# Regulatory documents via LDA (adapted documents)

*Sebastian Knigge*

*18 8 2019*

## 1 Setup

Following libraries are used in the code:

```r
library(dplyr)
library(tidytext)
library(pdftools)
library(tidyr)
library(stringr)
library(tidytext)
library(udpipe)
library(topicmodels)
library(ggplot2)
library(wordcloud)
library(tm)
library(SnowballC)
library(RColorBrewer)
library(RCurl)
library(XML)
library(openxlsx)
library(keras)
```

## 2 Import data

The Documents had to be preprocessed. For the documents wp2.5 all list of contents had to be deleted, because they were the same in each of these documents. No more adjustments had to be made.

In this code reulatory documents are red in and processed via LDA. This first part focusses on reading in the pdf documents.

```r
# getting the right order
setwd('..')
documents <- read.xlsx("Docs_classes.xlsx")[,2]
classes <- read.xlsx("Docs_classes.xlsx")[,c(1,3)]
documents <- paste0(documents,".pdf")
documents %>% as.data.frame() %>% stargazer(summary=FALSE, header = FALSE, title="Document Titles")
```

```r
# getting the right directory
library(here)
setwd("../")
path <- getwd() %>%
  file.path("TextDocs")
setwd(path)
```

Following functions are used to set up and analyze the pdfs.

Table 1: Document Titles

| | . |
|---|---|
| 1 | admin-wp1.1_analysis_legal_institutional_environment_final.pdf |
| 2 | admin-wp1.2_good_practices_final.pdf |
| 3 | admin-wp2.1_estimation_methods1.pdf |
| 4 | admin-wp2.2_estimation_methods2.pdf |
| 5 | admin-wp2.3-estimation_methods3.pdf |
| 6 | admin-wp2.4_examples.pdf |
| 7 | admin-wp2.5_alignment.pdf |
| 8 | admin-wp2.5_editing.pdf |
| 9 | admin-wp2.5_greg.pdf |
| 10 | admin-wp2.5_imputation.pdf |
| 11 | admin-wp2.5_macro_integration.pdf |
| 12 | admin-wp2.5_macro_integration.pdf |
| 13 | admin-wp2.6_good_practices.pdf |
| 14 | admin-wp2.6_guidelines.pdf |
| 15 | admin-wp3.1_quality1.pdf |
| 16 | admin-wp3.2_quality2.pdf |
| 17 | admin-wp3.3_quality.pdf |
| 18 | admin-wp3.4_quality.pdf |
| 19 | admin-wp3.5_quality_measures.pdf |
| 20 | admin-wp3_coherence.pdf |
| 21 | admin-wp3_growth_rates.pdf |
| 22 | admin-wp3_suitability1.pdf |
| 23 | admin-wp3_suitability2.pdf |
| 24 | admin-wp3_suitability3.pdf |
| 25 | admin-wp3_uncertainty.pdf |
| 26 | admin-wp5_frames.pdf |
| 27 | admin-wp5_frames_examples.pdf |
| 28 | admin-wp5_frames_recommendation.pdf |

```r
read_pdf_clean <- function(document){
  # This function loads the document given per name
  # and excludes the stop words inclusive numbers
  pdf1 <- pdf_text(file.path(path, document)) %>%
    strsplit(split = "\n") %>%
    do.call("c",.) %>%
    as_tibble() %>%
    unnest_tokens(word,value) %>%
    # also exclude all words including numbers and special characters
    filter(grepl("^[a-z]+$", word))
  # load stopword library
  data(stop_words)
  # add own words to stop word library - here the numbers from 1 to 10
  new_stop_words <- tibble(word=as.character(0:9),
                           lexicon=rep("own",10)) %>%
                    bind_rows(stop_words)
  # stop words are excluded via anti_join
  pdf1 %>%
    anti_join(new_stop_words)
}

plot_most_freq_words <- function(pdf, n=7){
  # plots a bar plot via ggplot
  pdf %>% count(word) %>% arrange(desc(n)) %>% head(n) %>%
    ggplot(aes(x=word,y=n)) +
    geom_bar(stat="identity")+
    # no labels for x and y scale
    theme(axis.title.y=element_blank(),
          axis.title.x=element_blank())
}
```

