

AIDM7390 Data Mining and Knowledge Discovery for Digital Media
Group Project
Twitter Data Analysis

Front page:

AIDM7390

Title

100 Days Of Code

Student IDs

20465106 CHEN Xiaoqi (leader)

20426550 CAI Runlin

20465769 GUO Yuju

20449496 HUANG Zefei

1. Introduction:

AI has become a hot topic in society. More and more people pay attention to learning programming, and even self-learning code, hoping to hold the opportunity of computer science trend to improve their thinking ability and find a high-paying career. We hope to explore what programming are those people who are interested in computer science learning or what programming they want to learn, and what specific content they are interested in. So we can judge the application scenarios of computer science and the current career direction. This is also closely related to our profession.

We choose “100DaysOfCode” as our hashtag. Actually, it includes two rules, like people need to code minimum an hour every day for the next 100 days, and tweet your progress every day with the #100DaysOfCode hashtag. This topic is like a large community to share code communication, programming discussion. Many people check in on this hashtag and share their mood of code learning every day. Therefore, there is a lot of data for us to observe and analyze which part of the code people like and what specific content they are discussing.

2. Procedures:

- After we got the API, we first set the app with the name “zooloretto” and then got the access token.
- In R, we installed a “rtweet” package so that R can send requests to Twitter, then automatically retrieved the data and organize the returned data into tidy structures.
- We set the hashtag of “#100DaysOfCode” to 10000 visits, excluding retweets, and users using English as their twitter content.
- We used the access token in R to get the data we want, and we can see in the table what content the returned tweets contain.
- We use the “ggplot2” package to visualize the number and time of tweets in the past few days, and count them every 3 hours, so we can observe the #100DaysOfCode activity. And we plot frequency of tweets for five related users and all tweet statuses under this hashtag.
- We installed several packages like “tm”, “SnowballC”, “wordcloud”, “wordcloud2” and “RColorBrewer” for text mining and word cloud.
- We used “RSentiment” packages and “syuzhet” packages for sentiments analysis. We can

get the emotion of the tweet, such as positivity, enjoyment, sadness, depression, etc.

