# Text analysis in R

SOLVEIG BJØRKHOLT

**Statistisk sentralbyrå**
Statistics Norway

# Agenda

Who am I?

Text gathering

Text prepping

Text usage

Example from using unstructured text

Statistisk sentralbyrå
Statistics Norway

# Who am I?

- Solveig Bjørkholt

- Certified RStudio Trainer

- Working at Statistics Norway

https://github.com/Zunny369

**Statistisk sentralbyrå**
Statistics Norway

# Fetch it

- You might have a .txt, .pdf, .docx, .json, .csv or any sort of file.

  ◦ readtext

```
list_of_files <- list.files("./path/to/dir", full.names = TRUE)
text <- readtext(list_of_files, encoding = "UTF-8")
```

**Statistisk sentralbyrå**
Statistics Norway

- For this example, I will use text from the Lucerne Group event webpage

```r
library(rvest)
```

```r
webpage <- read_html("https://www.meetup.com/Lucerne-R-User-Group/events/276438125/") %>%
    html_nodes("div.flex-item.flex-item--2.eventContent > div > div > section:nth-child(1) > div") %>%
    html_text()
```

```
## [1] "DetailsLucerne RUG: Text Analytics + temporal-spatial visualizationThe Lucerne R User Group is pleased to announce its first event after the winter break.We meet online via Zoom (link below).We are looking forward to two new shiny presentations.
Solveig Bjørkholt (Statistics Norway) is a certified R instructor and will talk about text analysis. Next, our second speaker Fabian Mundt (University of Lucerne/PH-Karlsruhe) will talk about »Exploring temporal-spatial visualizations«.Time plan:- 18.00 -
18.05 - Virtual reception- 18.05 - 18.35 - Presentation by Fabian Mundt- 18.35 - 19.05 - Presentation by Solveig Bjørkholt- 19.05 - 19:40 - Discussion and Q&AZoom-Link https://unilu.zoom.us/j/91657568596?pwd=eVRGLzZPdnlVakdtMXBVaGZ0aTdydz09Meeting-ID:
[masked]Kenncode:[masked]Follow us on Twitter: @lucerne_rOur GitHub repo: https://github.com/Lucerne-R-User-Group"
```
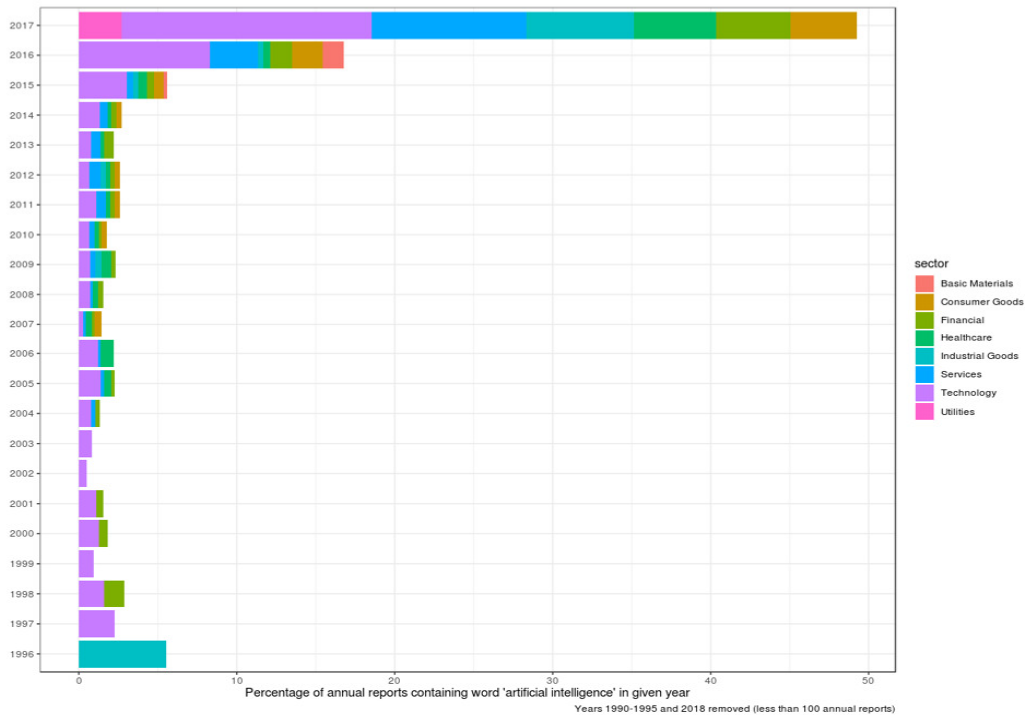
Statistisk sentralbyrå
Statistics Norway

# Clean it

- What contains important meaning in your text?

