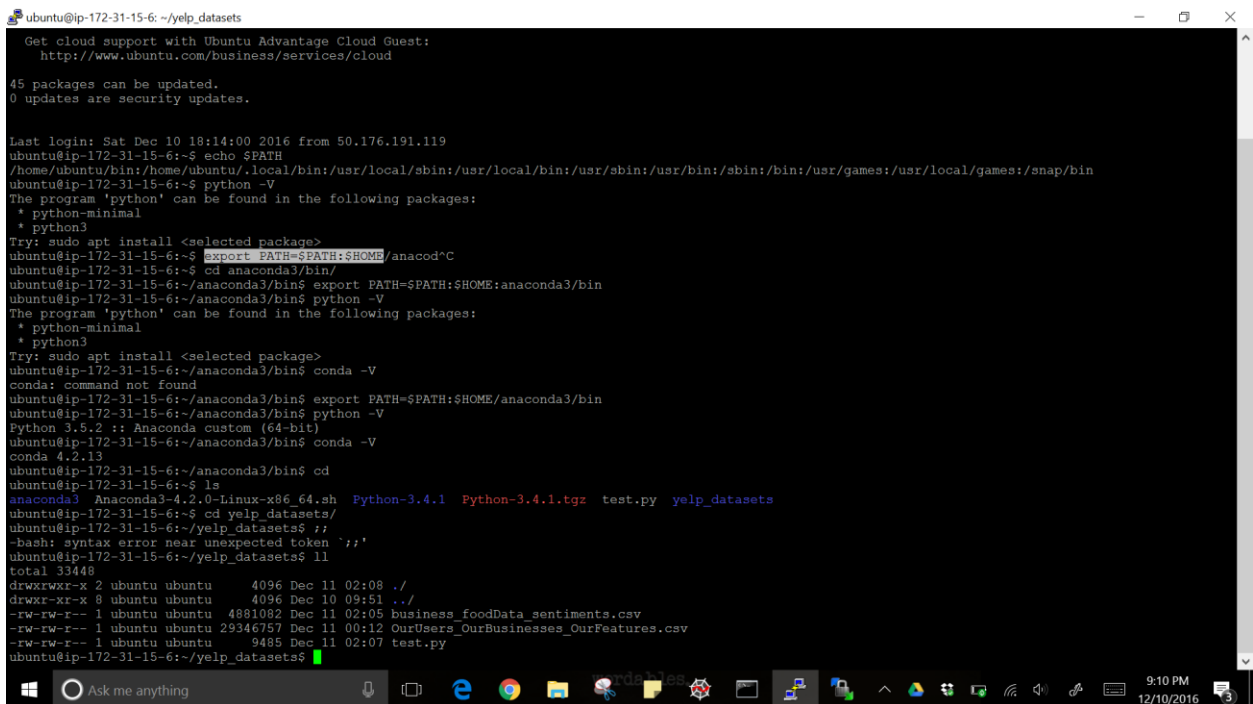
User\_id 1 was trained on 80 records and tested on 21 records and we got a train accuracy 73% and a test accuracy 66%. The actual average stars were 2.9 for the user and the model predicted an average stars of 3 for new businesses.

## Python Code for Predicting

```
print("----- Predicting new businesses and Writing -----")
resultFile = open("FilteredUsers_NewBusiness_Predictions.csv",'w')
wr = csv.writer(resultFile, delimiter=',', lineterminator='\n' )
wr.writerow(['user_id', 'business_id', 'prediction'])
all_users_new_businesses_predictions = defaultdict(dict)
for user, business_list in all_users_new_businesses_features.items():
    for new_business, new_business_features in business_list.items():
        predicted_label = users_models[user].predict([new_business_features])
        wr.writerow([user, new_business, predicted_label[0]])
        all_users_new_businesses_predictions[user] = { new_business : predicted_label[0] }

resultFile.close()
print("done")
```