Now we can read in all documents in a for loop:

```r
setwd(path)
# inital set up for the corpus
pdf1 <- read_pdf_clean(documents[1])
corpus <- tibble(document=1, word=pdf1$word)
# adding the documents iteratively
for (i in 2:length(documents)){
  pdf_i <- read_pdf_clean(documents[i])
  corpus <- tibble(document=i, word=pdf_i$word) %>% bind_rows(corpus,.)
}
```

# 3   LDA

The LDA model is applied. First the document term matrix has to be set up.

```r
dtm <- corpus %>% count(document, word, sort = TRUE) %>%
  select(doc_id=document, term=word, freq=n) %>%
  document_term_matrix()
c(N,M) %<-% dim(dtm)
```

Using the function LDA sets up the model and prediction/evaluation is done via predict(). But first of all it

shall be verified whether the Predict function actually delivers the same classification as the export of the gamma matrix directly from the LDA model. Therefore both gamma matrices of the single functions are compared. Table 2 displays the output of the gamma matrix received by the predict() function and Table 3 displays the gamma matrix returned by the LDA model itself.

```
set.seed(123)
documents_lda <- LDA(dtm, method = "Gibbs",
                     k = 7, control = list(seed = 1234))

prediction5 <- predict(documents_lda, newdata=dtm, type="topic")

prediction5 <- merge(prediction5, classes, by.x="doc_id", by.y="No")

prediction5 %>%
  select(doc_id,topic_001,topic_002,topic_003,topic_004,topic_005, topic_006, topic_007) %>%
  mutate_each(funs(as.numeric), doc_id,topic_001,topic_002,topic_003,topic_004,topic_005, topic_006, top
  arrange(desc(-doc_id)) %>%
  round(2) %>%
  stargazer(summary=F, rownames = F, header = F, title="Gamma matrix for predict function", label="pred
```

Table 2: Gamma matrix for predict function

| doc_id | topic_001 | topic_002 | topic_003 | topic_004 | topic_005 | topic_006 | topic_007 |
|--------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 1 | 0 | 0.960 | 0 | 0.020 | 0 | 0.020 | 0.010 |
| 2 | 0.010 | 0.790 | 0 | 0.080 | 0.010 | 0.100 | 0.010 |
| 3 | 0.250 | 0.020 | 0.010 | 0.020 | 0.190 | 0.440 | 0.070 |
| 4 | 0.360 | 0.020 | 0.020 | 0.030 | 0.030 | 0.500 | 0.040 |
| 5 | 0.800 | 0.010 | 0.040 | 0.010 | 0.060 | 0.060 | 0.010 |
| 6 | 0.680 | 0.020 | 0.090 | 0.010 | 0.140 | 0.040 | 0.020 |
| 7 | 0.850 | 0 | 0.050 | 0.010 | 0.040 | 0.020 | 0.030 |
| 8 | 0.820 | 0.020 | 0.040 | 0.010 | 0.010 | 0.080 | 0.020 |
| 9 | 0.790 | 0 | 0.050 | 0.030 | 0.050 | 0.050 | 0.030 |
| 10 | 0.930 | 0 | 0.030 | 0 | 0.010 | 0.010 | 0.010 |
| 11 | 0.840 | 0.010 | 0.060 | 0.010 | 0.040 | 0.030 | 0.010 |
| 12 | 0.840 | 0.010 | 0.060 | 0.010 | 0.050 | 0.030 | 0.010 |
| 13 | 0.020 | 0.140 | 0.010 | 0.050 | 0.010 | 0.750 | 0.020 |
| 14 | 0.330 | 0.020 | 0.020 | 0.010 | 0.020 | 0.590 | 0.010 |
| 15 | 0.010 | 0.010 | 0.010 | 0.070 | 0.060 | 0.820 | 0.020 |
| 16 | 0.020 | 0.010 | 0.040 | 0 | 0.880 | 0.040 | 0.010 |
| 17 | 0.020 | 0.010 | 0.040 | 0.010 | 0.090 | 0.060 | 0.770 |
| 18 | 0.020 | 0.010 | 0.030 | 0.010 | 0.890 | 0.040 | 0.010 |
| 19 | 0.090 | 0.010 | 0.040 | 0.030 | 0.670 | 0.100 | 0.050 |
| 20 | 0.040 | 0.010 | 0.040 | 0.720 | 0.050 | 0.130 | 0.010 |
| 21 | 0.010 | 0 | 0.930 | 0 | 0.040 | 0.010 | 0 |
| 22 | 0.010 | 0.010 | 0.030 | 0.010 | 0.930 | 0.010 | 0.010 |
| 23 | 0.010 | 0.010 | 0.060 | 0.020 | 0.850 | 0.020 | 0.040 |
| 24 | 0.020 | 0.010 | 0.040 | 0.010 | 0.890 | 0.020 | 0.010 |
| 25 | 0.070 | 0 | 0.870 | 0.010 | 0.030 | 0.020 | 0.010 |
| 26 | 0.020 | 0.020 | 0 | 0.120 | 0.020 | 0.150 | 0.670 |
| 27 | 0 | 0.150 | 0 | 0.550 | 0.010 | 0.030 | 0.250 |
| 28 | 0.010 | 0.050 | 0.010 | 0.760 | 0.010 | 0.080 | 0.070 |

```
ext_gamma_matrix <- function(model=documents_lda){
  # get gamma matrix for chapter probabilities
  chapters_gamma <- tidy(model, matrix = "gamma")
  # get matrix with probabilities for each topic per chapter
  spreaded_gamma <- chapters_gamma %>% spread(topic, gamma)
  spreaded_gamma %>%
    mutate_each(funs(as.numeric), document,1,2,3,4,5,6,7) %>%
  arrange(desc(-document))
}

ext_gamma_matrix(documents_lda) %>%
  round(2) %>%
  stargazer(summary=F, rownames = F, header=F, title="Gamma matrix extracted from model", label="extrac
```