- We used Latent Dirichlet Allocation analysis was aiming at detecting sub-topics from the general debate.
- We also installed several packages like “ggmap” packages and “maps” packages for visualizing. We added some geographic information to the twitter data using the “rtweet::lat_lng()” function, and then we used the “ggplot2::map_data()” function to get the “world” data. Finally we layered the twitter data onto the map with “ggplot2::geom_point()” by specifying the *long* and *lat* to *x* and *y*. We can see where twitter users are on the map and which countries they belong to, so we can see which countries are more interested in the topic

```
1 install.packages('rtweet')
2 library('rtweet')
3
4 token <- create_token(app = 'Zooloretto',
5                       consumer_key = 'BpRSq1Toabc2RsB0lRV3HqraB',
6                       consumer_secret = 'NdSZNMgcSS23lvxv4zHPWaqmp04DptBLh777bK1c7wbVTDpu1',
7                       access_token = '1327201654018490368-IuksF3YAZmMCIn087m5iqvdu4Mnjdm',
8                       access_secret = '7hA7xB6N49FcBmv8ux2ui0aQqilQZteikrjZos1Pwvfkf',
9                       set_renv = TRUE)
10
11 Code <- search_tweets('#100DaysOfCode', n=10000, include_rt = FALSE, lang='en')
12 Code
13
14
15
16
17
18
19
20
21 #####
22 # Information of official account of #100DaysOfCode #
23 #####
24
25 daysoc <- lookup_users('_100DaysOfCode')
26 daysoc$name
27 daysoc$description
28 daysoc$followers_count
29
30 #####
31 # search tweet #
32 #####
33
34 data.frame(Code)
35 colnames(Code)
36 Code[100,]$text
37 Code[100,]$screen_name
38 Code[100,]$created_at
39 Code[100,]$retweet_count
40
41 #####
42 # Twitter rate limits #
43 #####
44
45 LCode <- search_tweets('#100DaysOfCode', n = 1000000, retryonratelimit = TRUE)
46 LCode
47
48
49
50
51
52 install.packages('ggplot2')
53 library('ggplot2')
54
55 ts_plot(Code, '3 hours') +
56   ggplot2::theme_minimal() +
57   ggplot2::theme(plot.title = ggplot2::element_text(face = 'bold')) +
58   ggplot2::labs(
59     x = NULL, y = NULL,
60     title = "Frequency of #100DaysOfCode Twitter statuses from past 4 days",
61     subtitle = "Twitter status (tweet) counts aggregated using three-hour intervals",
62     caption = "\nSource: Data collected from Twitter's REST API via rtweet"
63   )
64
65 #####
66 # Get 5 users timelines #
67 #####
68
69 codeuser <- get_timelines(c("#100DaysOfCode", "kallaway", "amanhimself", "ossia", "FreeCodeCamp"), n = 3200)
70
71 codeuser %>%
72   dplyr::filter(created_at > "2020-11-1") %>%
73   dplyr::group_by(screen_name) %>%
74   ts_plot("days", trim = 1L) +
75   ggplot2::geom_point() +
76   ggplot2::theme_minimal() +
77   ggplot2::theme(
78     legend.title = ggplot2::element_blank(),
79     legend.position = "bottom",
80     plot.title = ggplot2::element_text(face='bold')) +
81   ggplot2::labs(
