

LDA Topicmodel

Riccardo Ruta riccardo.ruta@studenti.unimi.it

06/05/2022

Contents

| | |
|--|----|
| PART I - CREATE THE DTM | 1 |
| PART II - FIND THE BEST NUMBER OF TOPICS K | 3 |
| PART III - ANALISYS OF THE TOPICS | 10 |
| PART IV - RE-ANALISYS OF THE TOPICS as GENRE | 12 |

PART I - CREATE THE DTM

1) Import and clean dataset

```
# import the data
tw <- read_csv("data/large_files/politicians_all_final_tweets.csv")

## Rows: 396611 Columns: 10
## -- Column specification -----
## Delimiter: ","
## chr (9): tw_screen_name, nome, tweet_testo, creato_il, url, party_id, genere...
## dbl (1): creato_il_code
##
## 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.
```

```
# Adjust datetime (Run code in this order!)
Sys.setlocale("LC_TIME", "C")
tw$date <- as.Date(strptime(tw$creato_il,"%a %b %d %H:%M:%S %z %Y", tz = "CET"))
tw$date <- na.replace(tw$date, as.Date(tw$creato_il))

# Create week variable
tw <- tw %>% mutate(week = cut.Date(date, breaks = "1 week", labels = FALSE))

# Create month variable
tw <- tw %>% mutate(month = cut.Date(date, breaks = "1 month", labels = FALSE))

# Remove missing from tweets column (using remove_na tidyverse)
tw <- tw %>% drop_na(tweet_testo)
```

```
# Remove space from genere variable
a <- unique(tw$genere)
tw$genere <- gsub(a[3], "male", tw$genere)

# Select variables for the analysis
dataset <- tw %>% select(nome, tweet_testo, genere, party_id, chamber, status, date, week, month )
```

2) Create the corpus

```
corp <- corpus(dataset, text = "tweet_testo")
```

3) Create the Dfm removing stopwords and trimming

```
my_word <- as.list(read_csv("data/it_stopwords_new_list.csv", show_col_types = FALSE)) #read.csv with u
my_list <- c(" ", "c'è", "+", " ", my_word$stopwords, stopwords('italian'))
rev_dfm <- dfm(corp, remove = c(stopwords("italian"), my_list),
               tolower = TRUE, stem = FALSE, remove_punct = TRUE, remove_numbers = TRUE)
```

```
## Warning: 'dfm.corpus()' is deprecated. Use 'tokens()' first.
```

```
## Warning: '...' should not be used for tokens() arguments; use 'tokens()' first.
```

```
## Warning: 'remove' is deprecated; use dfm_remove() instead
```

```
## Warning: 'stem' is deprecated; use dfm_wordstem() instead
```

```
rev_dfm <- dfm_trim(rev_dfm, min_termfreq = 0.95, termfreq_type = "quantile",
                   max_docfreq = 0.1, docfreq_type = "prop")
rev_dfm <- dfm_group(rev_dfm, groups = month)

kable(topfeatures(rev_dfm, 20), col.names = "TopFeatures")
```

| | TopFeatures |
|--------------------------|-------------|
| governo | 25991 |
| grazi | 20761 |
| lavoro | 18274 |
| paes | 16444 |
| anni | 16281 |
| <U+0001F1EE><U+0001F1F9> | 15196 |
| president | 14215 |
| grand | 13662 |
| italiani | 11993 |
| italia | 11957 |
| l'italia | 11728 |
| via | 11495 |
| <U+0001F534> | 10852 |
| politica | 9931 |
| cittadini | 9331 |
| bene | 9269 |
| forza | 8474 |
| ministro | 8413 |
| buon | 8307 |
| insiem | 8142 |

4) Convert the Document Feature Matrix (Dfm) in a Topic Model (Dtm)

```
dtm <- quanteda::convert(rev_dfm, to = "topicmodels")
```

PART II - FIND THE BEST NUMBER OF TOPICS K

1) First try K = 20:80

```
load("data/results1.Rda")
```

```
## Finding the best K
top1 <- c(20:80)

## let's create an empty data frame
results1 <- data.frame(first=vector(), second=vector(), third=vector())

system.time(
  for (i in top1)
  {
    set.seed(123)
    lda1 <- LDA(dtm, method= "Gibbs", k = (i), control=list(verbose=50L, iter=100))
    topic <- (i)
    coherence <- mean(topic_coherence(lda1, dtm))
    exclusivity <- mean(topic_exclusivity(lda1))
    results1 <- rbind(results1 , cbind(topic, coherence, exclusivity ))
  }
}
```

```
)
# save(results1,file="data/results1.Rda")
```

```
kable(head(results1))
```

| topic | coherence | exclusivity |
|-------|------------|-------------|
| 20 | 0.1674464 | 9.352493 |
| 21 | 0.8376643 | 9.386225 |
| 22 | 0.7783789 | 9.396400 |
| 23 | 0.3991629 | 9.420433 |
| 24 | 0.2202299 | 9.457425 |
| 25 | -3.0470840 | 9.470567 |

```
kable(tail(results1))
```

| | topic | coherence | exclusivity |
|----|-------|------------|-------------|
| 59 | 78 | -3.9815896 | 9.821740 |
| 60 | 79 | -2.6484726 | 9.807149 |
| 61 | 80 | -3.8557027 | 9.840021 |
| 62 | 20 | 0.1674464 | 9.352493 |
| 63 | 21 | 0.8376643 | 9.386225 |
| 64 | 22 | 0.7783789 | 9.396400 |

