



R-LADIES JOZI NOV MEETUP

## ABOUT R-LADIES

- + Worldwide organisation whose mission is to promote gender diversity in the R community
- + R-Ladies' primary focus is on supporting minority gender R enthusiasts to achieve their programming potential, by building a collaborative global network
- + Our aim is to create a warm and welcoming place for R users to exchange ideas, learn and collaborate

## CONNECT WITH US



[https://www.meetup.com/  
rladies-johannesburg/](https://www.meetup.com/rladies-johannesburg/)



@RLadiesJozi  
@RLadiesGlobal  
@WeAreRLadies



[https://github.com/orgs/rladies/  
/teams/jozi/repositories](https://github.com/orgs/rladies/teams/jozi/repositories)



[https://rladies-  
community-  
slack.herokuapp.com/](https://rladies-community-slack.herokuapp.com/)

# DECEMBER MEETUP



## Year-end event:

*The Computers* documentary  
11 December  
Braamfontein

More info about the project:  
<http://eniacprogrammers.org/>

# TOPICS FOR 2019

Data science and ethics

Data visualisation



Please share ideas!

# THANK YOU



CoE for Mathematical and  
Statistical Sciences



International Relations



*Opinions expressed and conclusions arrived at are those of the author and are not necessarily to be attributed to the CoE-MaSS.*



# INTRODUCTION TO TEXT MINING IN R

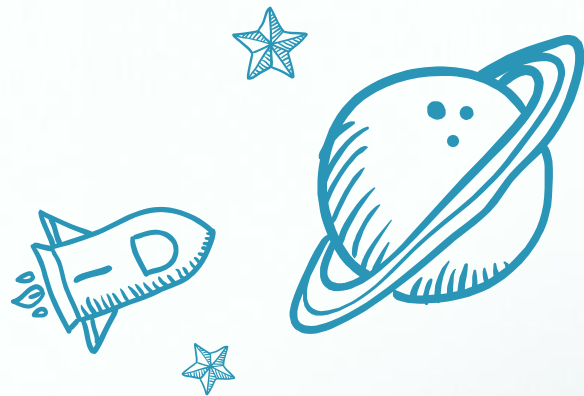
RETHA LANGA



# CONTENTS

- + 3 main ideas for tonight
- + Text mining
- + *Dubula ibhunu* trial case study
- + A practical example
- + Resources

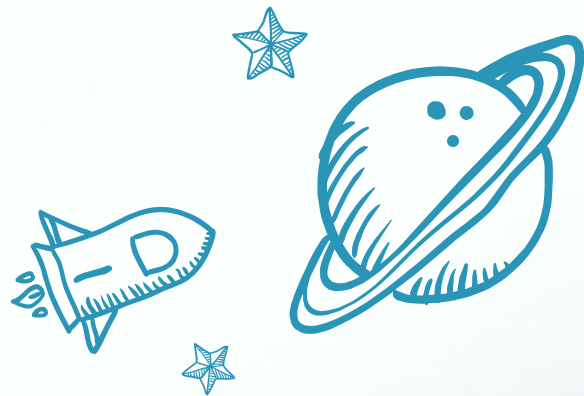




# KNOW YOUR DATA

Knowing your data and having solid industry/domain knowledge will lead to better questions and richer findings

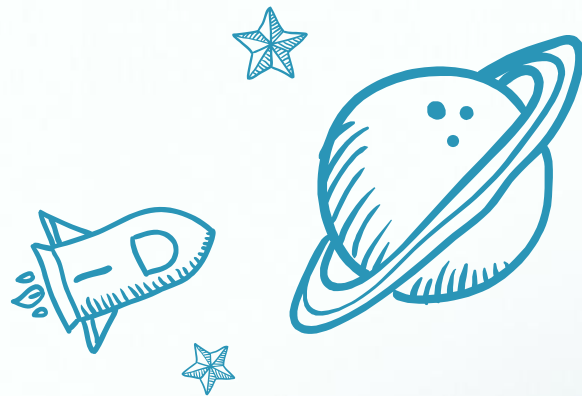




# THINK ABOUT YOUR ?

*"The art and science of asking questions is the source of all knowledge" – Thomas Berger*





# BROADEN YOUR HORIZONS

Applying ways of doing things from different fields can  
unlock some seriously cool – and ethical – stuff.