## AWS EC2 Instance



```
ubuntu@ip-172-31-15-6: ~/yelp_datasets
Get cloud support with Ubuntu Advantage Cloud Guest:
http://www.ubuntu.com/business/services/cloud

45 packages can be updated.
0 updates are security updates.

Last login: Sat Dec 10 18:14:00 2016 from 50.176.191.119
ubuntu@ip-172-31-15-6:~$ echo $PATH
/home/ubuntu/bin:/home/ubuntu/.local/bin:/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/games:/usr/local/games:/snap/bin
ubuntu@ip-172-31-15-6:~$ python -V
The program 'python' can be found in the following packages:
* python-minimal
* python3
Try: sudo apt install <selected package>
ubuntu@ip-172-31-15-6:~$ export PATH=$PATH:$HOME/anaconda3
ubuntu@ip-172-31-15-6:~$ cd anaconda3/bin/
ubuntu@ip-172-31-15-6:~/anaconda3/bin$ export PATH=$PATH:$HOME/anaconda3/bin
ubuntu@ip-172-31-15-6:~/anaconda3/bin$ python -V
The program 'python' can be found in the following packages:
* python-minimal
* python3
Try: sudo apt install <selected package>
ubuntu@ip-172-31-15-6:~/anaconda3/bin$ conda -V
conda: command not found
ubuntu@ip-172-31-15-6:~/anaconda3/bin$ export PATH=$PATH:$HOME/anaconda3/bin
ubuntu@ip-172-31-15-6:~/anaconda3/bin$ python -V
Python 3.5.2 :: Anaconda custom (64-bit)
ubuntu@ip-172-31-15-6:~/anaconda3/bin$ conda -V
conda 4.2.13
ubuntu@ip-172-31-15-6:~/anaconda3/bin$ cd
ubuntu@ip-172-31-15-6:~$ ls
anaconda3  Anaconda3-4.2.0-Linux-x86_64.sh  Python-3.4.1  Python-3.4.1.tgz  test.py  yelp_datasets
ubuntu@ip-172-31-15-6:~$ cd yelp_datasets/
ubuntu@ip-172-31-15-6:~/yelp_datasets$ ;;
-bash: syntax error near unexpected token `;;'
ubuntu@ip-172-31-15-6:~/yelp_datasets$ ll
total 33448
drwxrwxr-x 2 ubuntu ubuntu    4096 Dec 11 02:08 ./
drwxr-xr-x 8 ubuntu ubuntu    4096 Dec 10 09:51 ../
-rw-rw-r-- 1 ubuntu ubuntu 4881082 Dec 11 02:05 business_foodData_sentiments.csv
-rw-rw-r-- 1 ubuntu ubuntu 29346757 Dec 11 00:12 OurUsers_OurBusinesses_OurFeatures.csv
-rw-rw-r-- 1 ubuntu ubuntu  9485 Dec 11 02:07 test.py
ubuntu@ip-172-31-15-6:~/yelp_datasets$
```



