

# LDA Documentation

*Sebastian Knigge*

*6 8 2019*

## 1 Setup

Following packages were used in this script:

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

## 2 Example 1 (6 Books)

### 2.1 Get data

#### 2.1.1 Definition of functions

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

```
sampling_books <- function(seed=1234, n=20){
  # sample n books from the whole library
  set.seed(seed)
  gutenbergr_works() %>%
    # select works with title
    dplyr::filter(!is.na(title)&!is.na(gutenbergr_bookshelf)) %>%
    # set the sample size
    sample_n(n) %>%
    # set a special download link
    gutenbergr_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(gutenbergr_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 %>%
    regexr(pattern="\n|\\r")-1
  sub_inds[sub_inds<0] <- nchar(titles)[sub_inds<0]
  sub_inds <- pmin(sub_inds, 45)
  titles %>%
    substr(1,sub_inds)
}

get_titles <- function(x, n_books){
  # get the sampled gutenbergs_ids
  unique_ids <- x %>%
    select(gutenberg_id) %>%
    unique() %>% unlist()
  # get the titles
  titles <- gutenbergs() %>%
    dplyr::filter(gutenberg_id %in% unique_ids) %>%
    select(gutenberg_id, title, author) %>%
    mutate(title=shorten_titles(title))
  # get the number of gutenbergs_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)
}

```

Note:

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

### 2.1.2 Sample corpus

Now we can use all these functions to get to the initial corpus sample. In this example 6 books are chosen.

```
n_books <- 6
by_chapter <- set_up_books(n_books=n_books, seed=222)
get_titles(by_chapter, n_books)

## Warning in get_titles(by_chapter, n_books): --- 5 books have 0 chapters ---
## $titles
## # A tibble: 1 x 3
##   gutenber_id title author
##   <int> <chr> <chr>
## 1      2095 Clotelle: A Tale of the Southern States Brown, William Wells
##
## $len
## [1] 1
```

The function `set_up_books()` returns a warning that several books seem to consist of only one chapter. In order to get a corpus consisting out of 6 books, the function `append_by_chapter()` is used, which fills up the corpus to the desired number of books.

```
appended_by_chapter <- append_by_chapter(x=by_chapter, n_books = n_books)
word_counts <- exclude_stop_words(appended_by_chapter)
```

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

In table 4 the sampled titles for the book sample with the seed 222 are displayed. It appears through the function `append_by:chapter()` one book was added, called “My Novel” — Volume 04“.

```
titles <- get_titles(appended_by_chapter, n_books)
titles$titles %>% stargazer(summary=FALSE, font.size = "footnotesize",
                           header=FALSE, title="Book-titles", rownames=FALSE,
                           label="titles:6books")
```

Table 1: Book-titles

gutenberg_id	title	author
2095	Clotelle: A Tale of the Southern States	Brown, William Wells
6315	The Awakening of Helena Richie	Deland, Margaret Wade Campbell
6971	Judaism	Abrahams, Israel
7635	The Disowned — Volume 05	Lytton, Edward Bulwer Lytton, Baron
10319	Dave Darrin’s Third Year at Annapolis; Or, Le	Hancock, H. Irving (Harrie Irving)
21039	Boycotted, and Other Stories	Reed, Talbot Baines

We also want to check if the book categories are different. See Table (2) for comparison.

```
gbids <- titles$titles$gutenberg_id
categories <- gutenbergs_works() %>%
  filter(gutenberg_id %in% gbids) %>%
  select(gutenberg_id, gutenberg_bookshelf)
categories %>% stargazer(summary=FALSE, font.size = "footnotesize",
                        header=FALSE, title="Book-categories", rownames=FALSE,
                        label="categories:6books")
```

Table 2: Book-categories

gutenberg_id	gutenberg_bookshelf
2095	African American Writers
6315	Bestsellers, American, 1895-1923
6971	Judaism
7635	Historical Fiction
10319	Children's Book Series
21039	School Stories

Obviously the corpus is very diverse. LDA should not have problems to cluster the chapters of the books.