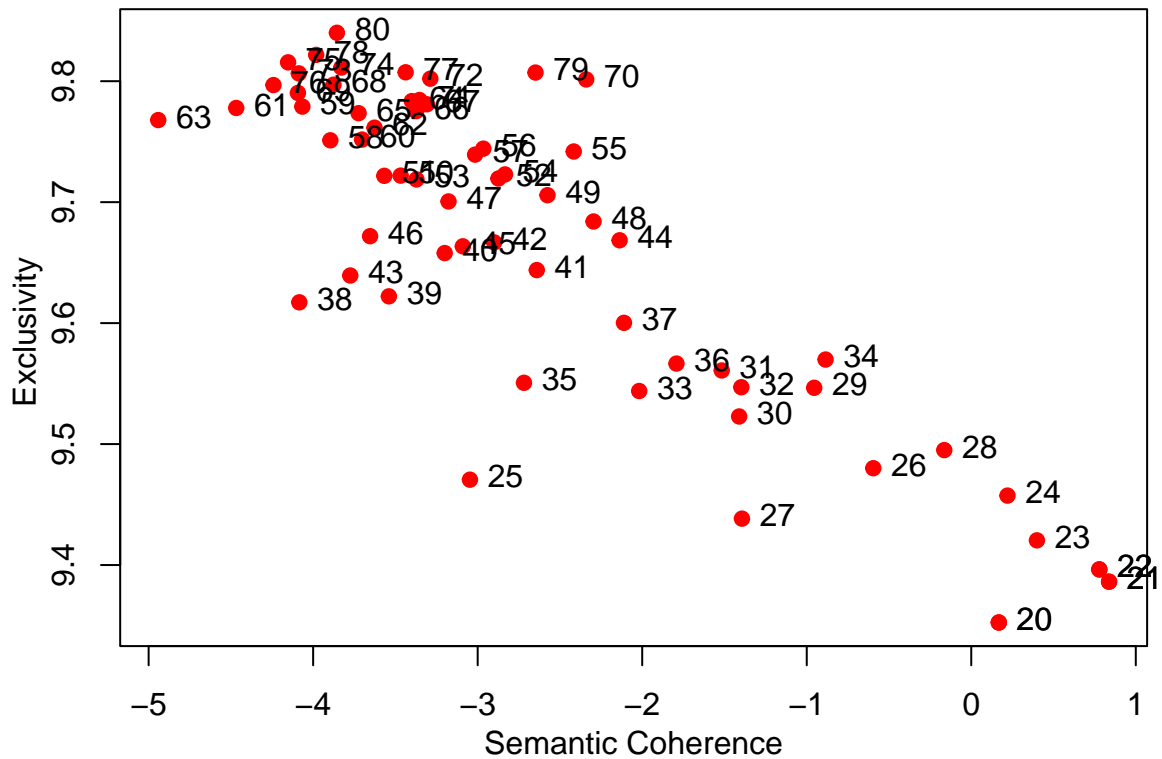
```
str(results1)
```

```
## 'data.frame':   64 obs. of  3 variables:
## $ topic      : num  20 21 22 23 24 25 26 27 28 29 ...
## $ coherence  : num  0.167 0.838 0.778 0.399 0.22 ...
## $ exclusivity: num  9.35 9.39 9.4 9.42 9.46 ...
```

```
plot1 <- as.ggplot(~plot(results1$coherence, results1$exclusivity,
                        main="Scatterplot K=20:80",xlab="Semantic Coherence",
                        ylab="Exclusivity ", pch=19,
                        col=ifelse(results1$coherence<=-155.8,"black","red")) +
  text(results1$coherence, results1$exclusivity,
        labels=results1$topic, cex= 1, pos=4))
```

```
# ggsave("figs/plot_k_20-80.png", plot = plot1)
plot1
```

Scatterplot K=20:80



From this first try it's seems that the correct number of K is

2) Second try K = 70:90

```
load("data/results2.Rda")
```

```
top2 <- c(70:90)
top2
```

```
results2 <- data.frame(first=vector(), second=vector(), third=vector())
results2
```

```
system.time(
  for (i in top2)
  {
    set.seed(123)
    lda2 <- LDA(dtm, method= "Gibbs", k = (i), control=list(verbose=50L, iter=100))
    topic <- (i)
    coherence <- mean(topic_coherence(lda2, dtm))
    exclusivity <- mean(topic_exclusivity(lda2))
    results2 <- rbind(results2 , cbind(topic, coherence, exclusivity ))
  }
)
```

```
# save(results2,file="data/results2.Rda")
```

```
kable(results2)
```