82     x=NULL, y=NULL,
83     title="Frequency of Twitter statuses posted by 100DaysOfCode and following",
84     subtitle="Twitters status (tweet) count aggregated by day from Nov 2020",
85     caption="\nSource: Data collected from Twitter's REST API via rtweet"
86   )
87
```

```

88 # #####
89 # Text mining and word cloud #
90 # #####
91
92 install.packages('tm')
93 install.packages('SnowballC')
94 install.packages('wordcloud')
95 install.packages('RColorBrewer')
96
97 library('tm')
98 library('SnowballC')
99 library('wordcloud')
100 library('RColorBrewer')
101
102 Codecode.v <- VectorSource(Code$text)
103 Codecode.c <- SimpleCorpus(Codecode.v)
104
105 inspect(Codecode.c)
106
107 Codecode.c.p <- tm_map(Codecode.c, content_transformer(tolower))
108 Codecode.c.p <- tm_map(Codecode.c.p, removeNumbers)
109 Codecode.c.p <- tm_map(Codecode.c.p, removeWords, stopwords('english'))
110 Codecode.c.p <- tm_map(Codecode.c.p, removeWords, c("day"))
111 Codecode.c.p <- tm_map(Codecode.c.p, removePunctuation)
112 Codecode.c.p <- tm_map(Codecode.c.p, stripwhitespace)
113
114 inspect(Codecode.c.p)
115
116 dtm <- TermDocumentMatrix(Codecode.c.p)
117 m <- as.matrix(dtm)
118 v <- sort(rowSums(m), decreasing = TRUE)
119 d <- data.frame(word = names(v), freq = v)
120 head(d,10)
121
122 set.seed(1234)
123 wordcloud(words = d$word, freq = d$freq, min.freq = 1,
124           max.words = 200, random.order = FALSE, random.color = TRUE, rot.per = 0.25,
125           colors = brewer.pal(13,"Paired"))
126
127 # #####
128 # word cloud 2 #
129 # #####
130
131 install.packages("wordcloud2")
132 library("wordcloud2")
133 wordcloud2(d,size = 1, color = "random-light",shape = "triangle-forward")
134
135 # #####
136 # count #
137 # #####
138
139 top10 <- head(d, 10)
140 top10
141 barplot(freq ~ word, data = top10, width =2,border = NA,las=2,
142        main = "Top 10 most frequent words",cex.main=1,col = terrain.colors(10))
143
144 # #####
145 # network #
146 # #####
147
148 install.packages("topicmodels")
149 install.packages("lubridate")
150 install.packages("SentimentAnalysis")
151 install.packages("ggpubr")
152 install.packages("dplyr")
153 install.packages("tidytext")
154 install.packages("quanteda")
155 install.packages('textdata')
156 library("topicmodels")
157 library("lubridate")
158 library("SentimentAnalysis")
159 library("ggpubr")
160 library("dplyr")
161 library("tidytext")
162 library("quanteda")
163 library('textdata')
164
165 text_code <- Codecode.c.p
166 text_df <- data.frame(text_clean = get("content",text_code),stringsAsFactors = FALSE)
167 Code$text <- text_df$text_clean
168 toks <- tokens(Code$text)
169
170 set.seed(30)
171 fcmat <- fcm(toks, context = "document", tri = FALSE)
172 feat <- names(topFeatures(fcmat, 30))
173 fcm_select(fcmat, pattern = feat) %>%
174 textplot_network(min_freq = 0.5)
175