Table 3: Gamma matrix extracted from model

| document | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 1 | 0 | 0.95 | 0 | 0.02 | 0 | 0.02 | 0.01 |
| 2 | 0 | 0.79 | 0.01 | 0.09 | 0 | 0.1 | 0.01 |
| 3 | 0.26 | 0.02 | 0.02 | 0.02 | 0.19 | 0.44 | 0.06 |
| 4 | 0.36 | 0.01 | 0.02 | 0.04 | 0.03 | 0.5 | 0.05 |
| 5 | 0.78 | 0.01 | 0.04 | 0.01 | 0.07 | 0.08 | 0.02 |
| 6 | 0.67 | 0.03 | 0.09 | 0.02 | 0.14 | 0.04 | 0.01 |
| 7 | 0.82 | 0.01 | 0.05 | 0.01 | 0.05 | 0.03 | 0.04 |
| 8 | 0.8 | 0.02 | 0.03 | 0.02 | 0.02 | 0.08 | 0.03 |
| 9 | 0.75 | 0.01 | 0.06 | 0.02 | 0.06 | 0.07 | 0.03 |
| 10 | 0.92 | 0.01 | 0.03 | 0.01 | 0.01 | 0.01 | 0.02 |
| 11 | 0.85 | 0.01 | 0.06 | 0.01 | 0.05 | 0.02 | 0.01 |
| 12 | 0.83 | 0.01 | 0.07 | 0.01 | 0.04 | 0.03 | 0.01 |
| 13 | 0.02 | 0.15 | 0.01 | 0.06 | 0.01 | 0.74 | 0.02 |
| 14 | 0.34 | 0.02 | 0.02 | 0.02 | 0.02 | 0.57 | 0.01 |
| 15 | 0.01 | 0.02 | 0.01 | 0.08 | 0.07 | 0.8 | 0.02 |
| 16 | 0.02 | 0.01 | 0.03 | 0.01 | 0.88 | 0.04 | 0.01 |
| 17 | 0.03 | 0.01 | 0.04 | 0.02 | 0.1 | 0.05 | 0.75 |
| 18 | 0.03 | 0 | 0.03 | 0 | 0.87 | 0.04 | 0.01 |
| 19 | 0.1 | 0.01 | 0.04 | 0.04 | 0.66 | 0.1 | 0.06 |
| 20 | 0.04 | 0.01 | 0.05 | 0.69 | 0.05 | 0.15 | 0.02 |
| 21 | 0.01 | 0 | 0.91 | 0 | 0.04 | 0.02 | 0 |
| 22 | 0.01 | 0.01 | 0.04 | 0.01 | 0.91 | 0.01 | 0.01 |
| 23 | 0.01 | 0.01 | 0.07 | 0.02 | 0.82 | 0.03 | 0.04 |
| 24 | 0.03 | 0.01 | 0.03 | 0.01 | 0.89 | 0.02 | 0.01 |
| 25 | 0.08 | 0.01 | 0.85 | 0.01 | 0.02 | 0.02 | 0.01 |
| 26 | 0.03 | 0.03 | 0.01 | 0.12 | 0.02 | 0.14 | 0.65 |
| 27 | 0 | 0.16 | 0 | 0.53 | 0.01 | 0.03 | 0.26 |
| 28 | 0.01 | 0.06 | 0.01 | 0.74 | 0.01 | 0.1 | 0.08 |