# Machine Learning Modelling

```
ubuntu@ip-172-31-15-6: ~/yelp_datasets
1, 36.1239576, -115.2079142, 3.5, 11, 5], 'o9C4hmK1QfYpDgE0PytGkw': [0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 36.2644291578, -115.257308033, 3.5, 213, 4], 'inLlapU-uOJ
AH9a114mX9g': [0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 36.1718303, -115.1420069, 4.0, 17, 4], 'D3e0Caowj 1X4o6C9h2N1w': [0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 33.3777879, -111
.9364361, 2.0, 21, 4], 'ivVVCpUm9mmU6K7xR41vEg': [0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 33.3044287, -111.9836263, 3.5, 43, 4], 'zDjScXELRtHjV84peJL-zQ': [0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 45.510394, -73.564563, 4.5, 6, 5], '9ZBKg9mXQ1lRJLip7e08A': [0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 45.6294231, -73.5864417, 4.0, 7, 4], 'Hs0X8-bf
h5GnEH0ySNiVwJQ': [0, 1, 1, 0, 0, 0, 0, 0, 0, 40.4418087, -79.9964699, 4.0, 4, 4], 'oYLENNKc6Q XKG9ixiMiumQ': [0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 33.4215567, -11
1.9169322, 4.0, 21, 4], 'FVQUcEfGEMAKvDdpfpr2EQ': [0, 0, 0, 0, 0, 0, 0, 0, 0, 45.551316, -73.5993917, 2.5, 6, 3], '4gcd2_d3h2eJmVUwUhmOFg': [0, 1, 1, 1, 0
, 1, 0, 0, 0, 1, 55.9490807, -3.186781, 5.0, 37, 4], 'MkQz_wPoaELPVyGhwaf7Ug': [0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 33.7150963, -112.1136988, 3.0, 88, 4], 'cdxkag8
rkxElCkLdAwEiKw': [0, 1, 0, 0, 0, 0, 0, 0, 0, 36.118221, -115.172532, 2.5, 12, 3], 'MaLTBw9QgxK8oeLRQFleTg': [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 36.1086896, -1
15.1731693, 2.0, 5, 3], 'tOHV1S6uhBxLB5Iy2C5hMQ': [0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 33.4649802, -112.0481577, 4.0, 27, 4], 'OK62awudOUWkwF3ivxHZNA': [0, 1, 1, 0
, 0, 0, 0, 0, 0, 0, 36.1010195, -115.0634182, 2.0, 7, 4], '2Ju9HrZABDcV2iEzhUfhmw': [0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 36.1277253, -115.1744147, 4.0, 10, 4], '9c
ANTUpowLcOXly53ThAwg': [0, 1, 0, 0, 0, 0, 0, 0, 0, 33.4217395, -111.883236, 3.5, 4, 4], 'o5PqKOR7IhyU8rlgr6HlrQ': [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 33.637747