Process of analyzing collections of textual materials in order to capture key concepts and themes and uncover hidden relationships and trends

Source: IBM

## RESEARCH ARTICLE

### Quantitative Analysis of Culture Using Millions of Digitized Books

Jean-Baptiste Michel,<sup>1,2,3,4,5,†</sup> Yuan Kui Shen,<sup>2,6,7</sup> Aviva Presser Aiden,<sup>2,8,9</sup> Adrian Veres,<sup>2,6,9</sup> Matthew K. Gray,<sup>10</sup> The Google Books Team,<sup>10</sup> Joseph P. Pickett,<sup>11</sup> Dale Holberg,<sup>12</sup> Dan Clancy,<sup>10</sup> Peter Norvig,<sup>10</sup> Jon Orwant,<sup>10</sup> Steven Pinker,<sup>7</sup> Martin A. Nowak,<sup>1,3,14</sup> Erez Lieberman Aiden<sup>1,2,3,4,14,15,16,17,†</sup>

### Lexical shifts, substantive changes, and continuity in State of the Union discourse, 1790–2014

Alix Rule<sup>a</sup>, Jean-Philippe Cointet<sup>b</sup>, and Peter S. Bearman<sup>a,1</sup>



Article

### Using Sentiment Analysis to Monitor Electoral Campaigns: Method Matters—Evidence From the United States and Italy

Andrea Ceron<sup>1</sup>, Luigi Curini<sup>1</sup>, and Stefano M. Iacus<sup>1</sup>

Social Science Computer Review  
2015, Vol. 33(1) 3–20  
© The Author(s) 2014  
Reprints and permission:  
sagepub.com/journalsPermissions.nav  
DOI: 10.1177/0894439314521983  
scc.sagepub.com  
SAGE

### Content analysis of 150 years of British periodicals

Thomas Lansdall-Welfare<sup>a</sup>, Saatviga Sudhahar<sup>a</sup>, James Thompson<sup>b</sup>, Justin Lewis<sup>c</sup>, FindMyPast Newspaper Team<sup>d,1</sup>, and Nello Cristianini<sup>a,2</sup>