The tables below summarize which document refers to which topic, according to the LDA model.

# 4    Wordclouds

To check what topics tackle which context, we produce wordclouds using the TFIDF and the TF itself.

Table 4: Documents for Topic 1

| Topic | doc_id | Group |
|-------|--------|-------|
| 1 | 10 | 4 |
| 1 | 11 | 4 |
| 1 | 12 | 4 |
| 1 | 5 | 2 |
| 1 | 6 | 3 |
| 1 | 7 | 4 |
| 1 | 8 | 4 |
| 1 | 9 | 4 |

Table 5: Documents for Topic 2

| Topic | doc_id | Group |
|-------|--------|-------|
| 2 | 1 | 1 |
| 2 | 2 | 1 |

Table 6: Documents for Topic 3

| Topic | doc_id | Group |
|-------|--------|-------|
| 3 | 21 | 6 |
| 3 | 25 | 6 |

Table 7: Documents for Topic 4

| Topic | doc_id | Group |
|-------|--------|-------|
| 4 | 20 | 5 |
| 4 | 27 | 7 |
| 4 | 28 | 7 |

Table 8: Documents for Topic 5

| Topic | doc_id | Group |
|-------|--------|-------|
| 5 | 16 | 5 |
| 5 | 18 | 5 |
| 5 | 19 | 5 |
| 5 | 22 | 6 |
| 5 | 23 | 6 |
| 5 | 24 | 6 |

Table 9: Documents for Topic 6

| Topic | doc_id | Group |
|-------|--------|-------|
| 6 | 13 | 3 |
| 6 | 14 | 4 |
| 6 | 15 | 5 |
| 6 | 3 | 2 |
| 6 | 4 | 2 |

Table 10: Documents for Topic 7

| Topic | doc_id | Group |
|-------|--------|-------|
| 7 | 17 | 5 |
| 7 | 26 | 7 |

```r
plot_wordcloud <- function(corpus, selection="ALL", max.words=25, i, freq="tfidf"){
  # setting up a tibble which returns tfidf and tf and frequency for
  # the whole corpus
  tfidf <- corpus %>% count(document, word, sort = TRUE) %>%
    bind_tf_idf(word, document, n)
  # include all documents for selection if selection="ALL"
  if (all(selection=="ALL")) {
    selection <- corpus %>%
      select(document) %>%
      unique() %>%
      unlist() %>%
      sort()}
  # filter for all selected documents
  # use either ft or tfidf
  if (freq=="tfidf"){
    dtm_selected <- tfidf %>% filter(document%in%selection) %>%
      select(word, tf_idf) %>% count(word, wt=tf_idf, sort=TRUE)
  } else {
    dtm_selected <- tfidf %>% filter(document%in%selection) %>%
      select(word, tf) %>% count(word, wt=tf, sort=TRUE)
  }
  wordcloud(words = dtm_selected$word, freq = dtm_selected$n, min.freq = 1,
            max.words=max.words, random.order=FALSE,
            colors=brewer.pal(8, "Dark2"), scale=c(3,0.2),
            main="Title", use.r.layout = TRUE)
  text(x=0.5, y=1, paste("Topic", i))
}
```

For getting specific and more individual words for each cloud, we use the TFIDF in the first step.

