

# LDA Documentation

*Sebastian Knigge*

*6 8 2019*

## Initial Setup

This is the essential step for setting up the LDA model. These functions include the sampling procedure from the *gutenbergr* library

```
# loading packages
library(gutenbergr)
library(dplyr)
library(tidyr)
library(stringr)
library(tidytext)
library(udpipe)
library(topicmodels)
library(ggplot2)

sampling_books <- function(seed=1234, n=20){
  # sample n books from the whole library
  set.seed(seed)
  gutenberg_works() %>%
    # select works with title
    dplyr::filter(!is.na(title)) %>%
    # set the sample size
    sample_n(n) %>%
    # set a special download link
    gutenberg_download(
      mirror = "http://mirrors.xmission.com/gutenberg/"
    )
}
```

- good separation for 4 topics:
  - seed=12345
  - seed=54321
- for 6 books:
  - seed=222
  - seed 101
- for 10 books:
  - seed=54321
  - seed=123456

```
sampling_books <- function(seed=1234, n=20){
  # sample n books from the whole library
  set.seed(seed)
  gutenberg_works() %>%
    # select works with title
    dplyr::filter(!is.na(title)) %>%
    # set the sample size
    sample_n(n) %>%
    # set a special download link
```

```

    gutenbergl_download(
      mirror = "http://mirrors.xmission.com/gutenberg/"
    )
  }

set_up_books <- function(n_books=4, seed=1992){
  # initial book sample
  books <- sampling_books(n=n_books, seed=seed)
  by_chapter <- books %>%
    group_by(gutenberg_id) %>%
    # split in chapters
    mutate(chapter = cumsum(str_detect(text, regex("^chapter ", ignore_case = TRUE)))) %>%
    ungroup() %>%
    # exclude books without chapters
    dplyr::filter(chapter > 0)
  return(by_chapter)
}

shorten_titles <- function(titles){
  # shorten very long book titles by setting
  # a subset of characters of the first line
  # of the title
  sub_inds <- titles %>%
    regex(pattern="\\n|\\r")-1
  sub_inds[sub_inds<0] <- nchar(titles)[sub_inds<0]
  titles %>%
    substr(1,sub_inds)
}

get_titles <- function(x, n_books){
  # get the sampled gutenbergl_ids
  unique_ids <- x %>%
    select(gutenberg_id) %>%
    unique() %>% unlist()
  # get the titles
  titles <- gutenbergl_works() %>%
    dplyr::filter(gutenberg_id %in% unique_ids) %>%
    select(gutenberg_id, title, author) %>%
    mutate(title=shorten_titles(title))
  # get the number of gutenbergl_ids
  len <- nrow(titles)
  if(n_books!=len) warning(paste("---- ",n_books-len,
                                " books have 0 chapters ---- "))

  # the output as a list
  ret <- list(
    titles=titles,
    len=len
  )
  return(ret)
}

append_by_chapter <- function(x=by_chapter, n_books, seed_index=1){
  # append the books matrix until
  # we get the desired number of books n_books

```

```

titles <- get_titles(x, n_books)
n <- titles$len
while (n<n_books) {
  book2add <- sampling_books(n=1, seed=seed_index)
  by_chapter_add <- book2add %>%
    group_by(gutenberg_id) %>%
    # split in chapters
    mutate(chapter = cumsum(str_detect(text, regex("^chapter ", ignore_case = TRUE)))) %>%
    ungroup() %>%
    # exclude books without chapters
    dplyr::filter(chapter > 2)
  titles2add <- get_titles(by_chapter_add, 1)
  # adding the book to by_chapter if there are chapters in the
  # book plus it is not in the data already
  if (titles2add$len==1) if(!titles2add$titles$gutenberg_id%in%titles$titles$gutenberg_id) {
    x <- bind_rows(x, by_chapter_add)
  }
  n<-get_titles(x, n)$len
  seed_index <- seed_index+1
}
return(x)
}

exclude_stop_words <- function(x){
  # unite chapter and document title
  by_chapter_word <- x %>%
    unite(document, gutenberg_id, chapter) %>%
    # split into words
    unnest_tokens(word, text)
  # import tibble stop words
  data(stop_words)
  # find document-word counts
  word_counts <- by_chapter_word %>%
    # exclude stop words
    anti_join(stop_words) %>%
    # count each word by chapter
    count(document, word, sort = TRUE) %>%
    ungroup()
  return(word_counts)
}

convert_to_dtm <- function(x, minfq = 2){
  # get into a format lda can handle
  chapters_dtm <- x %>%
    select(doc_id=document, term=word, freq=n) %>%
    document_term_matrix() %>%
    # reduce by low frequencies
    dtm_remove_lowfreq(minfreq = minfq)
  return(chapters_dtm)
}

convert_to_dtm_2 <- function(x, n=n, minfq = 2, top=10000){

```

```

# get into a format lda can handle
chapters_dtm <- x %>%
  select(doc_id=document, term=word, freq=n) %>%
  document_term_matrix() %>%
  # reduce by low frequencies
  dtm_remove_tfidf(top=top)
return(chapters_dtm)
}