6873, -112.422728051, 3.5, 35, 4], 'nNXdKc ertaZyDha-EF29Q': [0, 1, 1, 1, 0, 0, 0, 0, 1, 55.942879, -3.282731, 3.5, 3, 4], 'VnVvJzps8QkDbclGmHnmQ': [0, 1
, 1, 0, 1, 0, 0, 0, 1, 0, 36.1025618, -115.169277, 2.5, 185, 4], '9wckJZ9Uryl8-3WNcaf03Q': [0, 1, 1, 0, 1, 0, 0, 0, 0, 33.6992556, -111.8883259, 2.5, 15,
4], 'dTPyszCleiE7vBVOBmEcYQ': [0, 0, 1, 1, 0, 0, 0, 0, 0, 55.9027316, -3.2862754, 3.0, 6, 2], '7ELJAClk35_Qq7ZjooOyeA': [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 45.
4458773, -73.8523782, 3.5, 3, 8], 'Kcr-bbIASd426moPKvD7og': [0, 1, 1, 1, 0, 0, 0, 0, 0, 40.495421, -80.055926, 2.0, 16, 4], 'HF0J0ZGcxnW0vNvdR7pvtA': [0,
1, 1, 0, 1, 0, 0, 0, 0, 35.152911, -80.832917, 3.5, 4, 4], 'uGIUzFrKhuUKS7zV1Y-7VQ': [0, 1, 0, 0, 0, 0, 0, 0, 0, 33.4957219827, -112.000841742, 4.5, 13
7, 4], 'BluvXwJUMS51o3e3nRwXFQ': [0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 36.1226902, -115.1701939, 3.0, 77, 4], 'uxWYJozmg36_EddkdcisW': [0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 4.
1, 40.4396782, -79.9959316, 4.0, 52, 5], 'sg4mAXmmn3v8T8owLWHP-Q': [0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 36.1429443, -115.2265555, 4.5, 17, 6], 'ZeOmG_i-9ziEjN2gic
HqW': [0, 1, 1, 0, 0, 0, 0, 0, 0, 35.9992533, -115.1040028, 3.5, 21, 4], 'FzQvXhvo99pGllBoeHi2w': [0, 0, 0, 0, 0, 0, 0, 0, 0, 55.9507884, -3.1962074,
4.0, 3, 6], 'RGT2VvOTjNDHlCRKXBmgCw': [0, 0, 1, 1, 0, 0, 0, 0, 0, 45.4841067, -73.6284956, 3.5, 6, 3], 'ilhL7fbaY29HDwVypx8JCG': [0, 1, 1, 0, 0, 0, 0, 0, 0,
0, 36.1749234, -115.136982, 2.0, 6, 3], 'CKjcewWeWvdJ7TzOQb2OIw': [0, 1, 1, 0, 1, 0, 0, 0, 0, 33.6967596, -111.9160606, 5.0, 172, 4], 'qInDSpmV-FdTNj6
UuZvTg': [0, 1, 1, 1, 1, 0, 0, 0, 33.4796453, -112.074152, 3.0, 250, 4], 'JNTpW2zBftG1Vr7iieazQ': [0, 0, 0, 0, 0, 0, 0, 0, 0, 33.3791178, -112.0653
467, 3.0, 4, 3], 'VPFNPMOFQkoQL2CBzagiA': [0, 1, 1, 0, 1, 1, 0, 0, 1, 36.0843559, -115.279599, 4.0, 14, 5], 'crAE_fwUKeJH2dLTa6Ocxw': [0, 1, 0, 0, 0, 0, 0, 0, 0,
0, 0, 1, 40.117699, -88.24061, 3.5, 7, 5], 'xhSxibm86e6LB6ic9s3QsA': [0, 1, 0, 0, 0, 0, 0, 0, 0, 36.1107037, -115.2786276, 3.0, 36, 4], 'QhwkFogGQA-Ar1
76U15F0Q': [0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 36.1030612595, -115.173450261, 4.0, 1818, 4], 'EA6zWRssCMuzb3LDaZob5w': [0, 1, 0, 0, 0, 0, 0, 0, 0, 36.125825130
