

# DSS Capstone - Connecting Users to Build a Social Network In Yelp!

*Usamah Khan*

*November 20, 2015*

## Introduction

As part of completing the Coursera Data Science Specialization, we were given a Yelp! dataset and tasked with exploring the data, determining a question and building a model or product associated with it. This report outlines the process undertaken and the methods applied.

The Dataset was supplied by Yelp from this address:

[https://d396qusza40orc.cloudfront.net/dsscaphone/dataset/yelp\\_dataset\\_challenge\\_academic\\_dataset.zip](https://d396qusza40orc.cloudfront.net/dsscaphone/dataset/yelp_dataset_challenge_academic_dataset.zip)

The folder contains 5 datasets of type .JSON and a PDF detailing terms of use. The datasets comprise the following information (Note - the names of the datasets as they were read in to R follow the names):

Business -> *business.data*

Check-in -> *checkin.data*

Review -> *review.data*

Tip -> *tip.data*

User -> *user.data*

After some exploratory analysis, performed in part for the Capstone quiz and for determining interesting observations in the data, the following question was isolated with the rationale below:

*Can Yelp be used to create a social network by recommending friends to you based on location, check-in, ratings and other factors?*

People are inclined to trust friends for recommendations and like to meet new people to go and eat with. This can be done by following those who you believe have good reviews but if there were a passive way to find recommendations like on Instagram's main feed users may find it interesting. The feedback that was recieved was positive but with one reviewer of five stating that they might not have a need for it. The question was undertaken and was found to be possible with certain limitations.

## Methods and Data

### Part I - Manipulating and Extracting the Data

Since the question is primarily about geographic data, the first thing to do was to determine how the data was too be isolated to a specific zone. The process for doing required isolating business ids in each of the 10 cities and link them to user data through reviews. i.e

Business Data -> Review Data -> User Data -> Isolated Users

The data for businesses was dirty and hence had to be cleaned since there were more cities and states than were actually possible. A method of K-means clustering was employed to accurately group businesses. Initial guesses were made based on longitudes and latitudes from the *geocode* package. With those extracted,

*NOTE - All code chunks presented only contains relevent snippets. No cleaning and intermediate steps are shown but can all be found for further inspection over at the project repository on GitHub*

```

coordinates <- as.data.table(cbind(business.data$longitude, business.data$latitude))

cities <- c('Edinburgh, UK', 'Karlsruhe, Germany', 'Montreal, Canada', 'Waterloo, Canada',
           'Pittsburgh, PA', 'Charlotte, NC', 'Urbana-Champaign, IL', 'Phoenix, AZ',
           'Las Vegas, NV', 'Madison, WI')

city.center <- geocode(cities) # Get initial guesses

city.clusters <- kmeans(coordinates, city.center) # Perform K-means

coordinates.clusters <- cbind(city.clusters$cluster, coordinates) # Get indexed cluster
                                                                # info

business.data <- cbind(coordinates.clusters$City, business.data) # Bind with business data

```

Now that all the data was separated it was determined that the scope of the problem would be related to a single city. The city of Montreal was isolated for this analysis. Hence Montreal Business Data was subset:

```

business.data.mtl <- subset(business.data, business.data$City == "Montreal, Canada")
business.data.mtl <- business.data.mtl[,c("City", "business_id", "longitude", "latitude")]

```

With this step completed it became relatively straightforward to get an inferred dataset of users based in Montreal. The code and the steps are outlined below:

```

# Subset from reviews, all businesses reviewed in Montreal

review.data.mtl <- review.data[which(review.data$business_id %in% business.data.mtl$
                                   business_id),]

# Subset from review.data.mtl, all users - Hence all users based out of Montreal

user.data.mtl <- user.data[which(user.data$user_id %in% review.data.mtl$user_id),]

```

This reduced the size of the *user.data* set down to 13861 from 366715. A more manageable size.

Determining a “home” location for users was the next step. Since no data exists for their home address and location (rightly so) an inference was made based on visited locations. The idea was to find the center of mass of all visited restaurants in the Montreal region by a user and assume that to be where the user likely spends most of their time at businesses and restaurants. A calculation on the mean of all the longitude and latitudes of businesses visited by a user was performed.

```

for(i in 1:nrow(user.data.mtl)) # Loop to isolate all businesses visited by User
{
  # to determine center of mass
  user.id <- subset(review.user.mtl, review.user.mtl$user_id == user.data.mtl[i,7])
  mx <- mean(user.id$longitude)
  my <- mean(user.id$latitude)

  user.id.location[i,1] <- user.data.mtl[i,7]
  user.id.location[i,2] <- mx
  user.id.location[i,3] <- my
}

user.id.location <- merge(user.data.mtl[,c(6,7,5,214,215)], user.id.location, by =
                        "user_id") # Add relevant data for users

```

## Part II - Building a Data product

With all the relevant data isolated, the next step was to determine how one could connect users. For this, product research was done on how, if this add-on existed, would it be used. The initial idea was to create a feed like on Instagram, however, it became clear that people prefer to have a visual resource to find friends, more like Tinder. So, a data product was built in shiny with extensive use of leaflet to give an example of how the data could be used.

