

Branch: master ▼

Find file

Copy path

## BzztPodtaxi / FinalBzzt.R



**aicasas** Add files via upload

b5b3acc 21 minutes ago

1 contributor

Raw

Blame

History



206 lines (165 sloc) | 7.2 KB

```
1
2 'Bzzt Final Project: Alexis Casas
3 Lines 5:110 are data cleaning, lines 110:133 are geocoding of ~26,00
4 saved as BzztGeocode.csv and BzztGeocodeD.csv
5 and can be ran quickly from lines 139 and on'
6
7 data1 <- read.csv('201901.csv', header=FALSE)
8 data2 <- read.csv('201902.csv', header=FALSE)
9 data3 <- read.csv('201903.csv', header=FALSE)
10
11 #Combining Datasets by columns
12 data <- rbind(data1, data2, data3)
13 rm(data1, data2, data3)
14
15 #Deleting rows with missing entries
16 data[data==""] <- NA
17 data[data=="None"] <- NA
18 which(is.na(data))
19 data<- data[complete.cases(data),]
20 which(is.na(data))
21
22 #install.packages("stringr")
23 library(stringr)
```

```

24
25 #Data cleaning
26 #Split the first column by ; for time stamp at different stops
27 data <- data.frame(data,do.call(rbind,str_split(data$V1,";")))
28 data <- data.frame(data,do.call(rbind,str_split(data$V5,";")))
29 data <- data.frame(data,do.call(rbind,str_split(data$V3,";")))
30 #Deleting first columns that were split
31 data <- data[,-c(1, 3, 5, 14, 18)]
32
33 #Filling in column datasets that were not made included
34 names(data)[10] <- "OriginAddress"
35 names(data)[14] <- "DestinationAddress"
36 names(data)[3] <- "Code"
37 names(data)[4] <- "ActualOriginTS"
38 names(data)[5] <- "ActualDestinationTS"
39 names(data)[6] <- "EstimatedOriginTS"
40 names(data)[7] <- "EstimatedDestinationTS"
41 names(data)[8] <- "RideRequestTS"
42 names(data)[9] <- "RideAcceptedTS"
43 names(data)[1] <- "EstimatedEndAddress"
44 names(data)[11] <- "ZipCode"
45 names(data)[12] <- "ActualCost"
46 names(data)[13] <- "EstimatedCost"
47 names(data)[2] <- "EstimatedOriginAddress"
48
49 #Taking care of columns entries which do not match the others (incco
50 data <- data[-grep("Stockholm", data$ZipCode), ]
51 data <- data[grep("[0-9]", data$EstimatedEndAddress), ]
52 data <- data[grep("[0-9]", data$EstimatedOriginAddress), ]
53
54
55 #Data cleaning for Date and Time Format
56
57 data$ActualOriginTS <- sub('([^\s]+).*', '\\1', data$ActualOriginTS)
58 data$ActualDestinationTS <- sub('([^\s]+).*', '\\1', data$ActualDesti
59 data$EstimatedOriginTS <- sub('([^\s]+).*', '\\1', data$EstimatedOrig
60 data$EstimatedDestinationTS <- sub('([^\s]+).*', '\\1', data$Estimate
61 data$RideRequestTS <- sub('([^\s]+).*', '\\1', data$RideRequestTS)

```

```

62 data$RideAcceptedTS <- sub('([^\+]+\).*', '\\1', data$RideAcceptedTS)
63
64 #Filling in NA
65 data[data==""] <- NA
66 data[data=="None"] <- NA
67 which(is.na(data))
68 data<- data[complete.cases(data),]
69
70 library(lubridate)
71 library(dplyr)
72
73 data$ActualOriginTS <- strptime(data$ActualOriginTS, format="%Y-%m-%d")
74 data$ActualDestinationTS <- strptime(data$ActualDestinationTS, format="%Y-%m-%d")
75 data$EstimatedOriginTS <- strptime(data$EstimatedOriginTS, format="%Y-%m-%d")
76 data$EstimatedDestinationTS <- strptime(data$EstimatedDestinationTS, format="%Y-%m-%d")
77 data$RideRequestTS <- strptime(data$RideRequestTS, format="%Y-%m-%dT%H:%M:%S")
78 data$RideAcceptedTS <- strptime(data$RideAcceptedTS, format="%Y-%m-%dT%H:%M:%S")
79
80
81 #Encoding Weekdays and Time to see which are the bussiest
82 data$Weekday <- (weekdays(data$ActualOriginTS))
83 data$Hour <- hour(data$ActualDestinationTS)
84 data$Month <- month(data$ActualDestinationTS)
85
86 #Plotting Number of rides taken on each day
87 data$Weekday <- factor(data$Weekday, levels=c("Sunday", "Monday", "Tuesday", "Wednesday", "Thursday", "Friday", "Saturday"))
88 freq_table <- table(data$Weekday)
89 margin.table(freq_table, 1)
90 plot(freq_table)
91 barplot(freq_table, main="Daily Ride Frequency",
92         xlab="Day of Week")
93
94 #Plotting the frequency density at different hours of the day
95 install.packages('ggpubr')
96 library(ggpubr)
97 d <- density(data$Hour)
98 plot(d, main="Ride Density: Hour of Day")
99 polygon(d, col="grey", border="blue")