## 4.1  Wordclouds using tfidf

```r
par(mfrow=c(2,4))
par(mar=c(1,1,0.5,1))
```

```r
plot_wordcloud(corpus, selection=ind1[,1], i=1) %>% unlist() %>% as.integer()
```
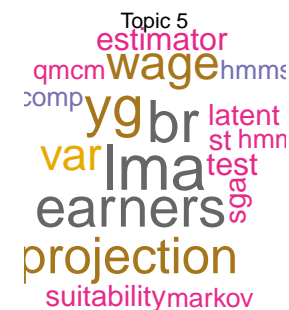
```
## integer(0)
```

```r
plot_wordcloud(corpus, selection=ind2[,1], i=2) %>% unlist() %>% as.integer()
```

```
## integer(0)
```

```r
plot_wordcloud(corpus, selection=ind3[,1], i=3) %>% unlist() %>% as.integer()
```

```
## integer(0)
```

```r
plot_wordcloud(corpus, selection=ind4[,1], i=4) %>% unlist() %>% as.integer()
```

```
## integer(0)
```

```r
plot_wordcloud(corpus, selection=ind5[,1], i=5) %>% unlist() %>% as.integer()
```

```
## integer(0)
```

```r
plot_wordcloud(corpus, selection=ind5[,1], i=6) %>% unlist() %>% as.integer()
```

```
## integer(0)
```

```r
plot_wordcloud(corpus, selection=ind5[,1], i=7) %>% unlist() %>% as.integer()
```

```
## integer(0)
```



## 4.2 Wordclouds using tf

The same can be done using the regular term frequency.

```r
par(mfrow=c(2,4))
par(mar=c(1,1,0.5,1))
```

```
plot_wordcloud(corpus, selection=ind1[,1], i=1, freq="tf")
plot_wordcloud(corpus, selection=ind2[,1], i=2, freq="tf")
plot_wordcloud(corpus, selection=ind3[,1], i=3, freq="tf")
plot_wordcloud(corpus, selection=ind4[,1], i=4, freq="tf")
plot_wordcloud(corpus, selection=ind5[,1], i=5, freq="tf")
plot_wordcloud(corpus, selection=ind5[,1], i=6, freq="tf")
plot_wordcloud(corpus, selection=ind5[,1], i=7, freq="tf")
```





# 5 Embedding via tfidf

Now it's interesting to see if embedding with tfidf will cluster other groups or the same. So we will reduce the Document Term Matrix to M*0.8 words which is a reduction by approx. 20%.

```
dtm_50 <- dtm %>% dtm_remove_tfidf(top=M*0.8)
set.seed(123)
documents_lda_2 <- LDA(dtm_50, method="Gibbs",
                       k = 7, control = list(seed = 1234))
```

```
prediction5_2 <- predict(documents_lda_2, newdata=dtm_50, type="topic")
prediction5_2 <- merge(prediction5_2, classes, by.x="doc_id", by.y="No")
# compare topic 1 with topic 2, 3, 4 and 5
ind1_2 <- prediction5_2 %>% filter(topic==1) %>% select(doc_id, Group)
ind2_2 <- prediction5_2 %>% filter(topic==2) %>% select(doc_id, Group)
ind3_2 <- prediction5_2 %>% filter(topic==3) %>% select(doc_id, Group)
ind4_2 <- prediction5_2 %>% filter(topic==4) %>% select(doc_id, Group)
ind5_2 <- prediction5_2 %>% filter(topic==5) %>% select(doc_id, Group)
ind6_2 <- prediction5_2 %>% filter(topic==6) %>% select(doc_id, Group)
ind7_2 <- prediction5_2 %>% filter(topic==7) %>% select(doc_id, Group)
```