7, -115.226354905, 3.5, 7, 4], 'AzhRzeO08GKHv94ZengJw': [0, 1, 0, 0, 1, 0, 0, 0, 0, 33.5397085, -111.8897671, 3.0, 5, 5], 'uybZVw143IwaKAsRqk83Zg': [0, 1
, 1, 0, 1, 1, 0, 0, 0, 36.1598719621, -115.22466845, 4.5, 73, 4], 'Hn4X80Mughw7521mUipaw': [0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 49.0102472, 8.3849638, 3.5, 41
3], 'KXlIQIWMZSRKBEzxby2jDDA': [0, 1, 1, 0, 0, 1, 0, 0, 1, 33.6981334, -111.8891539, 3.5, 35, 4], 'lWWKIdzgYN3dkRYxk07SA': [0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 43.4680417,
-80.5229389, 4.5, 54, 4], 'I04AK0ljrh0j0Q0uLluJmW': [0, 1, 1, 0, 0, 0, 0, 0, 0, 33.7977043, -112.1187033, 2.0, 13, 3], 'uMxfxbqT27KDzU9UTXeE
rg': [0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 33.581757, -112.008454, 4.0, 382, 5], 'Np5XWA-SL9fpboELZJumeng': [0, 1, 1, 0, 0, 0, 0, 0, 0, 33.40724, -112.06422, 4.0,
7, 3], '4vjbfifZ-JexWlPVwmDAYw': [0, 1, 1, 0, 0, 0, 0, 0, 0, 33.4370178, -112.1846782, 4.0, 6, 4], 'OUyVyPKOhQWwq00lhc8_6Q': [0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 35.1374964,
-80.7390948, 3.5, 6, 4], 'mkXlYHuNudcucy-SSh5Ulg': [0, 0, 0, 0, 0, 0, 0, 0, 0, 33.058358, -112.0456421, 3.0, 10, 3], 'pBuTjTSnjz-r9-i 1518p
A': [0, 0, 1, 0, 0, 0, 0, 0, 0, 45.503697, -73.817233, 4.5, 4, 4], 'RUcgIIVxzDbarVNZIUO-Bg': [0, 1, 1, 0, 1, 0, 1, 0, 1, 32.881335, -111.754586, 4.5, 4
7, 4], 'dKwAVPE7Yvja8oI5r2NKD5w': [0, 0, 1, 1, 0, 0, 0, 0, 0, 33.465461, -112.017768, 3.5, 5, 4], 'crAGS1CyMRcBU1VxKcuwg': [0, 1, 1, 0, 0, 0, 0, 0, 0, 36.1135909,
-115.3053069, 2.5, 3, 3], 'jQ2xPF08NMUleQhZJKXabw': [0, 1, 1, 1, 0, 0, 0, 0, 0, 33.177871534, -111.585883168, 2.5, 31, 4], '8N_dJABpJW6VnJo
xRpcw': [0, 1, 0, 0, 0, 0, 0, 0, 0, 36.1029663086, -114.929458618, 3.5, 4, 4], '51gNc8sg9KCD5EJ3TUrww': [0, 1, 1, 0, 1, 0, 0, 0, 0, 33.6120764, -112.1
853986, 3.0, 4, 2], 'u2261KOMjUPDSrUK90Emq': [0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 36.1122889, -115.297923, 3.5, 9, 4]
----- Predicting new businesses and Writing -----
done

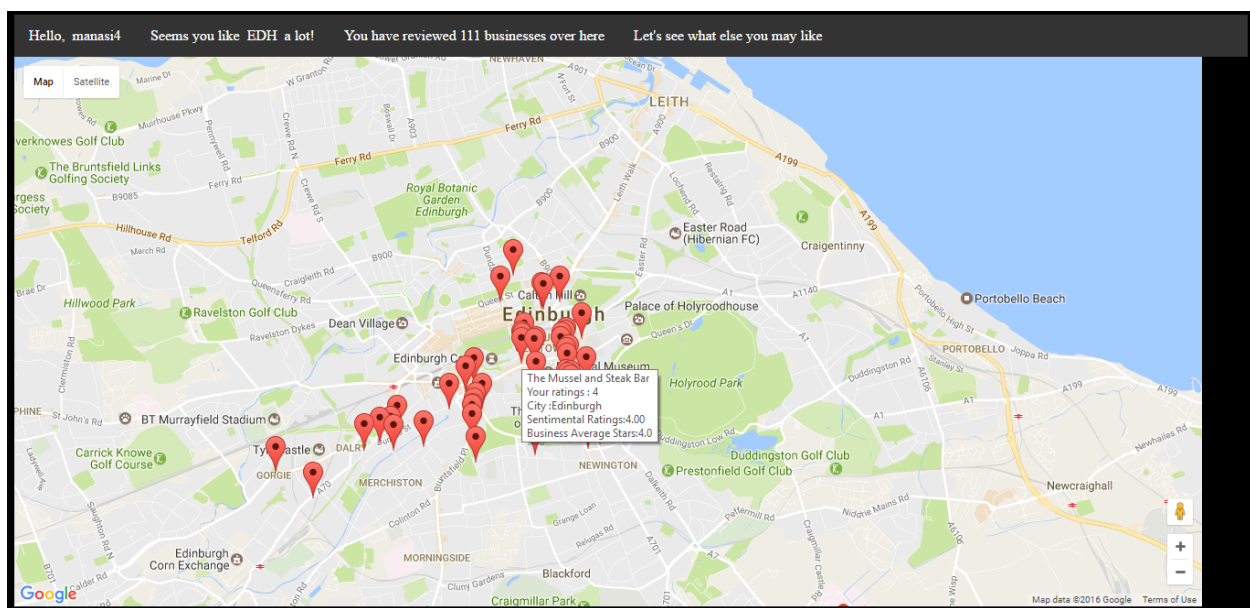
real    52m29.972s
user    52m19.020s
sys     0m9.920s
```