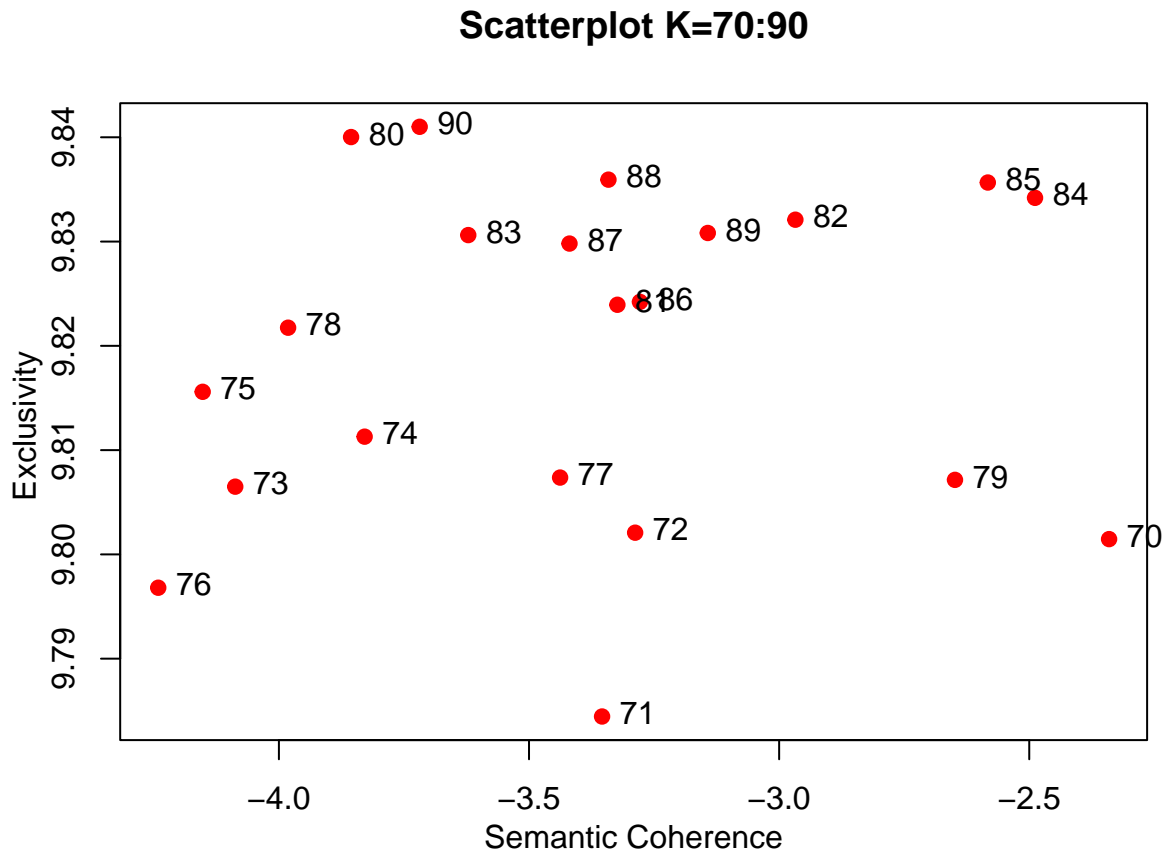
| topic | coherence | exclusivity |
|-------|-----------|-------------|
| 70 | -2.340422 | 9.801467 |
| 71 | -3.353960 | 9.784454 |
| 72 | -3.287822 | 9.802082 |
| 73 | -4.087261 | 9.806496 |
| 74 | -3.828622 | 9.811291 |
| 75 | -4.152354 | 9.815591 |
| 76 | -4.241225 | 9.796804 |
| 77 | -3.437873 | 9.807368 |
| 78 | -3.981590 | 9.821740 |
| 79 | -2.648473 | 9.807149 |
| 80 | -3.855703 | 9.840021 |
| 81 | -3.323339 | 9.823939 |
| 82 | -2.967526 | 9.832096 |
| 83 | -3.621427 | 9.830621 |
| 84 | -2.488446 | 9.834195 |
| 85 | -2.582846 | 9.835661 |
| 86 | -3.278275 | 9.824218 |
| 87 | -3.419067 | 9.829809 |
| 88 | -3.341212 | 9.835939 |
| 89 | -3.142610 | 9.830824 |
| 90 | -3.718580 | 9.841006 |

```
str(results2)
```

```
## 'data.frame':   21 obs. of  3 variables:
## $ topic      : num  70 71 72 73 74 75 76 77 78 79 ...
## $ coherence  : num  -2.34 -3.35 -3.29 -4.09 -3.83 ...
## $ exclusivity: num   9.8 9.78 9.8 9.81 9.81 ...

plot2 <- as.ggplot(~plot(results2$coherence, results2$exclusivity, main="Scatterplot K=70:90",
                          xlab="Semantic Coherence", ylab="Exclusivity ", pch=19,
                          col=ifelse(results2$coherence > -155.3,"red","black"))) +
  text(results2$coherence, results2$exclusivity,
        labels=results2$topic, cex= 1, pos=4))

# ggsave("figs/plot_k_70-90.png", plot = plot2)
plot2
```



In this case ... seems better than ...

3) Third try K = 10:40 with iteration = 1000

```
## Finding the best K
top_k <- c(10:40)
## let's create an empty data frame
results1 <- data.frame(first=vector(), second=vector(), third=vector())
system.time(
  for (i in top_k)
  {
    set.seed(123)
    lda1 <- LDA(dtm, method= "Gibbs", k = (i), control=list(verbose=50L, iter=1000))
    topic <- (i)
    coherence <- mean(topic_coherence(lda1, dtm))
    exclusivity <- mean(topic_exclusivity(lda1))
    results1 <- rbind(results1 , cbind(topic, coherence, exclusivity ))
  }
)
#save(results1,file="data/results_k_10-40.Rda")
```