Table 11: Documents for Topic 1

| Topic_embedding_0.8 | doc_id | Group |
|---|---|---|
| 1 | 26 | 7 |
| 1 | 27 | 7 |
| 1 | 28 | 7 |

Table 12: Documents for Topic 2

| Topic_embedding_0.8 | doc_id | Group |
|---|---|---|
| 2 | 17 | 5 |

Table 13: Documents for Topic 3

| Topic_embedding_0.8 | doc_id | Group |
|---|---|---|
| 3 | 16 | 5 |
| 3 | 18 | 5 |
| 3 | 19 | 5 |
| 3 | 21 | 6 |
| 3 | 22 | 6 |
| 3 | 23 | 6 |
| 3 | 24 | 6 |
| 3 | 25 | 6 |

Table 14: Documents for Topic 4

| Topic_embedding_0.8 | doc_id | Group |
|---|---|---|
| 4 | 15 | 5 |
| 4 | 20 | 5 |

Table 15: Documents for Topic 5

| Topic_embedding_0.8 | doc_id | Group |
|---|---|---|
| 5 | 13 | 3 |
| 5 | 14 | 4 |
| 5 | 3 | 2 |
| 5 | 4 | 2 |

Table 16: Documents for Topic 6

| Topic_embedding_0.8 | doc_id | Group |
|---|---|---|
| 6 | 1 | 1 |
| 6 | 2 | 1 |

Table 17: Documents for Topic 7

| Topic_embedding_0.8 | doc_id | Group |
|:---:|:---:|:---:|
| 7 | 10 | 4 |
| 7 | 11 | 4 |
| 7 | 12 | 4 |
| 7 | 5 | 2 |
| 7 | 6 | 3 |
| 7 | 7 | 4 |
| 7 | 8 | 4 |
| 7 | 9 | 4 |

```r
ext_gamma_matrix(documents_lda_2) %>%
  round(2) %>%
  stargazer(summary=F, rownames = F, header=F, title="Gamma matrix extracted from model for embedding wi
```
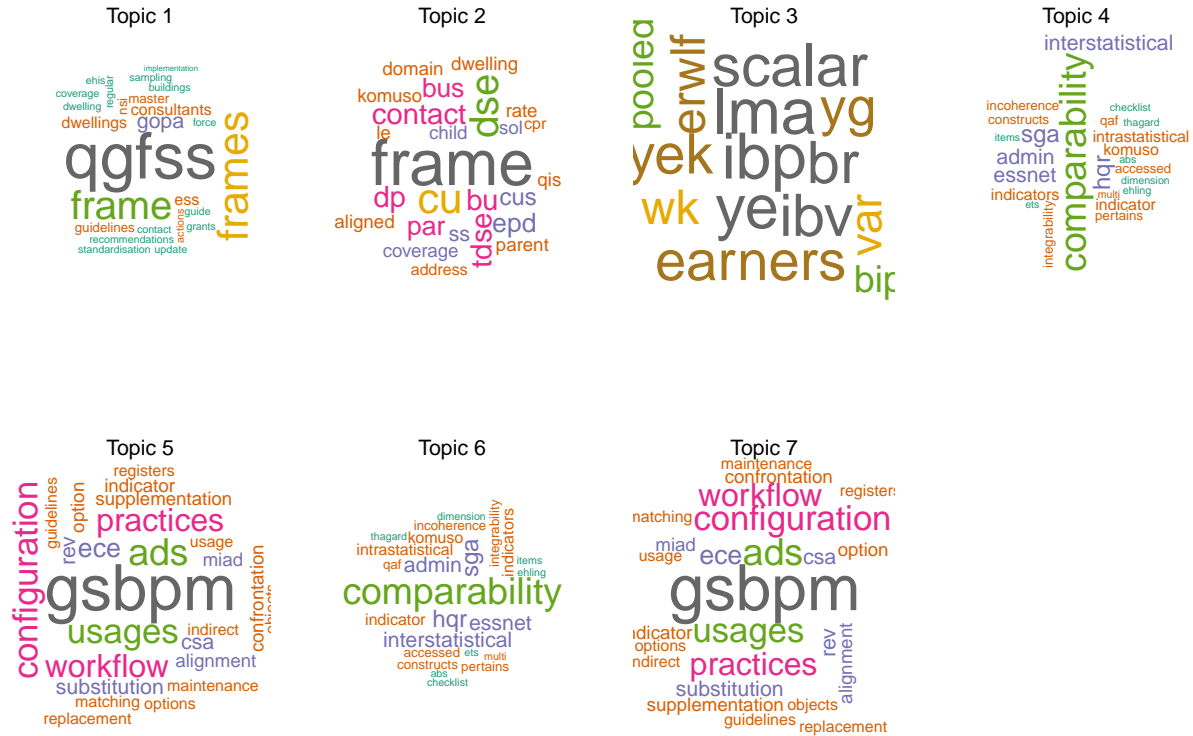
Table 18: Gamma matrix extracted from model for embedding with tfidf