PuTTY Fatal Error

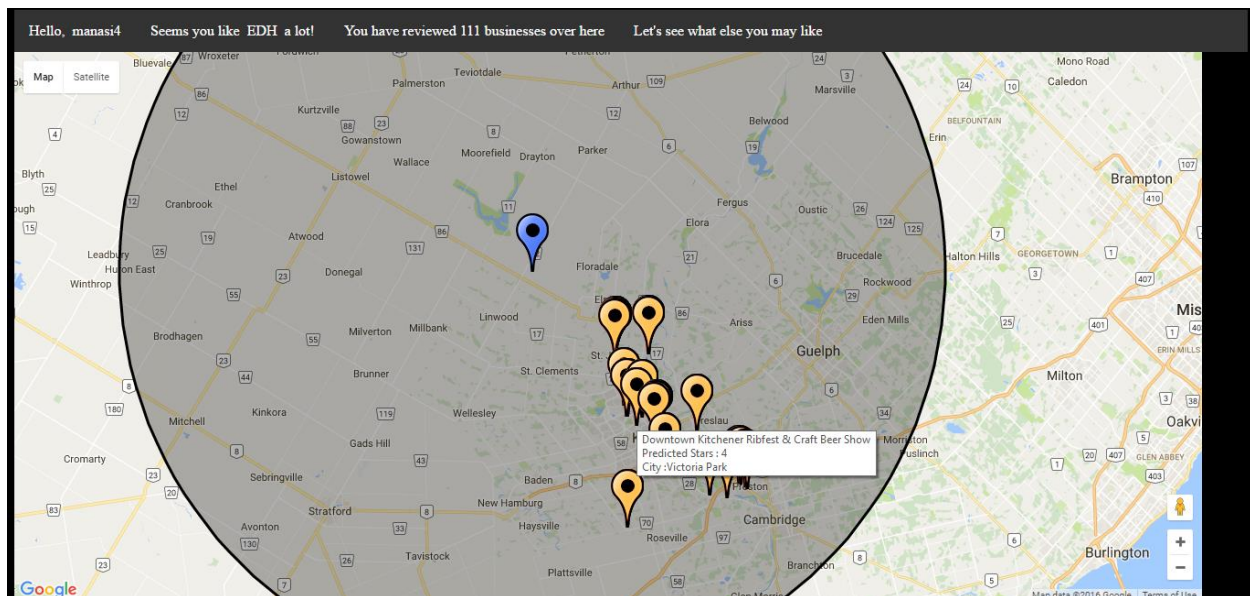
# Chapter 5

## 5. Application

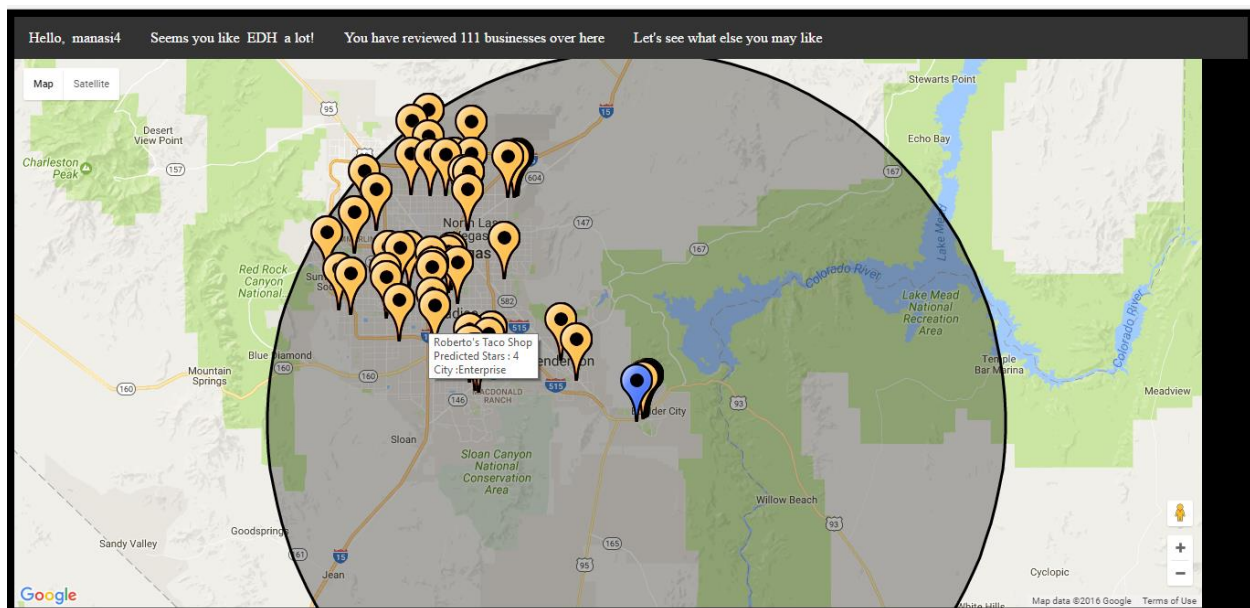
- Consider a user who had reviewed the food outlets in his region in Edinburgh, UK
- When the user logs in, he sees his rated businesses along with the sentiment score of his reviews and business average ratings
- This region is termed as hot spot of the user
- Google Maps API and Geometry Library were used