### Social media competitive analysis and text mining: A case study in the pizza industry

Wu He<sup>a,\*</sup>, Shenghua Zha<sup>b,1</sup>, Ling Li<sup>a,c,2</sup>

### Understanding Satisfied and Dissatisfied Hotel Customers: Text Mining of Online Hotel Reviews

Katerina Berezina, Anil Bilgihan, Cihan Cobanoglu & Fevzi Okumus



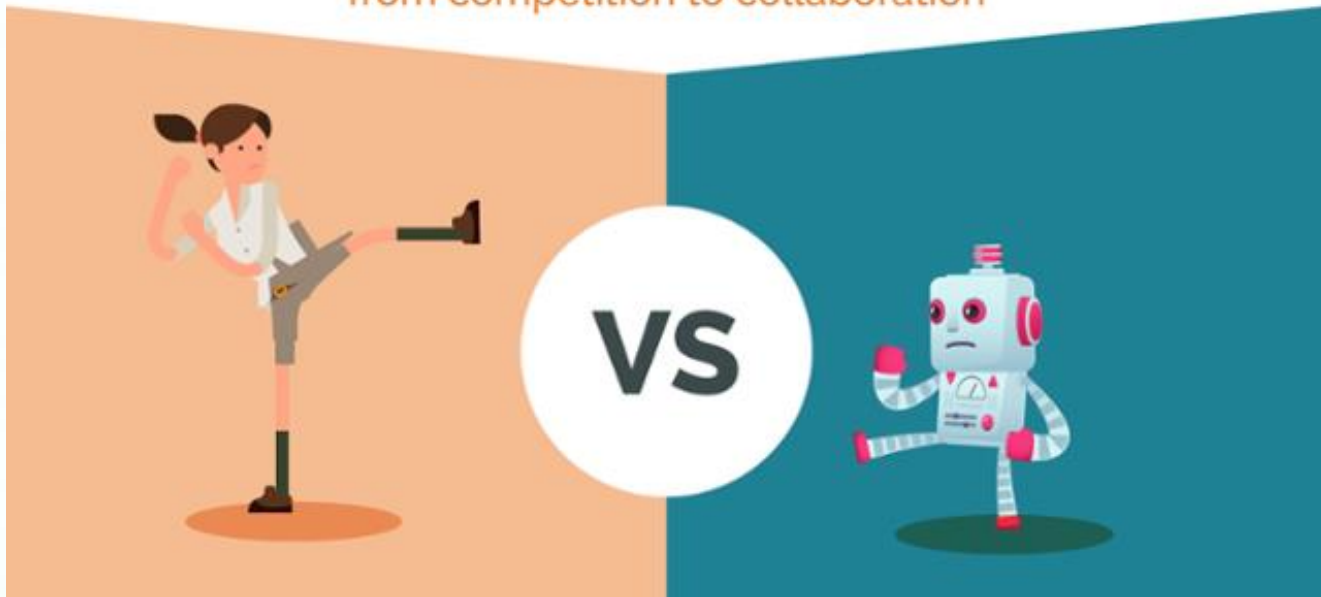


# *DUBULA IBHUNU TRIAL CASE STUDY*





# (Wo)man vs machine; from competition to collaboration





*HOW DID THE DIFFERENT ORGANISATIONS AND  
WITNESSES DEPLOY THE LAND ISSUE IN COURT?*



111 267

Words of testimony



3 organisations

ANC, AfriForum & TAU SA



11

Witness testimonies



Derek Hanekom



Collins Chabane



Dr Mongane Wally Serote



Gwede Mantashe



Julius Malema

Photos: Tourism.gov; SABC; City Press;  
Wikipedia

### 3 METHODOLOGICAL CHOICES

#### Bag of words

Words are treated as single tokens.

Only the frequencies of words per text are used and their positions are ignored.

#### Dictionary

Manual approach to compiling the list of keywords.

#### Sentiment analysis

Testimonies were characterised by a high level of emotion.

Witnesses often became emotional during the trial, especially when testifying about land reform.



### 3 THEMES





# *TURNING TO R*

# TEXT MINING STEPS

Importing text

Pre-  
processing

Creating a  
doc-term  
matrix  
(DTM)

Analysing  
your DTM

# IMPORTING TEXT

## STEP 1: Create a corpus

PC > Desktop > Textmining > Song >

| <input type="checkbox"/>            | Name        | Date modified        | Type          | Size  |
|-------------------------------------|-------------|----------------------|---------------|-------|
| <input checked="" type="checkbox"/> | song_corpus | 2018/11/24 12:28 ... | File folder   |       |
| <input type="checkbox"/>            | _Rhistory   | 2018/11/25 4:43 P... | RHISTORY File | 2 KB  |
| <input type="checkbox"/>            | Song        | 2018/11/25 4:23 P... | R Project     | 1 KB  |
| <input type="checkbox"/>            | Song_ladies | 2018/11/25 4:27 P... | R File        | 13 KB |