| document | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| 1 | 0.01 | 0 | 0 | 0.01 | 0.01 | 0.97 | 0 |
| 2 | 0.03 | 0.01 | 0.01 | 0.03 | 0.06 | 0.85 | 0 |
| 3 | 0.06 | 0.05 | 0.17 | 0.04 | 0.49 | 0.02 | 0.18 |
| 4 | 0.03 | 0.03 | 0.05 | 0.01 | 0.64 | 0.01 | 0.22 |
| 5 | 0.01 | 0.02 | 0.06 | 0.02 | 0.08 | 0.01 | 0.8 |
| 6 | 0.01 | 0.03 | 0.19 | 0.01 | 0.05 | 0.01 | 0.71 |
| 7 | 0.02 | 0.06 | 0.05 | 0.01 | 0.01 | 0.01 | 0.85 |
| 8 | 0.02 | 0.02 | 0.04 | 0.03 | 0.08 | 0.02 | 0.79 |
| 9 | 0.04 | 0.03 | 0.09 | 0.04 | 0.1 | 0.01 | 0.69 |
| 10 | 0.01 | 0.01 | 0.03 | 0 | 0.01 | 0 | 0.94 |
| 11 | 0.01 | 0.01 | 0.07 | 0.01 | 0.01 | 0.01 | 0.88 |
| 12 | 0.01 | 0.01 | 0.07 | 0.01 | 0.02 | 0.01 | 0.86 |
| 13 | 0.06 | 0.02 | 0.01 | 0.02 | 0.74 | 0.15 | 0.01 |
| 14 | 0.02 | 0.01 | 0.03 | 0.02 | 0.7 | 0.02 | 0.22 |
| 15 | 0.02 | 0.01 | 0.02 | 0.83 | 0.07 | 0.02 | 0.02 |
| 16 | 0.01 | 0.02 | 0.86 | 0.06 | 0.01 | 0.01 | 0.04 |
| 17 | 0.05 | 0.84 | 0.05 | 0.02 | 0.02 | 0.01 | 0.01 |
| 18 | 0.01 | 0.01 | 0.86 | 0.07 | 0.01 | 0 | 0.03 |
| 19 | 0.02 | 0.07 | 0.45 | 0.31 | 0.01 | 0 | 0.13 |
| 20 | 0.02 | 0.03 | 0.03 | 0.85 | 0.02 | 0.01 | 0.03 |
| 21 | 0.01 | 0.01 | 0.94 | 0.01 | 0.01 | 0.01 | 0.01 |
| 22 | 0.02 | 0.01 | 0.93 | 0.01 | 0.02 | 0.01 | 0.01 |
| 23 | 0.03 | 0.04 | 0.85 | 0.03 | 0.02 | 0.01 | 0.01 |
| 24 | 0.02 | 0.01 | 0.91 | 0.01 | 0.03 | 0.01 | 0.02 |
| 25 | 0.01 | 0.01 | 0.81 | 0.01 | 0.01 | 0 | 0.14 |
| 26 | 0.69 | 0.11 | 0.02 | 0.06 | 0.07 | 0.02 | 0.03 |
| 27 | 0.74 | 0.03 | 0.01 | 0.01 | 0.03 | 0.17 | 0.01 |
| 28 | 0.68 | 0.01 | 0.01 | 0.15 | 0.05 | 0.09 | 0.01 |