- When the user moves to different location, say Cambridge in US, he receives the predictions for the food businesses he has not reviewed yet
- These predictions will be based on businesses that he has reviewed in Edinburgh, UK



- The yellow markers show the new businesses
- Similarly, when he moves his location to Las Vegas, he sees a different prediction model



# Chapter 6

## 6. Challenges faced

- **Selecting features:**

Since there were 95 features for the businesses, it was a challenge to select the relevant features to predict with high accuracy. To overcome this problem, we analyzed each attribute and eliminated the features that had a lot of Null values and which were not relevant for Food related businesses.

- **Loading the data**

We first faced this issue while loading the reviews data feed into R for data analysis since it was around 2gb. We tried increasing the heap memory of R but still it was not able to process. To overcome this problem we used the *fread* function of *data.table* library which is meant to import big data from regular delimited files directly into R, without any detours.

- **Processing the review data for sentiment analysis**

To perform the sentiment analysis on each and every word of all the reviews, we developed a R code which was not able to process because of memory limits. Therefore, we implemented the sentiment analysis algorithm on Hadoop MapReduce on AWS EC2 instance.

- **Data standardization**

Since the algorithm that we implemented assumes that all features are centered around zero and have variance in the same order, we were getting lower accuracy. To resolve this issue we used the *preprocessing* module of *scikit-learn* library for mean removal and variance scaling.

# Chapter 7

## 7. Tools & Technologies Used

Technology	IDE/ Framework
Python	Spyder
R	R Studio
MapReduce	Eclipse/ Hadoop
Spring MVC	STS
Hadoop cluster	AWS EC2 instance
MongoDB	Mongo 3T Chef

# Chapter 8

## 8. Conclusion

We were able to build a prediction model which is able to help a yelper by predicting the ratings for the new businesses at a new location based on his/her historical data. This way we were able to mirror his/her personal likings on a new geographical location and were able to map his model.



# Chapter 9

## 9. Future scope

- All Business Categories  
We can use this model to predict the user rating for all businesses which are a part of the data set provided by Yelp
- Social Graph Mining  
We can analyze and understand the yelping pattern that can calculate the degree of similarity between the user and his friends. By this we can figure out which user is a trend setter for a particular business.
- #YELFIE  
According to the current '#Yelfie' trend we can analyze the sentiment of a user by image recognition which can be included as a feature in our current data set.

# Chapter 10

## 10. Work Distribution

Member Name	Work Responsibility/ Distribution
Bhavna, Sarthak, Manasi	Research applicable Machine Learning Algorithms, Frameworks, Libraries and Technical Papers
Kunal, Swarna	Featuring Engineering – Data cleaning, data conversation and feature analysis
Bhavna, Sarthak, Manasi, Kunal, Swarna	Designing our Predictive Model Architecture of data analysis
Bhavna, Kunal, Sarthak	Setting up the cloud cluster
Manasi, Swarna, Sarthak	Developing and Testing the neural networks and machine learning techniques used in the model
Bhavna, Kunal, Manasi, Sarthak	Deploy the solution on code and tune the accuracy
Swarna, Manasi, Bhavna, Sarthak, Kunal	Virtualization of the data, documentation and presentation



# Chapter 11

## 11. References

1. Data Analytics using Yelp Data [Nevil Patel, Suraj Ponugoti, Doan H Nguyen]
2. LIBSVM: A Library for Support Vector Machines [Chih-Chung Chang and Chih-Jen Lin, Department of Computer Science, National Taiwan University, Taipei, Taiwan]
3. Applications of Machine Learning to Predict Yelp [Kyle Carbon, Kacyn Fujii, Prasanth Veerina]