PC > Desktop > Textmining > Song > song\_corpus

| <input type="checkbox"/> | Name                     | Date modified        | Type          |
|--------------------------|--------------------------|----------------------|---------------|
| <input type="checkbox"/> | bezuidehout_court_2011   | 2018/09/07 10:03 ... | Text Document |
| <input type="checkbox"/> | chabane_court_2011       | 2018/09/07 10:00 ... | Text Document |
| <input type="checkbox"/> | combafriforum_court_2011 | 2018/09/07 10:00 ... | Text Document |
| <input type="checkbox"/> | comball_court_2011       | 2018/09/07 12:46 ... | Text Document |
| <input type="checkbox"/> | combanc_court_2011       | 2018/09/07 10:01 ... | Text Document |
| <input type="checkbox"/> | combtausa_court_2011     | 2018/09/07 10:03 ... | Text Document |
| <input type="checkbox"/> | goosen_court_2011        | 2018/09/07 10:03 ... | Text Document |
| <input type="checkbox"/> | gray_court_2011          | 2018/09/07 10:00 ... | Text Document |
| <input type="checkbox"/> | hanekom_court_2011       | 2018/09/07 10:02 ... | Text Document |
| <input type="checkbox"/> | kok_court_2011           | 2018/09/07 10:03 ... | Text Document |
| <input type="checkbox"/> | malema_court_2011        | 2018/09/07 10:02 ... | Text Document |
| <input type="checkbox"/> | mantashe_court_2011      | 2018/09/07 10:01 ... | Text Document |
| <input type="checkbox"/> | roets_court_2011         | 2018/09/07 10:00 ... | Text Document |
| <input type="checkbox"/> | serote_court_2011        | 2018/09/07 10:01 ... | Text Document |
| <input type="checkbox"/> | vanzyl_court_2011        | 2018/09/07 10:03 ... | Text Document |

A corpus in text mining is a database for holding and managing the text documents to be analysed – i.e. a text document collection.

# PACKAGES

```
library(readtext)↵  
library(stringr)↵  
library(quantda)↵  
library(tidytext)↵  
library(ggplot2)↵  
library(dplyr)↵  
library(tidyr)↵
```



# PRE-PROCESSING

Common steps include:

- + Using string operations to modify or remove parts of the text
- + Splitting the text into tokens (tokenization)
- + Normalization of the text (for example, converting to lower-case and stemming)
- + Removing stopwords, punctuation and numbers



## READING-IN YOUR DATA

```
21 ─  
22 court <- readtext("song_corpus/*_court_2011.txt")  
23 court  
24 ─  
25 random <- c("lordship", "lord", "sir", "indistinct", "inaudible")  
26 ─  
27 court$text <- tolower(court$text)  
28 random.regex <- paste(random, collapse = "|")  
29 random.regex  
30 ─  
31 court$text <- str_remove_all(court$text, pattern = random.regex)  
32 court  
33 ─  
34 ─
```

# PRE-PROCESSING AND CREATING A DTM

What is a Document–Term Matrix?

- + Document term matrix is one of the most common formats for representing a text corpus in a bag-of-words format.
- + A DTM is a matrix in which rows are documents, columns are terms, and cells indicate how often each term occurred in each document.

# PRE-PROCESSING AND CREATING A DTM

```
34 ─  
35 #Creating document-term matrix─  
36 ─  
37 court.corpus <- corpus(court)─  
38 court.dtm <- dfm(court.corpus,─  
39 ..... tolower = TRUE,─  
40 ..... stem = TRUE,─  
41 ..... remove_punct = TRUE,─  
42 ..... remove_numbers = TRUE,─  
43 ..... remove_hyphens = TRUE,─  
44 ..... remove_symbols = TRUE,─  
45 ..... remove = stopwords("english"))─  
46 ─
```

# WHAT DOES IT LOOK LIKE?

```
> court.dtm[, 1:10]
Document-feature matrix of: 15 documents, 10 features (27.3% sparse).
15 x 10 sparse Matrix of class "dfm"
      features
docs  read correct okay music general probabl biggest tool influenc human
bezuidenhout_court_2011.txt 13      4      9    22      5      4      1      8     21      2
chabane_court_2011.txt      2      3      6      3      3      0      0      0      0      0
combafriforum_court_2011.txt 18     24     18     52      8     11      1      5      6      1
comball_court_2011.txt     80    243     84     78     54    26      5     16     31     20
combanc_court_2011.txt     41    128     57      3     33    10      3      2      4     17
combtausa_court_2011.txt    21     91      9     23     13      5      1      9     21      2
goosen_court_2011.txt       3     33      0      1      2      0      0      0      0      0
gray_court_2011.txt         4      4      3     50      2      5      0      5      4      0
hanekom_court_2011.txt     18     60     27      0     12    10      2      2      2      5
kok_court_2011.txt          1     24      0      0      0      0      0      1      0      0
malema_court_2011.txt       8     13      4      0      7      0      0      0      2      8
mantashe_court_2011.txt     5     16     10      0      8      0      0      0      0      3
roets_court_2011.txt       14     20     15      2      6      6      1      0      2      1
serote_court_2011.txt       8     36     10      0      3      0      1      0      0      1
vanzyl_court_2011.txt       4     30      0      0      6      1      0      0      0      0
```