## 5.1 Wordclouds

```
par(mfrow=c(2,4))
par(mar=c(1,1,0.5,1))
plot_wordcloud(corpus, selection=ind1_2[,1], i=1)
plot_wordcloud(corpus, selection=ind2_2[,1], i=2)
plot_wordcloud(corpus, selection=ind3_2[,1], i=3)
plot_wordcloud(corpus, selection=ind4_2[,1], i=4)
plot_wordcloud(corpus, selection=ind5_2[,1], i=5)
plot_wordcloud(corpus, selection=ind4_2[,1], i=6)
plot_wordcloud(corpus, selection=ind5_2[,1], i=7)
```



We want to give an overview over the clustered documents with the respective wordclouds.

Table 19: Documents for Topic 1

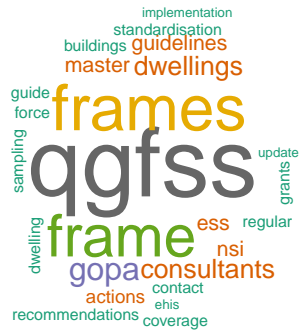| Topic | doc_id | Group |
|-------|--------|-------|
| 1 | 26 | 7 |
| 1 | 27 | 7 |
| 1 | 28 | 7 |

# Topic 1



Table 20: Documents for Topic 2

| Topic | doc_id | Group |
|-------|--------|-------|
| 2 | 17 | 5 |



Table 21: Documents for Topic 3

| Topic | doc_id | Group |
|-------|--------|-------|
| 3 | 16 | 5 |
| 3 | 18 | 5 |
| 3 | 19 | 5 |
| 3 | 21 | 6 |
| 3 | 22 | 6 |
| 3 | 23 | 6 |
| 3 | 24 | 6 |
| 3 | 25 | 6 |

Table 22: Documents for Topic 4

| Topic | doc_id | Group |
|-------|--------|-------|
| 4 | 15 | 5 |
| 4 | 20 | 5 |



Table 23: Documents for Topic 5

| Topic | doc_id | Group |
|-------|--------|-------|
| 5 | 13 | 3 |
| 5 | 14 | 4 |
| 5 | 3 | 2 |
| 5 | 4 | 2 |

Table 24: Documents for Topic 6

| Topic | doc_id | Group |
|-------|--------|-------|
| 6 | 1 | 1 |
| 6 | 2 | 1 |



Table 25: Documents for Topic 7

| Topic | doc_id | Group |
|-------|--------|-------|
| 7 | 10 | 4 |
| 7 | 11 | 4 |
| 7 | 12 | 4 |
| 7 | 5 | 2 |
| 7 | 6 | 3 |
| 7 | 7 | 4 |
| 7 | 8 | 4 |
| 7 | 9 | 4 |

Now we use the validation measure we used for the Example 1.

```r
validate_LDAclassification <- function(predict_table){
  # gamma_matrix is an object of the function ext_gamma_matrix()

  # First we'd find the topic that was most associated with
  # each chapter
  conversion <- predict_table %>%
    select(Group, topic) %>%
    group_by(Group) %>%
    top_n(1,topic) %>%
    unique()

  predict_table %>%
    left_join(conversion, by=c("topic")) %>%
    filter(Group.x!=Group.y) %>%
    nrow()/nrow(predict_table)
}
```

On both full bag of words and 80% embedding via Tfidf

```r
predict_table <- prediction5 %>% select(doc_id, topic) %>%
  merge( y=classes, by.x=1, by.y=1)

validate_LDAclassification(predict_table)
```

```
## [1] 0.5714286
```

```r
predict_table2 <- prediction5_2 %>% select(doc_id, topic) %>%
  merge( y=classes, by.x=1, by.y=1)

validate_LDAclassification(predict_table2)
```

```
## [1] 0.6785714
```