# Cleantweets

### Noblezada, Camasa, Cabia

### 2024-12-13

```r
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ---------------------- tidyverse 2.0.0 --
## v dplyr     1.1.4     v readr     2.1.5
## v forcats   1.0.0     v stringr   1.5.1
## v ggplot2   3.5.1     v tibble    3.2.1
## v lubridate 1.9.3     v tidyr     1.3.1
## v purrr     1.0.2
## -- Conflicts ---------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```r
library(tidytext)
library(dplyr)
library(stringr)
library(ggplot2)
library(sentimentr)
```

```r
tweetsDF <- read_csv("/cloud/project/ProjectDS/tweetsDF.csv")
```

```
## New names:
## Rows: 58086 Columns: 7
## -- Column specification
## ---------------------------------------------------- Delimiter: "," chr
## (4): screenName, text, statusSource, tweetSource dbl (1): ...1 dttm (2):
## created, Created_At_Round
## i Use `spec()` to retrieve the full column specification for this data. i
## Specify the column types or set `show_col_types = FALSE` to quiet this message.
## * `` -> `...1`
```

```r
tweetsDF$text <- iconv(tweetsDF$text, from = "UTF-8", to = "ASCII//TRANSLIT", sub = "")
tweetsDF$text <- tolower(tweetsDF$text)
tweetsDF$text <- gsub("https\\S+", "", tweetsDF$text)
tweetsDF$text <- gsub("#", "", gsub("\n", " ", tweetsDF$text))
tweetsDF$text <- gsub("([@?]\\S+)", "", tweetsDF$text)
tweetsDF$text <- gsub("\\?", "", tweetsDF$text)
tweetsDF$text <- gsub("\\b\\d{2}\\.\\d{2}\\.\\d{4}\\b", "", tweetsDF$text)
tweetsDF$text <- gsub("<a href=httptwitter.comdownloadiphone rel=nofollow>twitter for iphone<a>", "", tweetsDF$text)
tweetsDF$text <- gsub("<a href=([^>]*?) rel=nofollow>([^<]*?)<a>", "", tweetsDF$text)
tweetsDF$text <- gsub("<a href=httptwitter.comdownloadandroid rel=nofollow>twitter for android<a>", "", tweetsDF$text)
tweetsDF$text <- gsub("<a href= rel=nofollow>twitter web app<a>", "", tweetsDF$text)
tweetsDF$text <- gsub("30102022", "", tweetsDF$text)
tweetsDF$text <- gsub("\\s+", " ", tweetsDF$text)
```

```r
create_chunks <- function(df, start_row, end_row) {
return(df[start_row:end_row, ])
}

start_row <- 1
end_row <- 1000
chunk_data <- tweetsDF[start_row:end_row, ]



print(tweetsDF)
```

```
## # A tibble: 58,086 x 7
##      ...1 screenName    text  created             statusSource Created_At_Round
##     <dbl> <chr>         <chr> <dttm>              <chr>        <dttm>
## 1       1 whourj31      "a s~ 2022-10-30 23:59:43 "<a href=\"~ 2022-10-31 00:00:00
## 2       2 nnainot       "nah~ 2022-10-30 23:59:32 "<a href=\"~ 2022-10-31 00:00:00
## 3       3 febry_sri_M   " pr~ 2022-10-30 23:59:31 "<a href=\"~ 2022-10-31 00:00:00
## 4       4 telehuntwat~  "tra~ 2022-10-30 23:59:28 "<a href=\"~ 2022-10-31 00:00:00
## 5       5 Typing0824    "the~ 2022-10-30 23:59:20 "<a href=\"~ 2022-10-31 00:00:00
## 6       6 niccijsmith   "wha~ 2022-10-30 23:59:04 "<a href=\"~ 2022-10-31 00:00:00
## 7       7 502SPIDEY     "can~ 2022-10-30 23:58:56 "<a href=\"~ 2022-10-31 00:00:00
## 8       8 maeannesala~  "pra~ 2022-10-30 23:58:45 "<a href=\"~ 2022-10-31 00:00:00
## 9       9 bigvirtue1    "big~ 2022-10-30 23:58:37 "<a href=\"~ 2022-10-31 00:00:00
## 10     10 ashxxy        "the~ 2022-10-30 23:58:31 "<a href=\"~ 2022-10-31 00:00:00
## # i 58,076 more rows
## # i 1 more variable: tweetSource <chr>
```