#### 2.1.2.1 Reduction of the dimensionality

In the set up we have another parameter to adjust. The function *convert\_to\_dtm* takes the parameter *minfq*, which is used to reduce the “bag of words” (i.e. dimensionality). *minfq* is the minimum frequency for the bag of words dictionary. I will refer to this as “embedding”. 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] 8597
```

Let us compare it to the case if including all words.

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

```
## [1] 15186
```

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. 8597 words).

```
chapters_dtm_tfidf <- convert_to_dtm_2(word_counts, top=8597)
ncol(chapters_dtm_tfidf)
```

```
## [1] 8597
```

#### 2.1.2.2 Apply the LDA model on the full corpus

Set up the LDA model for the shrunked embedding corpus via frequency 2.

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

In comparison set up the LDA model for the full word embedding corpus.

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

Third LDA setup is for the shrinkage of the dictionary by TFIDF.

```
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("guttenberg_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(guttenberg_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(guttenberg_id, topic) %>%
    group_by(guttenberg_id) %>%
    # just keep the most frequent one
    top_n(1, n) %>%
    ungroup() %>%
    # keep title called census and topic
    transmute(consensus = guttenberg_id, topic)

  # check the fraction of missclassification
  chapter_classifications %>%
```

```

inner_join(book_topics, by = "topic") %>%
# mismatches
dplyr::filter(gutenberg_id != consensus)%>%
nrow()/nrow(chapter_classifications)
}

```

Now we exclude 3 times the beta matrix and evaluate the most likely result with the real results. Depending on how good the LDA will separate the books, this influences the goodness of the fit.

```

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()

```

Following matrix gives an overview:

```

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("misc. rate", "time")
performance_matrix %>% stargazer(summary=FALSE, header=F)

```

Table 3:

	freq2.embedding	all.embedding	tfidf
misc. rate	0.385	0.385	0.385
time	12.861	11.358	11.405

Surprisingly, the run-time to fit the LDA model for the embedding using all words, does not take way longer than the embedding using a lower frequency. The fitting of the reduced bag of words via tfidf takes even longer. Thus from here on the full embedding is used.

## 2.2 Evaluate model on testing set

First we split the data randomly in training and testing samples.

```

split_for_fit <- function(data, test_ratio=0.1, seed=1234){
  # function to set up the training
  # and the testing set from the input data which
  # is an object from the function convert_to_dtm()
  set.seed(seed)
  N <- nrow(data)
  n_test <- (N*test_ratio) %>% ceiling
  test_ind <- sample(1:N, n_test)
  train_ind <- (1:N)[-test_ind]
  ret <- list(train=data[train_ind,],
    test=data[test_ind,])
  return(ret)
}

```

The `fit_n_evaluate()` function will fit the LDA model and evaluate the goodness of fit for an object of the function `split_for_fit`. The section Validation gives insights in the procedure of how the “consensus” is set up.

```
fit_n_evaluate <- function(split, k=n_books){
  LDA_model <- LDA(split$train,
                    k = k, control = list(seed = 1234))
  # use the predict function of udpipe
  # the topic predict function already extract the most likely topics
  prediction <- predict(LDA_model, newdata=split$test) %>% .$topic
  # get "consensus" via maximum likelihood
  # first extract the gamma matrix of the model fitted on the training
  # data
  chapters_gamma <- ext_gamma_matrix(LDA_model)
  spreaded_gamma <- chapters_gamma %>% spread(topic, gamma)
  # get pdfs
  plotm <- spreaded_gamma %>%
    group_by(gutenberg_id) %>%
    # note: pdfs are unnormalized
    summarise_at(2:(titles$len+1), sum)
  topic_link <- plotm %>%
    apply(1, function(x) which.max(x[2:length(x)])) %>%
    cbind(plotm$gutenberg_id) %>%
    as.data.frame()
  # exclude the
  consensus <- split$test %>%
    rownames() %>%
    substr(1,regexpr("_",.)-1) %>%
    as.numeric() %>%
    as.data.frame() %>%
    # merge it to the topic
    merge(topic_link, by.y="V2", sort=FALSE) %>%
    select(..y)
  # missclassification rate will be returned
  sum(consensus!=prediction)/length(prediction)
}
```

We now can evaluate several fits of a model for different splits. Here is a parallelization for the individual for loops applied.