# ANALYSING YOUR DTM – A DICTIONARY APPROACH

```
51 ↵
52 #Creating a dictionary of keywords that belong in the three clusters (land, reconciliation & transformation)↵
53 ↵
54 ↵
55 my_Dict<-dictionary(list(land=c("land*"), farm=c("farm*"), farmer=c("farmer*"), ↵
56 .....transform=c("transform*"), freedom=c("freedom*"), ↵
57 .....revolution=c("revolu*"), ↵
58 .....cohesion=c("cohes*"), ↵
59 .....reconciliation=c("reconcil*"), redress=c("redress*"), ↵
60 .....mandela=c("mandela*"), dialogue=c("dialog*"), ↵
61 .....nationalisation=c("nationali*"), compensation=c("compens*"), ↵
62 .....rainbow=c("rainbow*"), agriculture=c("agricultur*"), ↵
63 .....charter=c("charter*"), ↵
64 .....constitution=c("constituti*"))↵
65 dict_court.dtm<-dfm_lookup(court.dtm, my_Dict, nomatch="unmatched")↵
66 dict_court.dtm↵
67 ↵
68 ↵
```



# THE RESULT

```
> dict_court.dtm
```

Document-feature matrix of: 15 documents, 18 features (44.4% sparse).

15 x 18 sparse Matrix of class "dfm"

|                              | features |      |        |           |         |            |          |                |         |         |
|------------------------------|----------|------|--------|-----------|---------|------------|----------|----------------|---------|---------|
| docs                         | land     | farm | farmer | transform | freedom | revolution | cohesion | reconciliation | redress | mandela |
| bezuidenhout_court_2011.txt  | 2        | 25   | 6      | 2         | 2       | 4          | 0        | 5              | 0       | 0       |
| chabane_court_2011.txt       | 1        | 0    | 0      | 0         | 1       | 0          | 0        | 0              | 0       | 0       |
| combafriforum_court_2011.txt | 11       | 45   | 18     | 0         | 2       | 0          | 3        | 9              | 0       | 5       |
| comball_court_2011.txt       | 82       | 219  | 126    | 15        | 63      | 40         | 3        | 23             | 13      | 23      |
| combanc_court_2011.txt       | 66       | 128  | 92     | 12        | 47      | 36         | 0        | 9              | 13      | 18      |
| combtausa_court_2011.txt     | 5        | 46   | 16     | 3         | 14      | 4          | 0        | 5              | 0       | 0       |
| goosen_court_2011.txt        | 0        | 2    | 2      | 0         | 9       | 0          | 0        | 0              | 0       | 0       |
| gray_court_2011.txt          | 11       | 7    | 1      | 0         | 0       | 0          | 0        | 7              | 0       | 5       |
| hanekom_court_2011.txt       | 2        | 70   | 50     | 1         | 11      | 1          | 0        | 3              | 0       | 3       |
| kok_court_2011.txt           | 0        | 1    | 1      | 1         | 0       | 0          | 0        | 0              | 0       | 0       |
| malema_court_2011.txt        | 51       | 44   | 29     | 8         | 18      | 28         | 0        | 0              | 0       | 11      |
| mantashe_court_2011.txt      | 7        | 6    | 5      | 0         | 5       | 2          | 0        | 2              | 0       | 1       |
| roets_court_2011.txt         | 0        | 38   | 17     | 0         | 2       | 0          | 3        | 2              | 0       | 0       |
| serote_court_2011.txt        | 5        | 8    | 8      | 3         | 12      | 5          | 0        | 4              | 13      | 3       |
| vanzyl_court_2011.txt        | 3        | 18   | 7      | 0         | 3       | 0          | 0        | 0              | 0       | 0       |