```

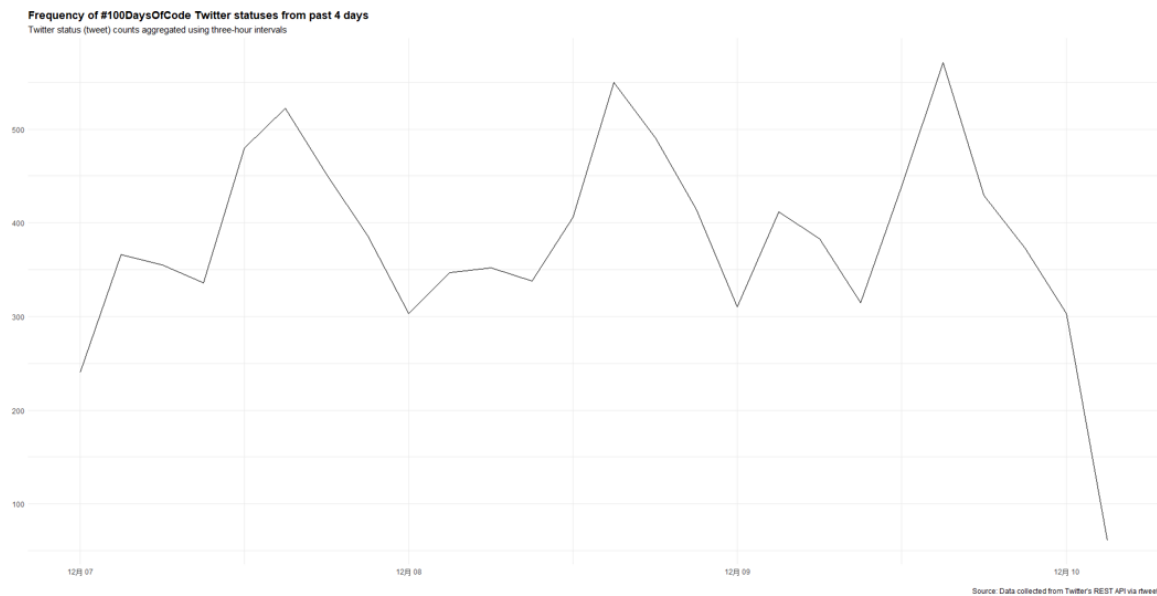
```

176 - #####
177 # Topic modelling with LDA #
178 - #####
179
180 set.seed(100)
181 samp <- sample(nrow(Code), 3500)
182
183 corpus_sub <- Corpus(VectorSource(Code$text[samp]))
184 dtm_sub = DocumentTermMatrix(corpus_sub)
185 doc.length = apply(dtm_sub, 1, sum)
186 dtm_sub = dtm_sub[doc.length > 0,]
187
188 k <- 3
189 DO <- LDA(dtm_sub, k, method = "Gibbs", control = list(nstart = 5, seed = list(2003, 5, 63, 100001, 765),
190 best = TRUE, burnin = 4000, iter = 2000, thin = 400))
191
192 topics <- as.matrix(topics(DO))
193 terms <- as.matrix(terms(DO, 10))
194 topics_prob <- as.matrix(DO@gamma)
195
196 topics_beta <- tidy(DO, matrix = "beta")
197
198 top_terms_b <- topics_beta %>%
199   group_by(topic) %>%
200   top_n(10, beta) %>%
201   ungroup() %>%
202   arrange(topic, -beta)
203
204 theme_set(theme_classic())
205 top_terms_b %>%
206   mutate(term = reorder(term, beta)) %>%
207   ggplot(aes(term, beta, fill = factor(topic))) + labs(x = 'words', y = NULL) +
208   geom_col(show.legend = TRUE) + facet_wrap(~ topic, scales = "free") +
209   theme(axis.text = element_text(angle = 30, vjust = 0.5, size = 8)) + coord_flip()
210
211 -
212 # Sentiment analysis #
213 - #####
214
215 install.packages('ROAuth')
216 install.packages('syuzhet')
217 install.packages('Rsentiment')
218 install.packages('kableExtra')
219 install.packages('knitr')
220 install.packages('RColorBrewer')
221 library('ROAuth')
222 library('syuzhet')
223 library('Rsentiment')
224 library('kableExtra')
225 library('knitr')
226 library('RColorBrewer')
227
228 mysentiment_code <- get_nrc_sentiment((Code$text))
229 Sentimentscores_code <- data.frame(colSums(mysentiment_code[,]))
230 names(Sentimentscores_code) <- 'Score'
231 Sentimentscores_code <- cbind('sentiment'=rownames(Sentimentscores_code), Sentimentscores_code)