```
# setting up how many cores to be used
useable_cores <- parallel::detectCores() - 1
# registering cluster
cl <- parallel::makeCluster(useable_cores)
doParallel::registerDoParallel(cl)
n <- 59
results <- foreach(i = 1:n, .combine = 'c', .export = ls(.GlobalEnv), .packages = c("dplyr", "udpipe",

chapters_dtm %>%
  split_for_fit(seed=12*i) %>%
  fit_n_evaluate()

}
```



```
## Warning in e$fun(obj, substitute(ex), parent.frame(), e$data): already
## exporting variable(s): chapters_dtm, ext_gamma_matrix, fit_n_evaluate,
## n_books, split_for_fit, titles
parallel::stopCluster(cl)
```

This is the output of the simulation over 59 splits and fits. The mean is the final result:

```
results %>% summary

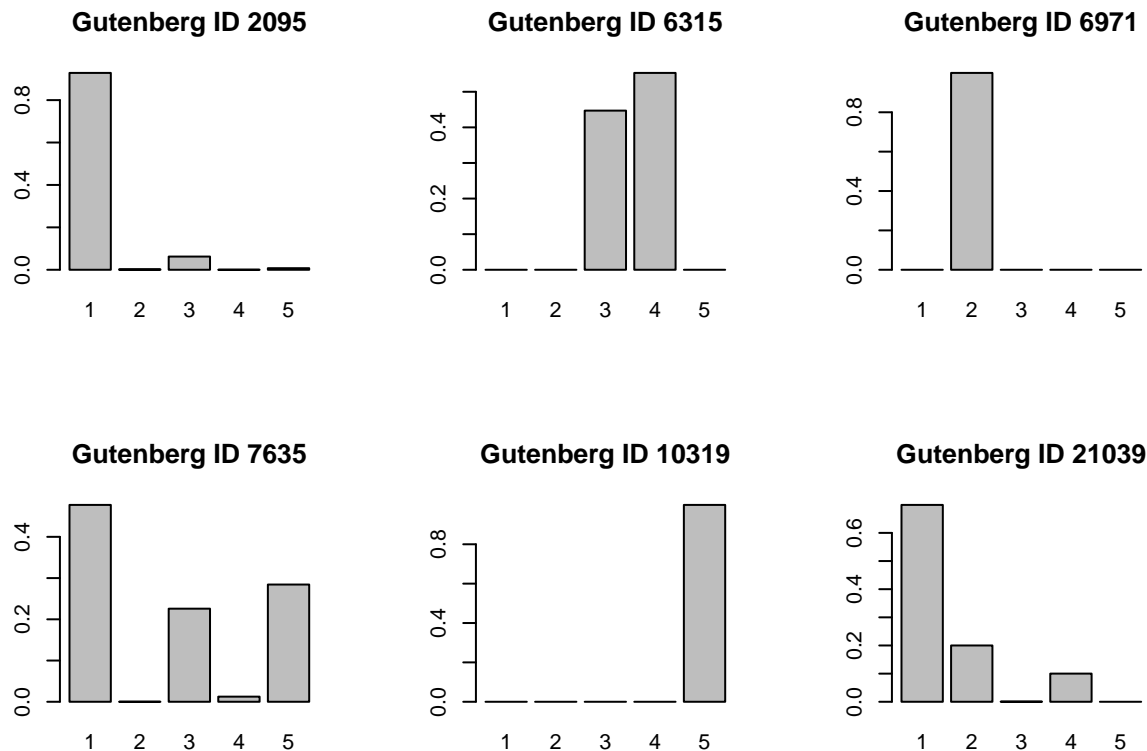
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.2500  0.5000  0.6667  0.6412  0.7500  1.0000
```

### 3 Validation

Since the LDA algorithm just clusters into  $k$  topics, it is necessary to evaluate which “topic” refers to which book, in order to calculate a missclassification rate. For each of the books we naively derive a distribution of the assignment to the topics. This is done by accumulating the distributions for each chapter of the book. The most likely assignment of the LDA model is chosen as the “correct” topic. Now calculating the missclassification rate is basically the fraction of those chapters/documents not coinciding with the most likely assignment.