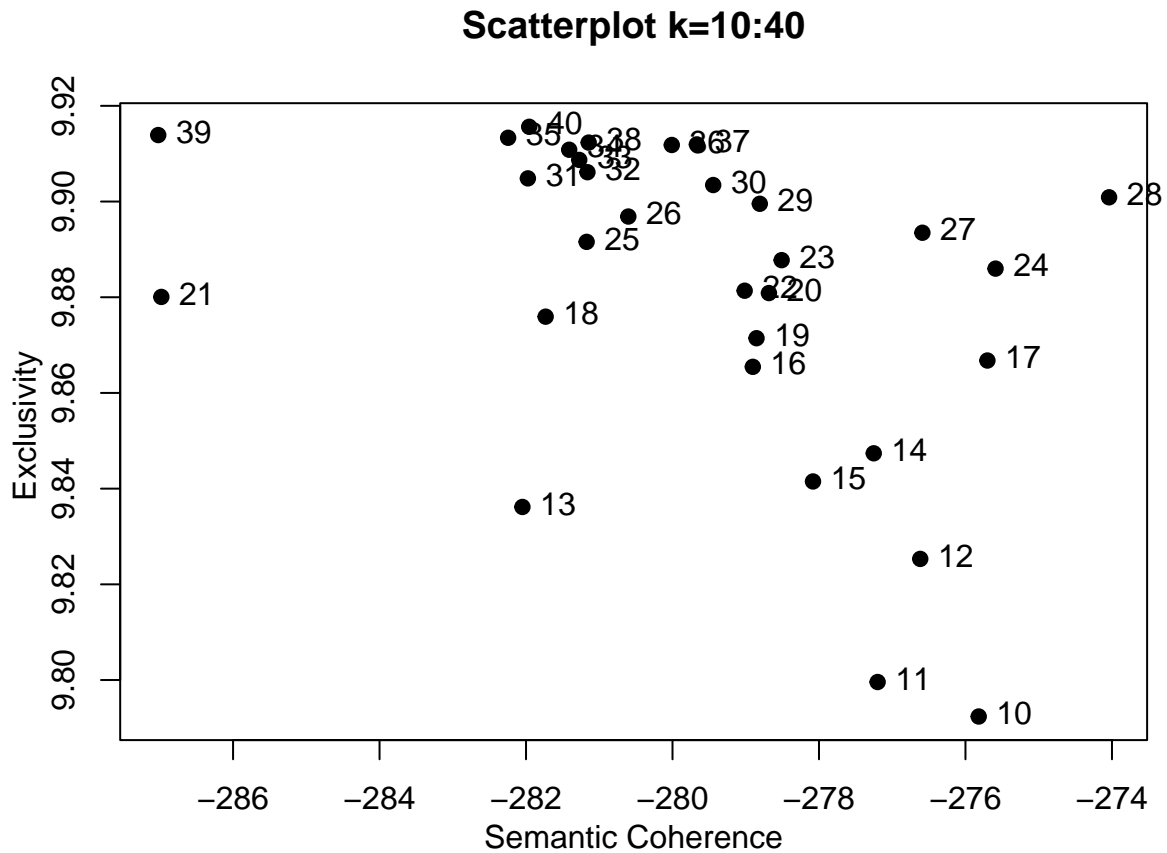
```
kable(results1)
```

| topic | coherence | exclusivity |
|-------|-----------|-------------|
| 10 | -275.8193 | 9.792377 |
| 11 | -277.1995 | 9.799585 |
| 12 | -276.6176 | 9.825317 |
| 13 | -282.0491 | 9.836170 |
| 14 | -277.2525 | 9.847386 |
| 15 | -278.0818 | 9.841498 |
| 16 | -278.9036 | 9.865459 |
| 17 | -275.7004 | 9.866774 |
| 18 | -281.7319 | 9.875940 |
| 19 | -278.8529 | 9.871439 |
| 20 | -278.6834 | 9.880872 |
| 21 | -286.9793 | 9.880067 |
| 22 | -279.0147 | 9.881363 |
| 23 | -278.5099 | 9.887750 |
| 24 | -275.5894 | 9.885981 |
| 25 | -281.1725 | 9.891572 |
| 26 | -280.6034 | 9.896865 |
| 27 | -276.5881 | 9.893478 |
| 28 | -274.0393 | 9.900883 |
| 29 | -278.8082 | 9.899531 |
| 30 | -279.4444 | 9.903447 |
| 31 | -281.9744 | 9.904841 |
| 32 | -281.1614 | 9.906131 |
| 33 | -281.2734 | 9.908701 |
| 34 | -281.4107 | 9.910822 |
| 35 | -282.2445 | 9.913325 |
| 36 | -280.0098 | 9.911821 |
| 37 | -279.6674 | 9.911934 |
| 38 | -281.1443 | 9.912321 |
| 39 | -287.0217 | 9.913895 |
| 40 | -281.9565 | 9.915620 |

```
str(results1)
```

```
## 'data.frame':   31 obs. of  3 variables:
## $ topic       : num  10 11 12 13 14 15 16 17 18 19 ...
## $ coherence   : num  -276 -277 -277 -282 -277 ...
## $ exclusivity: num   9.79 9.8 9.83 9.84 9.85 ...
```

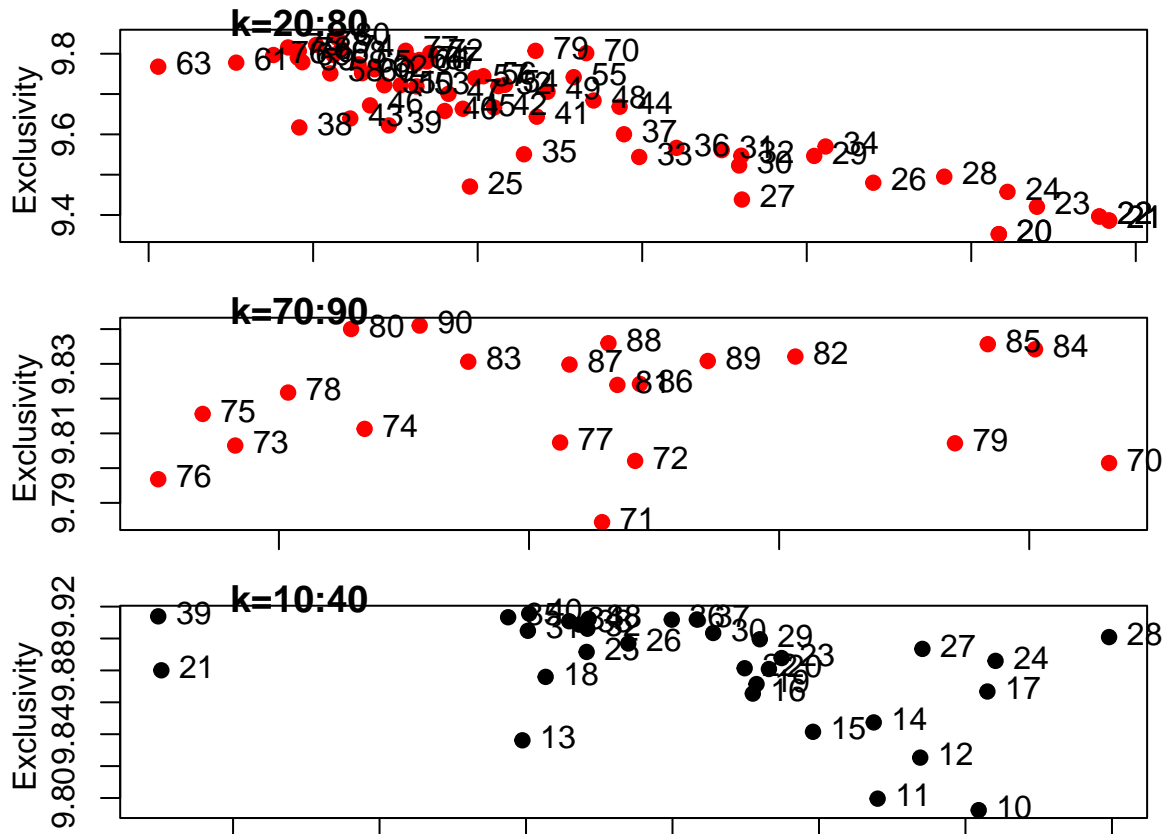
```
## k= 10-40
plot3 <- as.ggplot(~plot(results1$coherence, results1$exclusivity, main="Scatterplot k=10:40",
  xlab="Semantic Coherence", ylab="Exclusivity ", pch=19,
  col=ifelse(results1$coherence<=-153.9 | results1$exclusivity<9.9610 , "black", "red"))) +
  text(results1$coherence,
    results1$exclusivity, labels=results1$topic, cex= 1, pos=4))
plot3
```

```
#ggsave("figs/plot_k_10-40.png", plot = plot3)

grid <- plot_grid(plot1, plot2, plot3,
  labels = list("k=20:80", "k=70:90", "k=10:40"), label_x = .15, nrow = 3)

# ggsave("figs/plot_grid.png", plot = grid)
grid
```