232 rownames(Sentimentscores_code) <- NULL
233
234 ggplot(data = Sentimentscores_code, aes(x=sentiment, y=Score)) +
235   geom_bar(aes(fill=sentiment), stat = 'identity', width=0.7) + theme(legend.position = 'right') +
236   xlab('Sentiments') + ylab('Scores') + ggtitle('Sentiments of people behind the tweets on #100DaysOfCode')
237
238 -
239 #####
240 # Map #
241 - #####
242
243 install.packages('ggmap')
244 install.packages('maps')
245 install.packages('mapdata')
246 install.packages('igraph')
247 install.packages('gganimate')
248 install.packages('ggraph')
249 install.packages('ggalt')
250 install.packages('ggthemes')
251 library('ggmap')
252 library('maps')
253 library('mapdata')
254 library('igraph')
255 library('gganimate')
256 library('ggraph')
257 library('ggalt')
258 library('ggthemes')
259
260 CodeLoc <- rtweet::lat_lng(Code)
261 CodeLoc %>% names() %>% tail(2)
262 CodeLoc %>% dplyr::distinct(lng) %>% base::nrow()
263 CodeLoc %>% dplyr::distinct(lat) %>% base::nrow()
264 CodeLoc <- CodeLoc %>% dplyr::rename(long = lng)
265
266 world <- ggplot2::map_data('world')
267 world %>% glimpse(1000)
268
269 ggworldMap <- ggplot2::ggplot() +
270   ggplot2::geom_polygon(data = world,
271     aes(x = long,
272         y = lat,
273         group = group),
274     fill = "grey82",
275     color = "white",
276     alpha = 0.6)
277
278 gg_Code_title <- "#100DaysOfCode tweets worldwide"
279 gg_Code_cap <- "Tweets collected with the hashtags #100DaysOfCode"
280
281 gg_Code_map <- ggworldMap +
282   ggplot2::coord_quickmap() +
283   ggplot2::geom_point(data = CodeLoc,
284     aes(x = long, y = lat),
285     size = 0.7, # reduce size of points
286     color = "firebrick") + ggplot2::labs(title = gg_Code_title, caption = gg_Code_cap)
287
288 gg_Code_map

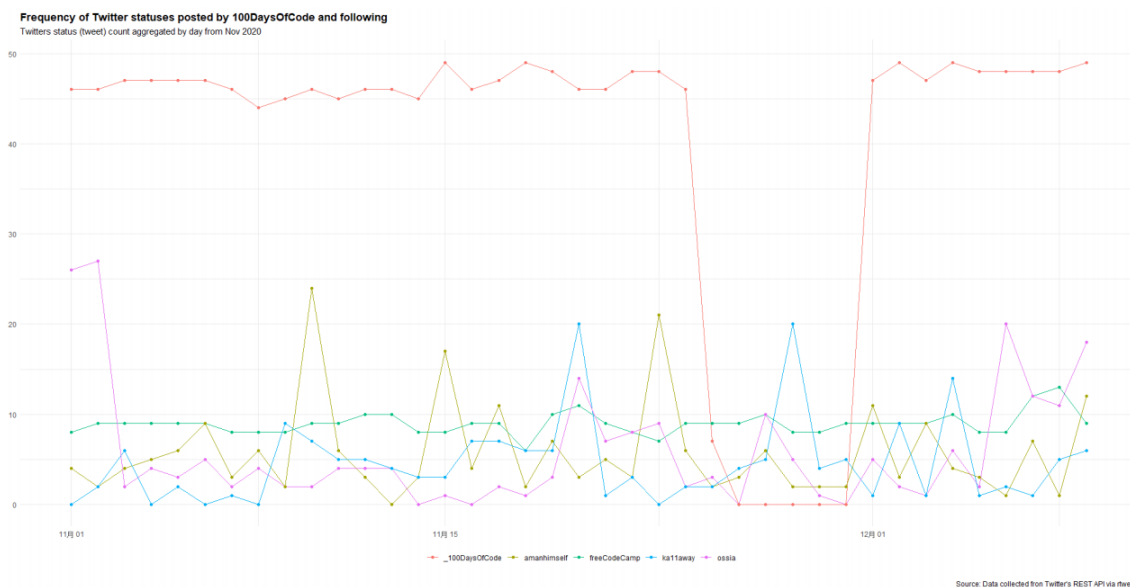
```