  ◦ Upper case? Lower case?

  ◦ Stopwords?

  ◦ Numbers?

  ◦ Punctation?

  ◦ Symbols?

  ◦ URLs?

  ◦ Email adresses?

  ◦ Stem? Lemmatize?

  ◦ Unigram, bigram, trigram?

Statistisk sentralbyrå
Statistics Norway

- Packages:
  - stringr
  - stringi

```r
library(stringr)

webpage <- str_to_lower(webpage)

webpage <- str_remove_all(webpage, "https://.*")

webpage <- str_remove_all(webpage, "[0-9]+")

webpage <- str_replace_all(webpage, "[[:punct:]]", " ")

webpage <- str_split_fixed(webpage, "link below", n = 2)

webpage <- str_squish(webpage)

webpage
```

```
## [1] "detailslucerne rug text analytics + temporal spatial visualizationthe lucerne r user group is pleased to announce its first event after the winter break we meet online via zoom"
## [2] "we are looking forward to two new shiny presentations solveig bjørkholt statistics norway is a certified r instructor and will talk about text analysis next our second speaker fabian mundt university of lucerne will talk about temporal spatial visualizations time plan virtual reception presentation by fabian mundt presentation by solveig bjørkholt discussion and q azoom link"
```

Statistisk sentralbyrå
Statistics Norway

# Vectorize it

```r
library(quanteda)
```

```r
dfm <- dfm(webpage)   # Make the string of text into a document feature matrix

dfm
```

- Bag of words

```r
dfm %>% convert(., to = "matrix") %>% knitr::kable()
```

| | detailslucerne | rug | text | analytics | + | temporal | spatial | visualizationthe | lucerne | r | user | group | is | pleased | to | announce | its | first | event | after | the |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| text1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| text2 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

- TF-IDF

```r
dfm %>% dfm_tfidf() %>% convert(., to = "matrix") %>% knitr::kable()
```

| | detailslucerne | rug | text | analytics | + | temporal | spatial | visualizationthe | lucerne | r | user | group | is | pleased | to | announce | its |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| text1 | 0.30103 0.30103 | 0 | 0.30103 | 0.30103 | 0 | 0 | 0.30103 | | 0 | 0 | 0.30103 | 0.30103 | 0 | 0.30103 | 0 | 0.30103 | 0.30103 |
| text2 | 0.00000 0.00000 | 0 | 0.00000 | 0.00000 | 0 | 0 | 0.00000 | | 0 | 0 | 0.00000 | 0.00000 | 0 | 0.00000 | 0 | 0.00000 | 0.00000 |

# What can you do with this?

## Descriptives

# Classification

## Unsupervised



## Supervised

# Example:
# Unstructured data as a source of information

# A world of data

# But what kind of data?

TIFF < TXT < PDF < XML

# Archives with scanned pictures

# From picture to text



=

# From text to table
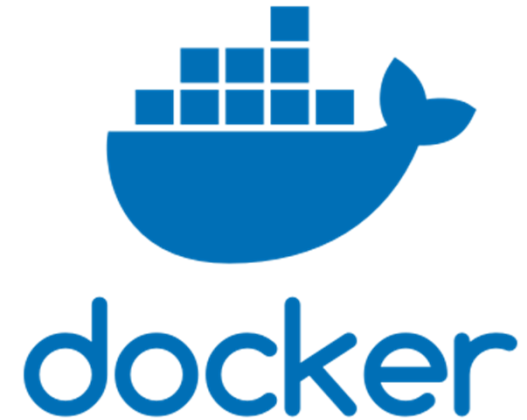
- *regex*

- dplyr

- stringr

- stringi

- lubridate

# From table to database

- dbplyr

- RSQLite

- elastic

- *Elasticsearch*

# From database to app

- Shiny

- data.table

- *Docker*





Statistisk sentralbyrå
Statistics Norway

**Årsregnskæppen**    Søk    Full arsrapport    Noter    Regnskap

## Velg foretak og år

```
URL til app:
sl-inno-p1:4770/arsrapporter/
```

🔗 Bookmark...

Skriv inn orgnr eller navn:

| 810034882: Sandnes El Forretning As | ▾ |

Velg et ar:

| 2019 | ▾ |

```
Orgnr: 810034882
Navn: Sandnes El Forretning As
År: 2019
```

**Statistisk sentralbyrå**
Statistics Norway

# Thanks!

Statistisk sentralbyrå
Statistics Norway