After this try i can state that 30(81) is the best choice

PART III - ANALISYS OF THE TOPICS

```
load("data/lda_k_30.Rda")
```

```
system.time(lda <- LDA(dtm, method= "Gibbs", k = 30, control = list(seed = 123)))
# save(lda, file = "data/lda_k_30.Rda")
```

Here i extract the most important terms from the model

```
terms <- get_terms(lda, 10)
dt1 <- terms[,1:10]
dt2 <- terms[,11:20]
dt3 <- terms[,21:30]

knitr::kable(dt1, col.names = c("Top terms 01","Top terms 02","Top terms 03",
                                "Top terms 04","Top terms 05","Top terms 06","Top terms 07","Top terms 08","Top terms 09","Top terms 10"),
              kable_styling(latex_options = "scale_down"))
```

| Top terms 01 | Top terms 02 | Top terms 03 | Top terms 04 | Top terms 05 | Top terms 06 | Top terms 07 | Top terms 08 | Top terms 09 | Top terms 10 |
|---------------|----------------------------|--------------|--------------|--------------|--------------------|------------------|------------------|------------------|--------------|
| #afghanistan | president | forza | governo | | #iostococonsalvini | #draghi | maggio | the | guerra |
| #tokyo2020 | #quirinal | grand | italiani | donn | governo | draghi | governo | @fratelliditalia | #ucraina |
| talebani | repubblica | buon | + | via | luglio | governo | #decretorilancio | of | pace |
| | #presidentedellarepubblica | pd | lavoro | violenza | #recoveryfund | #governodraghi | #fase2 | to | ucraina |
| agosto | #quirinale2022 | l'italia | #covid19 | giornata | president | lavoro | lavoro | and | putin |
| pass | draghi | politica | crisi | | #cont | paes | ministro | violenza | donn |
| | gennaio | | impres | minacc | #salvini | president | #bonafed | donn | marzo |
| @stampasgarbi | quirinal | c'è | diretta | | legg | @fratelliditalia | #silviaromano | novembr | |
| grazi | grand | casa | momento | mondo | cont | buon | #recoveryfund | covid | russia |
| afghanistan | #mattarella | @pdnetwork | l'italia | pensiero | l'italia | l'italia | ripartir | #morradimett | ucraino |

| Top terms 11 | Top terms 12 | Top terms 13 | Top terms 14 | Top terms 15 | Top terms 16 | Top terms 17 | Top terms 18 | Top terms 19 | Top terms 20 |
|--------------------|--------------|------------------|------------------|--------------|---------------|--------------|------------------|-----------------|--------------|
| #dpcm | pass | maggio | pass | giugno | giugno | governo | governo | #coronavirus | lavoro |
| governo | sindaco | april | draghi | #2giugno | #primalitalia | ministro | cont | #mes | paes |
| #iostococonsalvini | green | vaccinal | natal | scuola | roma | paes | #crisidigoverno | #covid19 | italia |
| ottobr | #greenpass | coprifuoco | green | minacc | bocca | c'è | crisi | #cont | donn |
| #cont | settembr | @stampasgarbi | vaccinati | + | luglio | cittadini | #cont | april | giornata |
| #mes | candidato | @fratelliditalia | dicembr | governo | piazza | president | paes | #forzalombardia | commission |
| covid | città | #nocoprifuoco | | #cont | @stampasgarbi | bene | maggioranza | mes | l'italia |
| @fratelliditalia | draghi | #pnrr | @fratelliditalia | paes | forza | | president | liquidità | italiani |
| de | piazza | pandemia | covid | @luigidimaio | | grazi | #governo | ripartir | #lega |
| jole | roma | draghi | @fattoquotidiano | cont | draghi | parlamento | @fratelliditalia | #fase2 | insiem |

```
knitr::kable(dt2, col.names = c("Top terms 11","Top terms 12","Top terms 13",
                                "Top terms 14","Top terms 15","Top terms 16","Top terms 17","Top terms 18","Top terms 19","Top terms 20"),
  kable_styling(latex_options = "scale_down")
```

```
knitr::kable(dt3, col.names = c("Top terms 21","Top terms 22","Top terms 23",
                                "Top terms 24","Top terms 25","Top terms 26","Top terms 27","Top terms 28","Top terms 29","Top terms 30"),
  kable_styling(latex_options = "scale_down")
```

Using 30 topics I imagined that..... COMMENT HERE

```
titles <- c("1", "2", "3", "4", "5", "6", "7", "8",
            "9", "10","11", "12","13", "14", "15", "16",
            "17", "18", "19", "20",
            "21", "22","23",
            "24", "25", "26", "27", "28", "29", "30")
```

```
table_titles <- rbind (titles, terms)
```

```
t1 <- table_titles[,1:10]
t2 <- table_titles[,11:20]
t3 <- table_titles[,21:30]
```

```
kable(t1)%>%
  kable_styling(latex_options = "scale_down")
```