3. Results:

We use time plot to show the frequency of the hashtag 100DaysOfCode Twitter statuses from past 4 days. It can be seen from the figure that the number of daily releases in the past 4 days has exceeded 3000, so this hashtag has a wealth of data for us to analyze the content of the code.

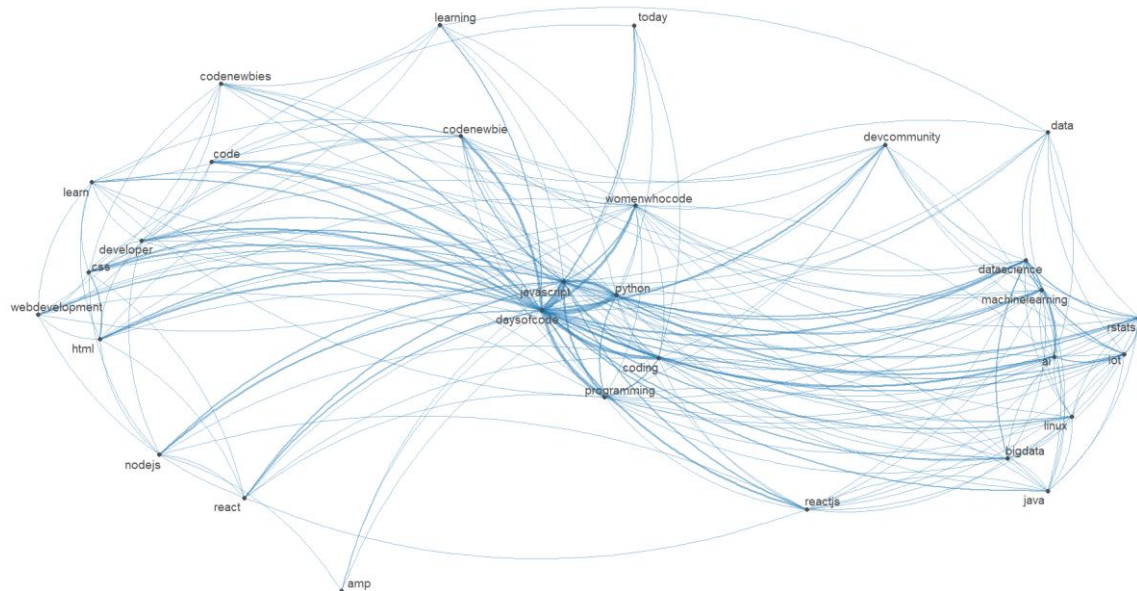


We do a tweet analysis on the official account “_100DaysOfCode” and the 4 people it followed. It can be found that “amahimself” and “_100DaysOfCode” are most closely related. When the number of “_100DaysOfCode” tweets increases, “amahimself” also increases most of the time.

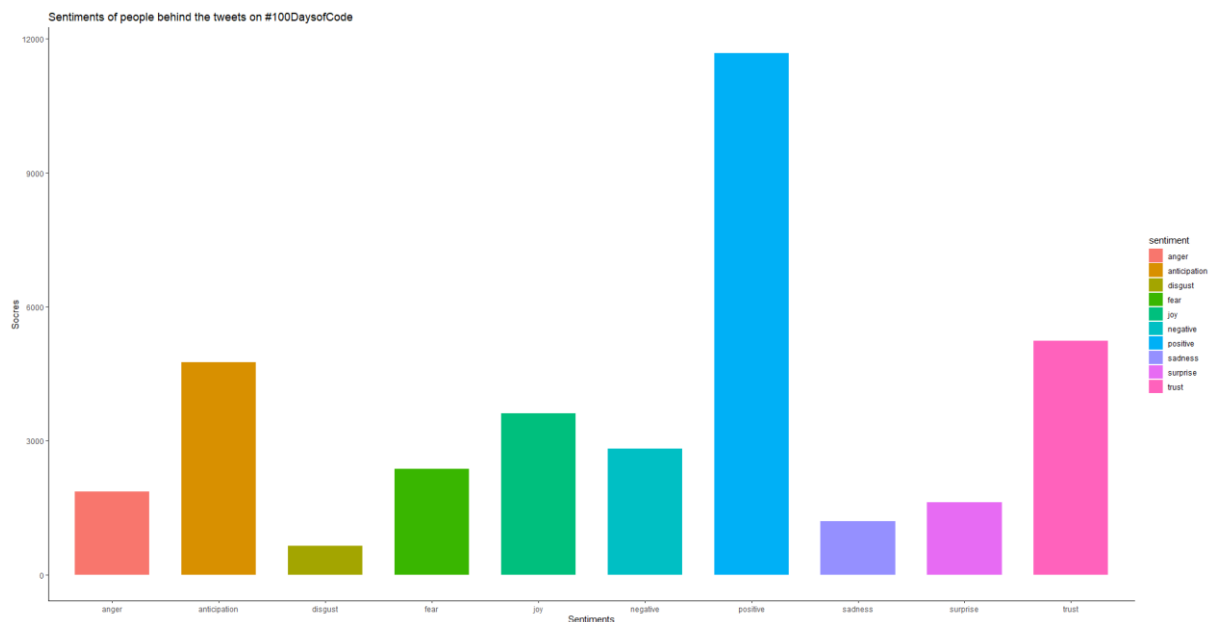


We use Twitter content to make word clouds. In addition to the core word *daysofcode*, *JavaScript* and *Python* appear the most frequently, indicating that many people are interested in these two languages. And we can find that the content of this hashtag is related to data