```

Now we can use all these functions to get to the initial corpus sample.

```

n_books <- 6
by_chapter <- set_up_books(n_books=n_books, seed=222)
appended_by_chapter <- append_by_chapter(x=by_chapter, n_books = n_books)
word_counts <- exclude_stop_words(appended_by_chapter)

```

```
## Joining, by = "word"
```

These are the sampled titles for the book sample with the seed 222.

```

titles <- get_titles(appended_by_chapter, n_books)
titles$titles %>% stargazer(summary=FALSE, font.size = "footnotesize", header=FALSE, title="Book-titles")

```

Table 1: Book-titles

gutenberg_id	title	author
11	Alice's Adventures in Wonderland	Carroll, Lewis
3096	Beatrice	Haggard, H. Rider (Henry Rider)
7705	"My Novel" — Volume 04	Lytton, Edward Bulwer Lytton, Baron
25603	Detailed Minutiae of Soldier life in the Army of Northern Virginia, 1861-1865	McCarthy, Carlton
47402	Along Alaska's Great River	Schwatka, Frederick
49675	Hawkins Electrical Guide v. 5 (of 10)	Hawkins, N. (Nehemiah)

In the set up we have another parameter to adjust. The minimum frequency for the bag of words dictionary. Let us set it to 2 in this case, meaning that we include a word only if the frequency is 2 or more.

```

chapters_dtm <- convert_to_dtm(word_counts, minfq=2)
ncol(chapters_dtm)

```

```
## [1] 10685
```

Let us compare it to the case if we include all words.

```

chapters_dtm_all <- convert_to_dtm(word_counts, minfq=0)
ncol(chapters_dtm_all)

```

```
## [1] 17961
```

We also want to compare this to a reduction of the word dictionary by the tfidf. For the sake of comparison the reduction is made to the same value as used above via minfreq=2 (i.e. 10685 words).

```

chapters_dtm_tfidf <- convert_to_dtm_2(word_counts, top=10685)
ncol(chapters_dtm_tfidf)

```

```
## [1] 10685
```

## Applicate the LDA model on the full corpus

```
tim1 <- Sys.time()
chapters_lda <- LDA(chapters_dtm,
                    k = n_books, control = list(seed = 1234))
tim2 <- Sys.time()
u_1 <- tim2-tim1
```

```
tim1 <- Sys.time()
chapters_lda_all <- LDA(chapters_dtm,
                       k = n_books, control = list(seed = 1234))
tim2 <- Sys.time()
u_all <- tim2-tim1
```

```
tim1 <- Sys.time()
chapters_lda_tfidf <- LDA(chapters_dtm,
                          k = n_books, control = list(seed = 1234))
tim2 <- Sys.time()
u_tfidf <- tim2-tim1
chapters_lda
```

## A LDA\_VEM topic model with 6 topics.

Now we evaluate the model all in once:

```
ext_gamma_matrix <- function(model){
  # get gamma matrix for chapter probabilities
  chapters_gamma <- tidy(model, matrix = "gamma")
  # split joint name of book and chapter
  chapters_gamma <- chapters_gamma %>%
    separate(document, c("gutenberg_id", "chapter"), sep = "_", convert = TRUE)
  # get matrix with probabilities for each topic per chapter
  gamma_per_chapter <- chapters_gamma %>%
    spread(topic, gamma)
  return(chapters_gamma)
}
```

```
validate_LDAClassification <- function(x){
  #First we'd find the topic that was most associated with
  # each chapter using top_n(), which is effectively the
  # "classification" of that chapter
  chapter_classifications <- x %>%
    group_by(gutenberg_id, chapter) %>%
    top_n(1, gamma) %>%
    ungroup()

  # We can then compare each to the "consensus"
  # topic for each book (the most common topic among its chapters),
  # and see which were most often misidentified.
  book_topics <- chapter_classifications %>%
    count(gutenberg_id, topic) %>%
    group_by(gutenberg_id) %>%
    # just keep the most frequent one
    top_n(1, n) %>%
    ungroup() %>%

```

```

# keep title called census and topic
transmute(consensus = gutenbergs_id, topic)

# check the fraction of missclassification
chapter_classifications %>%
  inner_join(book_topics, by = "topic") %>%
  # mismatches
  dplyr::filter(gutenbergs_id != consensus)%>%
  nrow()/nrow(chapter_classifications)
}

misc.rate_1 <- ext_gamma_matrix(chapters_lda) %>%
  validate_LDClassification()

misc.rate_all <- ext_gamma_matrix(chapters_lda_all) %>%
  validate_LDClassification()

misc.rate_tfidf <- ext_gamma_matrix(chapters_lda_tfidf) %>%
  validate_LDClassification()

performance_matrix <- data.frame(freq2.embedding=c(misc.rate_1, u_1),
  all.embedding=c(misc.rate_all, u_all),
  tfidf=c(misc.rate_tfidf, u_tfidf))
rownames(performance_matrix) <- c("missc. rate", "time")
performance_matrix %>% stargazer(summary=FALSE, header=F)

```

Table 2:

	freq2.embedding	all.embedding	tfidf
missc. rate	0.636	0.636	0.636
time	15.187	15.038	15.031