| Top terms 21 | Top terms 22 | Top terms 23 | Top terms 24 | Top terms 25 | Top terms 26 | Top terms 27 | Top terms 28 | Top terms 29 | Top terms 30 |
|-------------------|----------------------|------------------|---------------|--------------|------------------|------------------|------------------|---------------|----------------------|
| settembr | governo | #sanremo2022 | #coronavirus | grazi | natal | vaccini | governo | pass | #referendumgiustizia |
| #iovotono | #iostococonsalvini | febbraio | grazi | anni | cont | donn | agosto | green | giustizia |
| elettoral | #oggivotolega | draghi | misur | lavoro | governo | buon | anni | ottobr | pass |
| #processateanchem | #salvini | green | l'emergenza | grand | bilancio | @pdnetwork | vittim | #ddlzan | luglio |
| #referendum | salvini | pass | momento | diritti | @fratelliditalia | marzo | settembr | roma | #greenpass |
| parlamentari | #borgonzonipresident | #ucraina | emergenza | via | dicembr | @fratelliditalia | #iovotono | sindaco | gazebo |
| voto | #prescrizion | parlamento | coronavirus | legg | #mes | auguri | @fratelliditalia | legg | #ddlzan |
| @fratelliditalia | @fratelliditalia | @fratelliditalia | casa | president | #natal | draghi | covid | | riforma |
| scuola | #emiliaromagna | guerra | decreto | giovani | italiani | vaccinal | bonus | + | firm |
| referendum | #m5s | #greenpass | #iorestoacasa | città | mes | vaccino | scuola | @forza_italia | draghi |

| | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 | Topic 6 | Topic 7 | Topic 8 | Topic 9 | Topic 10 |
|--------|---------------|----------------------------|------------|----------|----------|--------------------|----------------|------------------|------------------|----------|
| titles | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| | #afghanistan | president | forza | governo | | #iostococonsalvini | #draghi | maggio | the | guerra |
| | #tokyo2020 | #quirinal | grand | italiani | donn | governo | draghi | governo | @fratelliditalia | #ucraina |
| | talebani | repubblica | buon | + | via | luglio | governo | #decretorilancio | of | pace |
| | | #presidentedellarepubblica | pd | lavoro | violenza | #recoveryfund | #governodraghi | #fase2 | to | ucraina |
| | agosto | #quirinale2022 | l'italia | #covid19 | giornata | president | lavoro | lavoro | and | putin |
| | pass | draghi | politica | crisi | | #cont | paes | ministro | violenza | donn |
| | | gennaio | | impres | minacc | #salvini | president | #bonafed | donn | marzo |
| | @stampasgarbi | quirinal | c'è | diretta | legg | @fratelliditalia | #silviaromano | novembr | | |
| | grazi | grand | casa | momento | mondo | cont | buon | #recoveryfund | covid | russia |
| | afghanistan | #mattarella | @pdnetwork | l'italia | pensiero | l'italia | l'italia | ripartir | #morradiomett | ucraino |

| | Topic 11 | Topic 12 | Topic 13 | Topic 14 | Topic 15 | Topic 16 | Topic 17 | Topic 18 | Topic 19 | Topic 20 |
|--------|--------------------|------------|------------------|------------------|--------------|---------------|------------|------------------|-----------------|------------|
| titles | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
| | #dpcm | pass | maggio | pass | giugno | giugno | governo | governo | #coronavirus | lavoro |
| | governo | sindaco | april | draghi | #2giugno | #primalitalia | ministro | cont | #mes | paes |
| | #iostococonsalvini | green | vaccinal | natal | scuola | roma | paes | #crisidigoverno | #covid19 | italia |
| | ottobr | #greenpass | coprifuoco | green | minacc | bocca | c'è | crisi | #cont | donn |
| | #cont | settembr | @stampasgarbi | vaccinati | + | luglio | cittadini | #cont | april | giornata |
| | #mes | candidato | @fratelliditalia | dicembr | governo | piazza | president | paes | #forzalombardia | commission |
| | covid | città | #nocoprifuoco | | #cont | @stampasgarbi | bene | maggioranza | mes | l'italia |
| | @fratelliditalia | draghi | #pnrr | @fratelliditalia | paes | forza | | president | liquidità | italiani |
| | de | piazza | pandemia | covid | @luigidimaio | | grazi | #governo | ripartir | #lega |
| | jole | roma | draghi | @fattoquotidiano | cont | draghi | parlamento | @fratelliditalia | #fase2 | insiem |

```
kable(t2)%>%
  kable_styling(latex_options = "scale_down")
```

```
kable(t3)%>%
  kable_styling(latex_options = "scale_down")
```

This method turns out to be cumbersome and not very informative because for many “topics” it is difficult for me to find a label.

So, I also decided to repeat the search using a much lower K trying to interpret the topics as a films genere.

PART IV - RE-ANALISYS OF THE TOPICS as GENRE

| | Topic 21 | Topic 22 | Topic 23 | Topic 24 | Topic 25 | Topic 26 | Topic 27 | Topic 28 | Topic 29 | Topic 30 |
|--------|-------------------|----------------------|------------------|--------------|-----------|------------------|------------------|------------------|---------------|----------------------|
| titles | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 |
| | settembr | governo | #sanremo2022 | #coronavirus | grazi | natal | vaccini | governo | pass | #referendumgiustizia |
| | #iovtotono | #iostococonsalvini | febbraio | grazi | anni | cont | donn | agosto | green | giustizia |
| | elettoral | #oggiivotolega | draghi | misur | lavoro | governo | buon | anni | ottobr | pass |
| | #processateanchem | #salvini | green | l'emergenza | grand | bilancio | @pdnetwork | vittim | #ddlzan | luglio |
| | #referendum | salvini | pass | momento | diritti | @fratelliditalia | marzo | settembr | roma | #greenpass |
| | parlamentari | #borgonzonipresident | #ucraina | emergenza | via | dicembr | @fratelliditalia | #iovtotono | sindaco | gazebo |
| | voto | #prescrizion | parlamento | coronavirus | legg | #mes | auguri | @fratelliditalia | legg | #ddlzan |
| | @fratelliditalia | @fratelliditalia | @fratelliditalia | casa | president | #natal | draghi | covid | | riforma |
| | scuola | #emiliaromagna | guerra | decreto | giovani | italiani | vaccinal | bonus | + | firm |
| | referendum | #m5s | #greenpass | #iorestocasa | città | mes | vaccino | scuola | @forza_italia | draghi |