## RECONCILIATION

|                | Frequency  | %           |
|----------------|------------|-------------|
| Cohesion       | 3          | 0.5         |
| Constitution   | 45         | 6.9         |
| Dialogue       | 66         | 10.1        |
| Mandela        | 23         | 3.5         |
| Rainbow        | 1          | 0.2         |
| Reconciliation | 23         | 3.5         |
| Redress        | 13         | 2           |
| <b>TOTAL</b>   | <b>174</b> | <b>26.6</b> |

# TRANSFORMATION

## THE FREEDOM CHARTER

ADOPTED AT THE CONGRESS OF THE PEOPLE AT  
KLIPTOWN, JOHANNESBURG, ON JUNE 25 AND 26, 1955.

WE, the People of South Africa, declare for all our country and the world to know:

that South Africa belongs to all who live in it, black and white, and that no government can justly claim authority unless it is based on the will of all the people;  
that our people have been robbed of their birthright to land, liberty and peace by a form of government founded on injustice and inequality;  
that our country will never be prosperous or free until all our people live in brotherhood, enjoying equal rights and opportunities; (Sesuto: -) (-) (-)  
that only a democratic state, based on the will of all the people, can secure to all their birthright without distinction of colour, race, sex or belief;

And therefore we, the People of South Africa, black and white together — equals, countrymen and brothers — adopt this Freedom Charter. And we pledge ourselves to strive together sparing neither strength nor courage, until the democratic changes here set out have been won.

### THE PEOPLE SHALL GOVERN!

Every man and woman shall have the right to vote for and to stand as a candidate for all bodies which make laws; (-) (-) (-)  
All people shall be entitled to take part in the administration of the country; (-) (-) (-)  
The rights of the people shall be the same, regardless of race, colour or sex; (-) (-) (-)

All people shall have equal rights to trade where they choose, to manufacture and to enter all trades, crafts and professions.

### THE LAND SHALL BE SHARED AMONG THOSE WHO WORK IT!

Restriction of land ownership on a racial basis shall be ended, and all the land re-

|                 | Frequency | %    |
|-----------------|-----------|------|
| Charter         | 20        | 3.1  |
| Freedom         | 63        | 9.6  |
| Nationalisation | 6         | 0.9  |
| Revolution      | 40        | 6.1  |
| Transformation  | 15        | 2.3  |
| TOTAL           | 144       | 22.1 |

## LAND

|              | Frequency | %    |
|--------------|-----------|------|
| Agriculture  | 29        | 4.4  |
| Compensation | 5         | 0.8  |
| Farm         | 93        | 14.2 |
| Farmer       | 126       | 19.3 |
| Land         | 82        | 12.6 |
| TOTAL        | 335       | 51.3 |

## ANC, TAU SA AND AFRIFORUM

|           | Reconciliation | Transformation | Land |
|-----------|----------------|----------------|------|
| ANC       | 115            | 120            | 208  |
| AfriForum | 38             | 3              | 66   |
| TAU SA    | 21             | 21             | 61   |