The idea was plot all data for the centroids of user activity, and from there give users an option to find other users within a specific radius.

To do this a UI.R and Server.R were written to include an interactive map with reactive inputs to users and distance. When clicked, the markers display popups with the following information: Name, Average Stars, Number of Reviews and Number of Fans. Enough information to be of interest to a user wanting to connect.

The second aspect of the methods was to allow further subsetting to the data based on checking all mutually reviewed businesses and star ratings. This was to be achieved by calling a user and isolating businesses reviewed and then referencing that with all users in the specific zone of reference. The app would then drop other users with little relevance to the user in question.

## Results

The purpose of this Capstone was to build a model, statistical or predictive, yet since the final course dealt with Developing Data Products, a lot of time was taken to build a tangible, dynamic environment for exploring the data in addition to a model. As a result, the result of the analysis came in the form of an app that can be accessed at ——— and played around with.

It was useful and showed a possible solution to the initial question laid out, however, with some caveats. First, to understand how the app works, the screenshot below shows an example of the usage and the zone of interest. The final data frame that was cleaned and used as input for all information on the map is in the following format:

*Note - this is a test set used for display:*

```
##           user_id
## 1 __0s_Oby_akeEKnkgTrz4w
## 2 __3cGi5oLgTK4EkFwJETSg
## 3 __53csc1JFS0I0lpxZXn_g
## 4 _pw6eL7Vgku-DssJL_g-A
## 5 _RoUSf53-Nw9qPgeyLpng
## 6 _wI5M_GcBLFhsExd474fw

##      name review_count fans average_stars longitude latitude
## 1  Roland           18    0           3.67 -73.56573 45.50711
## 2 Michael           34    0           3.09 -73.56316 45.50948
## 3    Ali             6    0            3.6 -73.57141 45.50842
## 4   Neil            40    2           3.65 -73.50997 45.53848
## 5  Nicole             2    0            4 -73.57439 45.52221
## 6   Aaron            10    0           4.4 -73.58080 45.52724
```

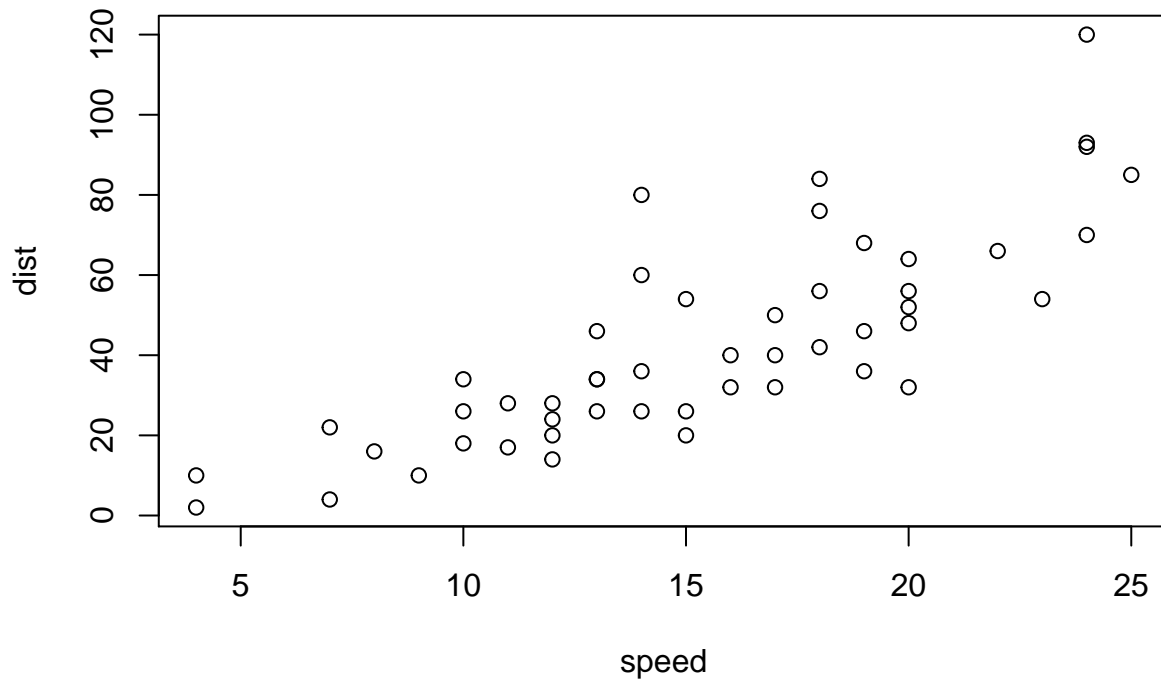
The



```
summary(cars)
```

```
##      speed      dist
##  Min.   : 4.0    Min.    :  2.00
## 1st Qu.:12.0    1st Qu.: 26.00
## Median :15.0    Median : 36.00
## Mean   :15.4    Mean    : 42.98
## 3rd Qu.:19.0    3rd Qu.: 56.00
## Max.   :25.0    Max.    :120.00
```

You can also embed plots, for example:



Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.