```

```
100
101
102 freq_table2 <- table(data$Hour)
103 margin.table(freq_table2, 1)
104 plot(freq_table2)
105 #Plot shows that there is a peak around 8am and 5-6pm (commuter hour)
106
107 #####Applying K means to actual vs estimated cost
108
109
110
111 #####Here I begin Geocoding:#####
112 #There are too many addresses to generate 26,000 different geocoordi
113 #I will look at just months 1:3 of data, specifically month 1 looks
114 sample <- data[grep("[1:3]", data$Month), ]
115
116 #install.packages("ggplot2")
117 library(ggplot2)
118 library(ggmap)
119
120
121 #Getting API
122 api <- "AIzaSyBFmuXtB3JzjnCB0M5op0yY4rcMfJ1wL5k" # Text file with th
123 register_google(key = api)
124 #concatenate the address
125 start_addy <- paste(sample$OriginAddress, "Stockholm, Sweden", sampl
126 end_addy <- paste(sample$DestinationAddress, "Stockholm, Sweden", sa
127 # geocode - check for warnings
128 addys_coords <- geocode(start_addy)
129 end_addy_coords <- geocode(end_addy)
130 sample <- data.frame(cbind(sample, addys_coords))
131 sampleD <- data.frame(cbind(sample, end_addy_coords))
132 #write.csv(sample, file = "BzztGeocode.csv")
133 #write.csv(sampleD, file = "BzztGeocodeD.csv")
134 #####Here is the start of k means with coordinates#####
135
136
```

```

137 #Cleaned data for Kmeans
138
139 #Kmeans for pick up location
140 bzzt <- read.csv("BzztGeocode.csv")
141 sample <- bzzt[,-c(1:16)]
142 sample$Weekday <- factor(sample$Weekday, levels=c("Sunday","Monday",
143 sample <- sample[,4:5]
144 which(is.na(sample))
145 sample<- sample[complete.cases(sample),]
146 #removing outlier
147 #sample <- sample[-grep("-[0-9]", sample$lon), ]
148 #The elbow mehtod will help determine the number of clusters. This i
149 wcss = vector() #within cluster, sum of square. finding the closenes
150 for (i in 1:10) wcss[i]= sum(kmeans(sample, i)$withinss)
151 plot(1:10, wcss, type='b', main=paste('The Elbow Method: Pick Up Loc
152       xlab = 'Number of Clusters',
153       ylab = 'WCSS')
154
155 #don't need feature scaling
156
157 #apply k-means
158 set.seed(42)
159 cluster <- kmeans(sample, 2)
160 sample$Borough <- factor(cluster$cluster)
161 ggplot(sample, aes(x=lat, y=lon, color=Borough)) + geom_point(shape=
162 #+ coord_cartesian(xlim = c(59.2, 59.4), ylim = c(17.8,18.3))
163 #+ coord_cartesian(xlim = c(57.2, 59.4), ylim = c(16,19))
164 #+ coord_cartesian(xlim = c(59.2, 59.4), ylim = c(17.8,18.3))
165 ggsave("pickupKmeans.png")
166
167 freq_table <- table(sample$Borough)
168 margin.table(freq_table, 1)
169
170 #install.packages('RgoogleMaps')
171 library(RgoogleMaps)
172 #data(sample)
173 col=as.numeric(sample$Borough)

```

```

174 par(pty="s")
175 plotmap(lat, lon, zoom = 13, col = col, pch=1, data = sample)
176
177 #Kmeans for drop off location
178 bzztD <- read.csv("BzztGeocodeD.csv")
179 sampleD <- bzztD[, -c(1:16)]
180 sampleD <- sampleD[, 4:5]
181 which(is.na(sampleD))
182 sampleD <- sampleD[complete.cases(sampleD),]
183
184 #The elbow method will help determine the number of clusters. This is
185 wcss = vector() #within cluster, sum of square. finding the closeness
186 for (i in 1:10) wcss[i] = sum(kmeans(sampleD, i)$withinss)
187 plot(1:10, wcss, type='b', main=paste('The Elbow Method: Drop Off Lo
188      xlab = 'Number of Clusters',
189      ylab = 'WCSS')
190 #don't need feature scaling
191
192 #apply k-means
193 set.seed(42)
194 cluster <- kmeans(sampleD, 2)
195 sampleD$Borough <- factor(cluster$cluster)
196 ggplot(sampleD, aes(x=lat, y=lon, color=Borough)) + geom_point(shape=1)
197 ggsave("dropoffKmeans.png")
198
199 freq_table <- table(sampleD$Borough)
200 margin.table(freq_table, 1)
201
202 col=as.numeric(sampleD$Borough)
203 par(pty="s")
204 plotmap(lat, lon, zoom = 13, col = col, pch=1, data = sampleD)
205
206

```