

R Data Mining

Week 11 Lab in R:

Text Mining

Dr. Liang (Leon) Chen



Instructions

- R script for Text Mining is provided in the next slide, please follow it and complete the lab in R.
- You do not need to type notes (starting at #), but it's a good manner to have them in you script.
- In order to see codes and notes clearly, I show the script in RStudio.
- Advanced materials are available to you at the end.

```

1 #define and choose the dataset: this is same with Step 2 in RM import configuration wizard
2 Lab11Data<-read.csv(file.choose(),header=FALSE, sep="\t", quote="", stringsAsFactors=FALSE)
3 #check the dimension of your data
4 dim(Lab11Data)
5 #assign column names for your data
6 colnames(Lab11Data) <- c("type", "text")
7 #select your first column as your factor
8 Lab11Data$type <- factor(Lab11Data$type)
9 #check the structure of your data
10 str(Lab11Data)
11 #install and then use the package tm for text mining
12 install.packages("tm")
13 library(tm)
14 #Load the data as a corpus and rename it
15 Lab11Doc<-Corpus(VectorSource(Lab11Data$text))
16 #Inspect the content of the document using inspect(Lab11Doc); too long, so not included here
17 #Data cleaning,remove stopwords,numbers, etc.
18 Lab11Doc1<-tm_map(Lab11Doc,removePunctuation)
19 Lab11Doc1<-tm_map(Lab11Doc1,tolower)
20 Lab11Doc1<-tm_map(Lab11Doc1,removewords, stopwords(kind="en"))
21 Lab11Doc1<-tm_map(Lab11Doc1,stripwhitespace)
22 Lab11Doc1<-tm_map(Lab11Doc1,removeNumbers)
23 # Build a term-document matrix
24 Lab11Doc2<-TermDocumentMatrix(Lab11Doc1)
25 #find words that occur at least 100 times using findFreqTerms()
26 findFreqTerms(Lab11Doc2, lowfreq = 100)
27 #Analyze the association between frequent terms (i.e. terms which correlate) using findAssocs() function
28 findAssocs(Lab11Doc2, terms = "call", corlimit = 0.1)
29 #generate a matrix
30 Lab11Matrix<-as.matrix(Lab11Doc2)
31 #generate a frequency matrix based on the number frequency of each word
32 Lab11_freqs = sort(rowSums(Lab11Matrix), decreasing=TRUE)
33 Lab11DF = data.frame(word=names(Lab11_freqs), freq=Lab11_freqs)
34 #view the first 20 most frequent terms
35 head(Lab11DF,20)
36 #visualize the first 20 most frequent terms
37 library(ggplot2)
38 Lab11DF_top<-Lab11DF[1:20,]
39 ggplot(Lab11DF_top, aes(x = word, y = freq))+geom_bar(stat = "identity") + coord_flip() + labs(title = "Most Frequent words")
40 #install and use wordcloud library for visualizaiton: bag of words
41 install.packages("wordcloud")
42 library(wordcloud)
43 set.seed(123)
44 #generate bag of words to show the top 200 words.
45 wordcloud(word=Lab11DF$word, freq=Lab11DF$freq, max.words=200, random.order=FALSE, colors=brewer.pal(8, "Dark2"))

```

You can use stopwords(kind = "en")
to see all the stopwords in R

Deliverable 1: Take a screenshot of the associations with term "call".

Deliverable 2: Take a screenshot of the top 20 most frequent words.
Are the top five words the same with those in RapidMiner (sort the wordlist by total occurrences in Deliverable 4 of RapidMiner lab)

Deliverable 3: Take a screenshot of the bag of words you generate.

Advanced Materials (Optional)

- You can follow [this site](#) to conduct Naive Bayes Classification for our dataset
- Please find the instruction for ngram package [here](#).
- [**Text Mining with R A Tidy Approach**](#)
- [**Text mining and word cloud fundamentals in R**](#)
- [**Simple Wordcloud**](#)
- [**Text Mining and N-Grams**](#)