I produced this pdf by filling out an R-Markdown, here i paste my complete script:

```
#####  
## EXERCISE 4 --> TOPIC MODEL ##  
#####  
  
library(quanteda.textmodels)  
library(quanteda)  
library(topicmodels)  
library(topicdoc)  
library(cowplot)  
library(ggplot2)  
library(ggplotify)  
library(knitr)  
library(tibble)  
library(kableExtra)  
  
getwd()  
setwd("C:/Users/Riccardo/Documents/UNIVERSITA'/MAGISTRALE/Big Data Analysis/LAB4/assignment_4")  
getwd()  
  
## Open the file with movie reviews inside the quanteda.textmodels package  
data("data_corpus_moviereviews",  
      package = "quanteda.textmodels")  
  
## Create the corpus  
corp <- tail(data_corpus_moviereviews, 500) +  
  head(data_corpus_moviereviews, 500)  
  
## Extract the actual texts of the movie reviews and adding them as docvars in the corpus  
docvars(corp, "texts") <- texts(corp)  
  
## Create the Dfm removing stopwords  
rev_dfm <- dfm(corp, remove = c(stopwords("english"), ("+"),("<"),(">"),  
                                ("rt"), ("00*"), ("="),  
                                ("`"), ("d"), ("r"), ("t"), ("m"), ("l"),  
                                ("b"), ("g"), ("$"), ("j"), ("o"), ("u"),  
                                ("e")), tolower = TRUE, stem = FALSE,  
              remove_punct = TRUE, remove_numbers=TRUE)
```

```

## Timming the Dfm
rev_dfm <- dfm_trim(rev_dfm, min_termfreq = 0.95, termfreq_type = "quantile",
                    max_docfreq = 0.1, docfreq_type = "prop")

kable(topfeatures(rev_dfm, 20))

## Keeping only documents with number of tokens >0
rev_dfm[ntoken(rev_dfm) == 0,]
rev_dfm <- rev_dfm[ntoken(rev_dfm) > 0,]

str(rev_dfm@docvars)

## Convert the Document Feature Matrix (Dfm) in a Topic Model (Dtm)
dtm <- convert(rev_dfm, to = "topicmodels")

#####
## Finding the best K
top1 <- c(20:80)
top1

## let's create an empty data frame
results1 <- data.frame(first=vector(), second=vector(), third=vector())
results1

system.time(
  for (i in top1)
  {
    set.seed(123)
    lda1 <- LDA(dtm, method= "Gibbs", k = (i), control=list(verbose=50L,
                                                             iter=100))

    topic <- (i)
    coherence <- mean(topic_coherence(lda1, dtm))
    exclusivity <- mean(topic_exclusivity(lda1))
    results1 <- rbind(results1 , cbind(topic, coherence, exclusivity ))
  }
)
# save(results1,file="data/results1.Rda")
# load("data/results1.Rda")

kable(head(results1))
kable(tail(results1))
str(results1)

plot1 <- as.ggplot(~plot(results1$coherence, results1$exclusivity, main="Scatterplot K=20:80",
                        xlab="Semantic Coherence", ylab="Exclusivity ", pch=19, col=ifelse(results1$coherence<=-155.8,"black",
                        text(results1$coherence, results1$exclusivity,
                            labels=results1$topic, cex= 1, pos=4))

# ggsave("figs/plot1.png", plot = plot1)
plot1

## From this first try it's seems that the correct number of K is close to 71

```

```
#####
## Now i repeat the procedure with 70:90

top2 <- c(70:90)
top2

results2 <- data.frame(first=vector(), second=vector(), third=vector())
results2

system.time(
  for (i in top2)
  {
    set.seed(123)
    lda2 <- LDA(dtm, method= "Gibbs", k = (i), control=list(verbose=50L,
                                                              iter=100))

    topic <- (i)
    coherence <- mean(topic_coherence(lda2, dtm))
    exclusivity <- mean(topic_exclusivity(lda2))
    results2 <- rbind(results2 , cbind(topic, coherence, exclusivity ))
  }
)

# save(results2,file="data/results2.Rda")
# load("data/results2.Rda")

kable(results2)
str(results2)

plot2 <- as.ggplot(~plot(results2$coherence, results2$exclusivity, main="Scatterplot K=70:90",
                        xlab="Semantic Coherence", ylab="Exclusivity ", pch=19,
                        col=ifelse(results2$coherence > -155.3,"red","black")) +
  text(results2$coherence, results2$exclusivity, labels=results2$topic, cex= 1, pos=
)

# ggsave("figs/plot2.png", plot = plot2)
plot2
## 81 seems better than 71

#####
## now i repeat the cycle with the correct number of iteration (2000) for 75:85

top3 <- c(75:95)
top3

results3 <- data.frame(first=vector(), second=vector(), third=vector())
results3

system.time(
  for (i in top3)
  {
    set.seed(123)
    lda3 <- LDA(dtm, method= "Gibbs", k = (i), control=list(verbose=50L,
                                                              iter=2000))
  }
)
```

```

    topic <- (i)
    coherence <- mean(topic_coherence(lda3, dtm))
    exclusivity <- mean(topic_exclusivity(lda3))
    results3 <- rbind(results3 , cbind(topic, coherence, exclusivity ))
  }
)

# save(results3,file="data/results3.Rda")
# load("data/results3.Rda")

kable(results3)
str(results3)

## k=81
plot3 <- as.ggplot(~plot(results3$coherence, results3$exclusivity, main="Scatterplot k=75:95",
  xlab="Semantic Coherence", ylab="Exclusivity ", pch=19, col=ifelse(results3$coherence<=-153.9 | re

  text(results3$coherence, results3$exclusivity,
    labels=results3$topic, cex= 1, pos=4))
plot3
#ggsave("figs/plot3.png", plot = plot3)

grid <- plot_grid(plot1, plot2, plot3, labels = list("k=20:80", "k=70:90",
  "k=75:95"),
  label_x = .15, nrow = 3)

grid
# ggsave("figs/plot_grid.png", plot = grid)

## After this try i can state that 81 is the best choice

#####
load("data/lda.Rda")
# system.time(lda <- LDA(dtm, method= "Gibbs", k = 81, control = list(seed = 123)))
# save(lda, file = "data/lda.Rda")

## Here i extract the most important terms from the model

terms <- get_terms(lda, 10)

dt1 <- terms[,1:10]
dt2 <- terms[,11:20]
dt3 <- terms[,21:30]
dt4 <- terms[,31:40]
dt5 <- terms[,41:50]
dt6 <- terms[,51:60]
dt7 <- terms[,61:70]
dt8 <- terms[,71:81]

knitr::kable(dt1, col.names = c("Top terms 01","Top terms 02","Top terms 03",
  "Top terms 04","Top terms 05","Top terms 06","Top terms 07","Top terms 08","Top
  kable_styling(latex_options = "scale_down")