Derek Hanekom



Collins Chabane



Dr Mongane Wally Serote



Gwede Mantashe



Julius Malema

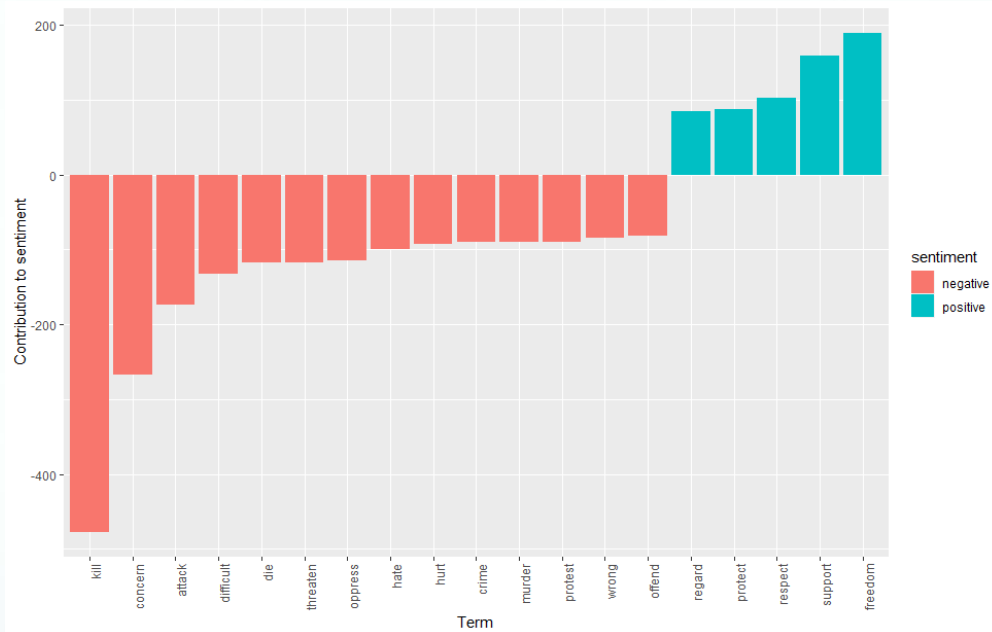
Photos: Tourism.gov, SABC, City Press;  
Wikipedia



## MALEMA AND THE REST OF THE ANC

|             | Reconciliation | Transformation | Land |
|-------------|----------------|----------------|------|
| Malema      | 29             | 76             | 112  |
| ANC veteran | 42             | 29             | 42   |

# SENTIMENT ANALYSIS





*LETS PRACTISE*



## 2018 STATE OF THE NATION AND OPPOSITION REPLY



# IMPORTING TEXT

## STEP 1: Create a corpus

View

PC > Desktop > Textmining > Rladies >

| <input type="checkbox"/>            | Name        | Date modified        | Type          | Size |
|-------------------------------------|-------------|----------------------|---------------|------|
| <input checked="" type="checkbox"/> | sona_corpus | 2018/11/25 12:38 ... | File folder   |      |
| <input type="checkbox"/>            | .Rhistory   | 2018/11/25 4:23 P... | RHISTORY File | 7 KB |
| <input type="checkbox"/>            | Example     | 2018/11/25 4:21 P... | R File        | 4 KB |
| <input type="checkbox"/>            | Rladies     | 2018/11/25 10:24 ... | R Project     | 1 KB |

ew

> Desktop > Textmining > Rladies > sona\_corpus

| Name                | Date modified        | Type          | Size  |
|---------------------|----------------------|---------------|-------|
| mairane_sona_2018   | 2018/11/25 10:37 ... | Text Document | 17 KB |
| ramaphosa_sona_2018 | 2018/11/25 12:38 ... | Text Document | 33 KB |



## RESOURCES

The quest for better questions:

<https://bit.ly/2OZUifa>

Ethical AI? It's All About Perspective:

<https://bit.ly/2AnvPuL>

A wolf in sheep's clothing: Disruption is overrated in terms of innovation:

<https://bit.ly/2BzVwdH>

Biography of an uncharted people:

<https://unchartedpeople.org/>

What is the role of the Arts and Humanities in the age of Data Science? Two short proposals:

<https://bit.ly/2ByRylj>

## RESOURCES

*Text mining with R: A tidy approach* by Julia Silge & David Robinson

<https://www.tidytextmining.com/>

Julia Silge:

<https://juliasilge.com/>

@juliasilge

Quanteda package materials:

<https://quanteda.io/>

Text Analysis in R article (hand-out)

Kaggle datasets:

<https://bit.ly/2THxFzq>

# RESOURCES

General resources:

[https://github.com/stepthom/text\\_mining\\_resources](https://github.com/stepthom/text_mining_resources)

Conversion Vignette:

<https://bit.ly/2P2UJ8i>

Plotting: <https://bit.ly/2PUm10k>



# THANK YOU!

## Any questions?

You can find me at:

- + @HeritageAffair
- + rethag@gmail.com