


```
ggplot  
)) +  
geom
```

```
library(wordcloud2)
library(NLP)

## ## Attaching packages: [NLP]
```

''' Neuchâtel package - NEI

```
## The following object is masked from 'package:ggplot2':  
##
```

```
##      annotate
```

```
library(tm)#used to do the text mining and text clearing  
library(readr)#used to read csv.file  
library(dplyr)#used to do some piping
```

```
## Warning: package 'dplyr' was built under R version 4.0.5

## 
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
## 
##     filter, lag

## The following objects are masked from 'package:base':
## 
##     intersect, setdiff, setequal, union

articles<-read.csv("articles.csv")#The data set is downloaded from Kaggle, and the link is provided in the reference list.



## Create a corpus from actual text



```
articles.corpus=Corpus(VectorSource(articles$title))
removeHTML=function(text){
 text=gsub(pattern='<.+\\>', '',text)
 text=gsub(pattern='</.+>', '',text)
 return(text)
}
```



## Clean corpus with tm functions



```
articles.corpus=articles.corpus%>%
 tm_map(content_transformer(removeHTML))%>%
 tm_map(removeNumbers)%>%
 tm_map(removePunctuation)%>%
 tm_map(stripWhitespace)%>%
 tm_map(content_transformer(tolower))%>%
 tm_map(removeWords,stopwords("english"))%>%
 tm_map(removeWords,stopwords("SMART"))
```


```

```
## Warning in tm_map.SimpleCorpus(., content  
## transformation drops documents
```

```
## Warning in tm_map.SimpleCorpus(., removeNumbers):  
  
## Warning in tm_map.SimpleCorpus(., removePunctuation)  
## documents  
  
## Warning in tm_map.SimpleCorpus(., stripWhitespace)  
## documents
```

```
## Warning in tm_map.SimpleCorpus(., content  
## drops documents
```

```
## Warning in tm_map.SimpleCorpus(., removeWords, stopwords("english")):  
## transformation drops documents  
  
## Warning in tm_map.SimpleCorpus(., removeWords, stopwords("SMART")):  
## transformation drops documents
```

Creat term document matrix

```
tdm=TermDocumentMatrix(articles.corpus)%>%#each row represent a word, an  
the cell correspond how many times the word appears in the document  
as.matrix()%>%convert it into a R matrix we can work with  
words=sort(rowSums(tdm),decreasing = TRUE)  
df=data.frame(word=names(words),freq=words)
```

Minor adjustments to data frame

```
df=df%>%  
  filter(nchar(as.character(word))>2,  
         word!="don'")
```

Create word cloud

 USP

1



4. Reference

1. <https://cran.r-project.org/web/packages/ggwordcloud/ggwordcloud.pdf>
 2. <https://cran.r-project.org/web/packages/wordcloud2/wordcloud2.pdf>
 3. <https://www.kaggle.com/datasets/hsankesara/medium-articles>