```



```

knitr::kable(dt2, col.names = c("Top terms 11","Top terms 12","Top terms 13",
                                "Top terms 14","Top terms 15","Top terms 16","Top terms 17","Top terms 18","Top
                                kable_styling(latex_options = "scale_down")

knitr::kable(dt3, col.names = c("Top terms 21","Top terms 22","Top terms 23",
                                "Top terms 24","Top terms 25","Top terms 26","Top terms 27","Top terms 28","Top
                                kable_styling(latex_options = "scale_down")

knitr::kable(dt4, col.names = c("Top terms 31","Top terms 32","Top terms 33",
                                "Top terms 34","Top terms 35","Top terms 36","Top terms 37","Top terms 38","Top
                                kable_styling(latex_options = "scale_down")

knitr::kable(dt5, col.names = c("Top terms 41","Top terms 42","Top terms 43",
                                "Top terms 44","Top terms 45","Top terms 46","Top terms 47","Top terms 48","Top
                                kable_styling(latex_options = "scale_down")

knitr::kable(dt6, col.names = c("Top terms 51","Top terms 52","Top terms 53",
                                "Top terms 54","Top terms 55","Top terms 56","Top terms 57","Top terms 58","Top
                                kable_styling(latex_options = "scale_down")

knitr::kable(dt7, col.names = c("Top terms 61","Top terms 62","Top terms 63",
                                "Top terms 64","Top terms 65","Top terms 66","Top terms 67","Top terms 68","Top
                                kable_styling(latex_options = "scale_down")

knitr::kable(dt8, col.names = c("Top terms 71","Top terms 72","Top terms 73",
                                "Top terms 74","Top terms 75","Top terms 76","Top terms 77","Top terms 78","Top
                                kable_styling(latex_options = "scale_down")

titles <- c("shakespeare", "armageddon", "3", "the usual suspects", "seven", "shrek", "7", "8", "9", "10",
            "11", "austin powers","griffin", "14", "15", "16", "the batman", "18", "19", "20",
            "godzilla", "22","the grat lebowski", "24", "25", "26", "27", "28", "29", "wonderful christm
            "patch adams", "32","tommy lee jones", "34", "35", "36", "37", "38", "39", "40",
            "mission impossible", "42","ted", "teen", "45", "eddie murphy", "47", "48", "49", "50",
            "the truman show", "52","song", "special effect", "55", "56", "disney cartoon", "john travol
            "61", "62","bill campbell", "64", "Saving private Ryan", "pulp fiction", "67", "aliens", "68",
            "rocky", "72","73", "74", "75", "76", "77", "78", "79", "80", "81"
            )

table_titles <- rbind (titles, terms)

t1 <- table_titles[,1:10]
t2 <- table_titles[,11:20]
t3 <- table_titles[,21:30]
t4 <- table_titles[,31:40]
t5 <- table_titles[,41:50]
t6 <- table_titles[,51:60]
t7 <- table_titles[,61:70]
t8 <- table_titles[,71:81]

kable(t1)%>%
  kable_styling(latex_options = "scale_down")

```

```

kable(t1)%>%
  kable_styling(latex_options = "scale_down")

kable(t1)%>%
  kable_styling(latex_options = "scale_down")

kable(t1)%>%
  kable_styling(latex_options = "scale_down")

kable(t1)%>%
  kable_styling(latex_options = "scale_down")

kable(t1)%>%
  kable_styling(latex_options = "scale_down")

kable(t1)%>%
  kable_styling(latex_options = "scale_down")

#####
## Repeat the procedure with K = 3:15 looking for film genre

genere <- c(3:15)

results.genere <- data.frame(first=vector(), second=vector(), third=vector())

system.time(
  for (i in genere)
  {
    set.seed(123)
    lda.genere <- LDA(dtm, method= "Gibbs", k = (i), control=list(verbose=50L, iter=100))
    topic <- (i)
    coherence <- mean(topic_coherence(lda.genere, dtm))
    exclusivity <- mean(topic_exclusivity(lda.genere))
    results.genere <- rbind(results.genere , cbind(topic,
                                                    coherence, exclusivity ))
  }
)

# save(results.genere,file="data/results_genere.Rda")
# load("data/results_genere.Rda")

plot.genere <- as.ggplot(~plot(results.genere$coherence,
                              results.genere$exclusivity, main="Scatterplot K=3:9",xlab="Semantic Coher",
                              col=ifelse(results.genere$coherence > -155.3,"red","black")) + text(results.genere$coherence,
                              results.genere$exclusivity, labels=results.genere$topic))

# ggsave("figs/plot_genere.png", plot = plot.genere)

plot.genere

```

```

compar <- plot_grid(plot3, plot.genere, label_x = .15, ncol = 2 )

compar
# ggsave("figs/plot_compar.png", plot = compar)

#####
## Looking for topics = genre 6 seems to be the best choice

load("data/lda_genere.Rda")
# system.time(lda.genere <- LDA(dtm, method= "Gibbs", k = 6, control = list(seed = 123)))
# save(lda.genere, file = "data/lda_genere.Rda")

## Extract the 15 most important terms from the model

terms2 <- get_terms(lda.genere, 15)

A <- kable(terms2, col.names = c("Top terms genre1", "Top terms genre2", "Top terms genre3", "Top terms
A

B <-kable(terms2, col.names = c("Action", "Thriller", "Comedy", "Fantasy", "Family", "Romance"))
B

```