



Homework 5

ZHENYU XU

15620161152286

Q1、 2

```
rm(list = ls())
library(RCurl)
library(XML)
library(bitops)
library(stringr)

url=paste(c("http://publicliterature.org/pdf/2ws1610.pdf","http://publicliterature.org/pdf/2ws2410.pdf","http://publicliterature.org/pdf/2ws3310.pdf") )

abs=lapply(url, FUN = function(x) htmlParse(x, encoding = "Latin-1"))

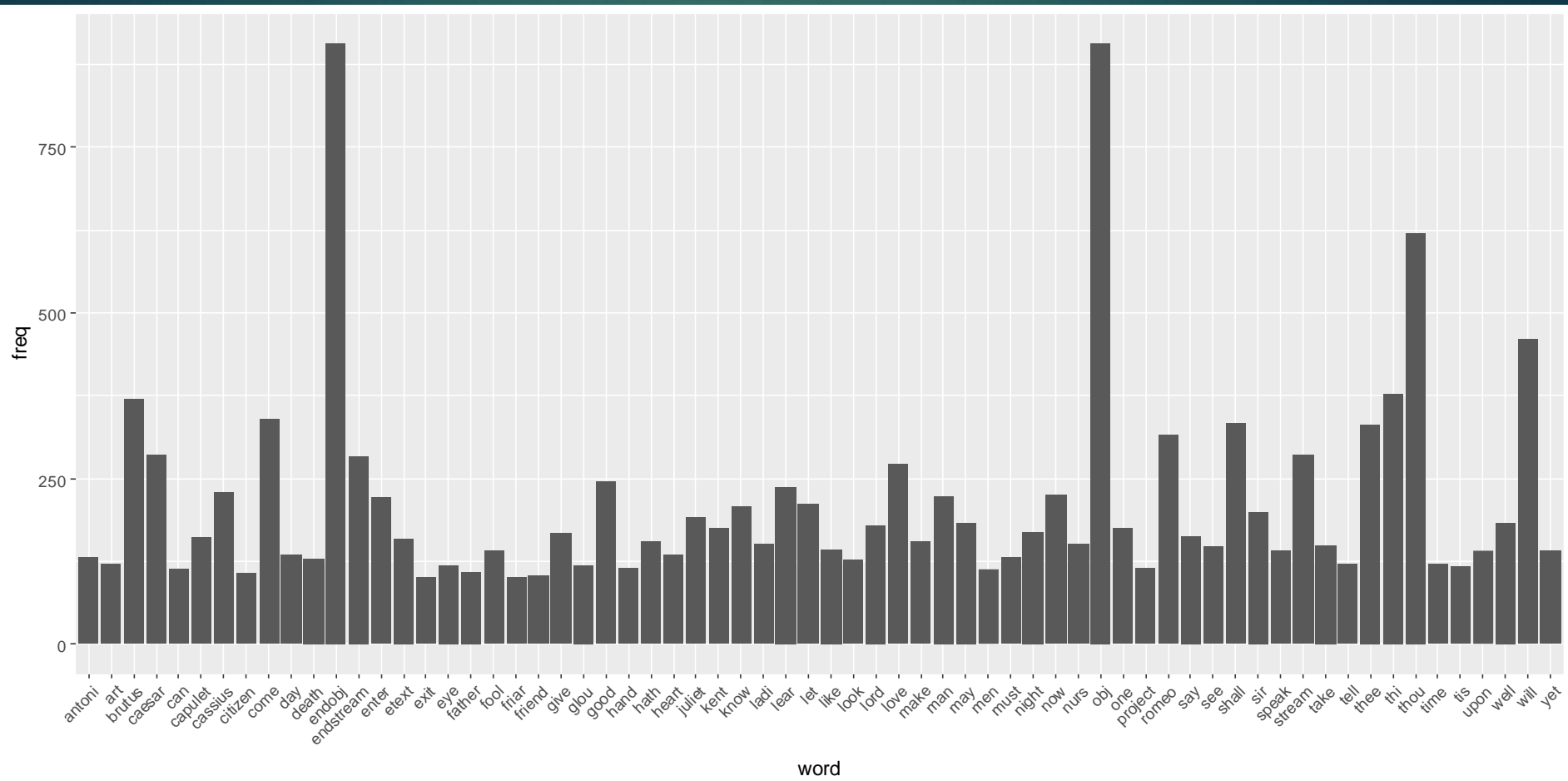
clean_txt = function(x) {
  cleantxt = xpathApply(x, "//body//text()
                        [not(ancestor :: script)][ not(ancestor :: style)]
                        [not(ancestor :: noscript)] " ,xmlValue)

  cleantxt = paste(cleantxt, collapse="\n")
  cleantxt = str_replace_all(cleantxt, "\n", " ")
  cleantxt = str_replace_all(cleantxt, "\r", "")
  cleantxt = str_replace_all(cleantxt, "\t", "")
  cleantxt = str_replace_all(cleantxt, "<br>", "")
  return(cleantxt)
}

cleantxt = lapply(abs,clean_txt)
vec_abs = unlist(cleantxt)
vec_abs
```

```
library(tm)
library(SnowballC)
abs = Corpus(VectorSource(vec_abs))
abs_dtm = DocumentTermMatrix(abs, control = list(
  stemming = TRUE, stopwords = TRUE, minWordLength = 3,
  removeNumbers = TRUE, removePunctuation = TRUE))
dim(abs_dtm)
inspect(abs_dtm)
#Find the words that occur more than 5 times
findFreqTerms(abs_dtm, 5)
#Remove sparse terms
removeSparseTerms(abs_dtm, 0.5)
inspect(removeSparseTerms(abs_dtm, 0.5))
library(ggplot2)
library(wordcloud)
freq = colSums(as.matrix(abs_dtm))
wf = data.frame(word=names(freq), freq=freq)
plot = ggplot(subset(wf, freq>100), aes(word, freq))
plot = plot + geom_bar(stat="identity")
plot = plot + theme(axis.text.x=element_text(angle=45, hjust=1))
plot
freq = colSums(as.matrix(abs_dtm))
dark2 = brewer.pal(8, "Dark2")
wordcloud(names(freq), freq, max.words=200, rot.per=0.1, colors=dark2)
dev.off()
```

shall
one fool
made
heaven brother
upon madam
romeo ladi
casca
endstream
part sweet
poor light life need
bring death
much take cassius good
mean better twocall find arm thing
montagu w hose know thi night may
keep everi love enter
show boy sir dear daughter face part doth king lie fearsinc
tybalt friar stand hath eye heat speak second gon way sword tell hast
henc like letter make look hand edg heart let brutus
back best w rong fire name must master mari dead
banish fellow put servant men exit god scenereceiv
see lord use yet glou son follow get alb gutenber till live
even run caesar stream gone done well thee
now say give endobj
thou nurs farewel thing
juliet place man
citizen bear
head antoni
tis
ob
ob
ob



```
hist(freq, col = "grey", breaks = 20, ylim = c(0, 5000), xlab = "freq of words")
```