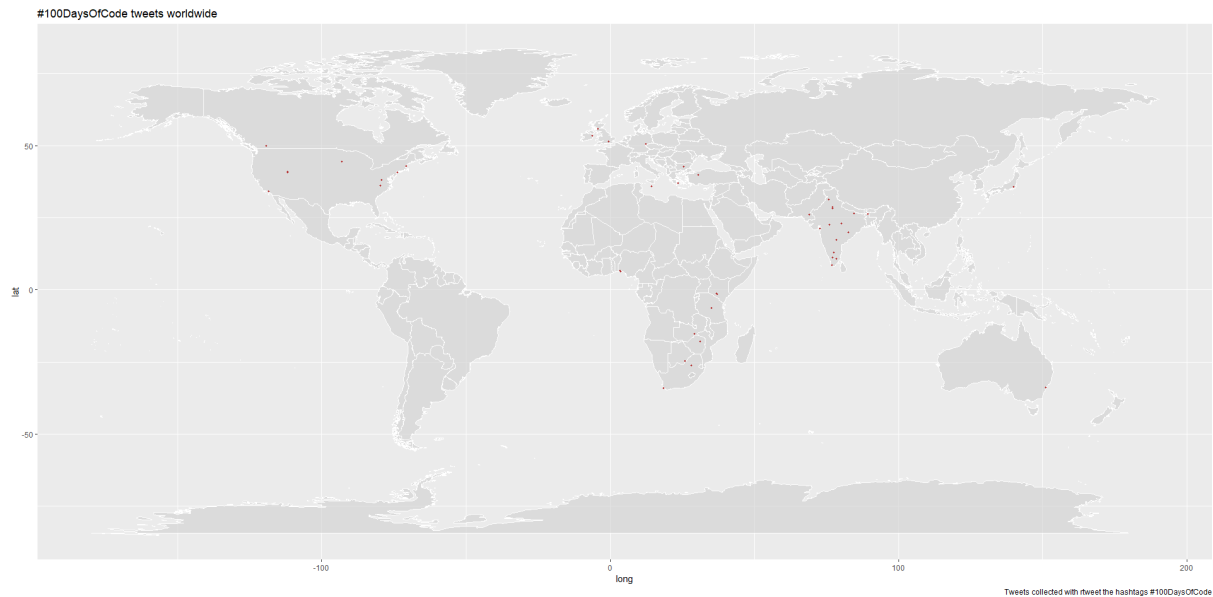
It can be seen from the node graph that the words with more occurrences are basically important nodes, and they are closely related to each other. The core word *daysofcode* has a deep connection with each node. In addition, for example, people who mention *JavaScript* usually mention *programming*, and people who mention *Python* often mention *devcommunity*. Explain that they all belong to the same "community".



We did a sentiment analysis to show everyone's tweet mood. Most people show positive attitudes on this topic, and many people also show trust that they can learn code well and look forward to learning and communicating.



We tried to explore the geographic location of the people who tweeted that with #100DaysOfCode, however, there are few tweets have latitude and longitude information, only 45 of our 9923 tweets have it. According to this map, we speculated that more people in India participated in this code event.



4. Conclusion:

From time plot we can see that there are many users participating in #100dayofdcode every day. By visualizing #100DaysOfCode's Twitter content, we can see that most of the people following this topic are interested in the computer field. Most people are interested in JavaScript and Python, and many of them are working on development, projects, or machine learning. So maybe we can engage in industries such as web production, data analysis, and artificial intelligence product development. We are also seeing a lot of new entrants in programming learning, which is consistent with the rapid development of data science. Interestingly, we found that a lot of female coder also actively tweets on this topic, indicating that the gender gap in the industry may be gradually decreasing.

Most of the people involved in hashtag maintain an optimistic and positive attitude towards code or communication. People who are enjoying fun and looking forward to learning than those people in a state of confusion or sadness.

Limitations and improvements:

According to the current word cloud, we still cannot judge whether users belong to code learners or workers, nor can we see their specific employment industries.

We need to do single-person Twitter content analysis for more users, so that we can analyze the individual in a comprehensive way, then we can know their learning purpose or employment direction.