```
LDA_model <- LDA(chapters_dtm_all,
                 k = n_books, control = list(seed = 1234))
# get "consensus" via maximum likelihood
# first extract the gamma matrix of the model fitted on the training
# data
chapters_gamma <- ext_gamma_matrix(LDA_model)
spreaded_gamma <- chapters_gamma %>% spread(topic, gamma)
# get pdfs
plotm <- spreaded_gamma %>%
  group_by(gutenberg_id) %>%
  # note: pdfs are unnormalized
  summarise_at(2:(titles$len+1), sum)
topic_link <- plotm %>%
  apply(1, function(x) which.max(x[2:length(x)])) %>%
  cbind(plotm$gutenberg_id) %>%
  as.data.frame()

par(mfrow=c(2,3))
for (i in 1:n_books){
  vec <- plotm[i,2:n_books] %>% unlist()
  barplot(vec/sum(vec), main=paste("Gutenberg ID", plotm[i,1]))
}
```



This does not exclude the fact that two books are assigned to the same topic. But still, from each of the plots of the “distributions” you can obtain how good the chapters of a single book are classified. The more of the mass of the distribution is on a single topic, the better the chapters of this topic are predicted.

## 4 Example 2 (10 books)

First we will sample to books again as done in the example 1.

```
n_books <- 10
by_chapter2 <- set_up_books(n_books=n_books, seed=123456)
appended_by_chapter2 <- append_by_chapter(x=by_chapter2, n_books = n_books)
word_counts2 <- exclude_stop_words(appended_by_chapter2)

## Joining, by = "word"

titles <- get_titles(appended_by_chapter2, n_books)
titles$titles %>% stargazer(summary=FALSE, font.size = "footnotesize",
                           header=FALSE, title="Book-titles Example 2", rownames=FALSE,
                           label="titles:6books")

chapters_dtm_all2 <- convert_to_dtm(word_counts2, minfq=0)
```

### 4.1 Evaluate model

The same code as for Example 1 is used to evaluate the model.

```
# registering cluster
cl <- parallel::makeCluster(useable_cores)
doParallel::registerDoParallel(cl)
```

Table 4: Book-titles Example 2

gutenberg_id	title	author
3263	The Portygee	Lincoln, Joseph Crosby
6267	The Weavers: a tale of England and Egypt of f	Parker, Gilbert
6315	The Awakening of Helena Richie	Deland, Margaret Wade Campbell
6971	Judaism	Abrahams, Israel
7492	The Fighting Chance	Chambers, Robert W. (Robert William)
7635	The Disowned — Volume 05	Lytton, Edward Bulwer Lytton, Baron
13154	Lippincott's Magazine of Popular Literature a	Various
13756	Story of Chester Lawrence	Anderson, Nephi
21039	Boycotted, and Other Stories	Reed, Talbot Baines
29685	Submarine Warfare of To-day	Domville-Fife, Charles W. (Charles William)

```
n <- 59
results <- foreach(i = 1:n, .combine = 'c', .export = ls(.GlobalEnv), .packages = c("dplyr", "udpipe",

  chapters_dtm_all2 %>%
    split_for_fit(seed=12*i) %>%
    fit_n_evaluate()

})

## Warning in e$fun(obj, substitute(ex), parent.frame(), e$data): already
## exporting variable(s): chapters_dtm_all2, ext_gamma_matrix, fit_n_evaluate,
## n_books, split_for_fit, titles
parallel::stopCluster(cl)
```

This is the output of the simulation over 59 splits and fits. The mean is the final result:

```
results %>% summary

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.5263  0.7368  0.7895  0.7770  0.8421  0.9474

LDA_model2 <- LDA(chapters_dtm_all2,
                  k = n_books, control = list(seed = 1234))
# get "consensus" via maximum likelihood
# first extract the gamma matrix of the model fitted on the training
# data
chapters_gamma <- ext_gamma_matrix(LDA_model2)
spreaded_gamma <- chapters_gamma %>% spread(topic, gamma)
# get pdfs
plotm <- spreaded_gamma %>%
  group_by(gutenberg_id) %>%
  # note: pdfs are unnormalized
  summarise_at(2:(titles$len+1), sum)
topic_link <- plotm %>%
  apply(1, function(x) which.max(x[2:length(x)])) %>%
  cbind(plotm$gutenberg_id) %>%
  as.data.frame()

par(mfrow=c(2,2))
```

```
for (i in 1:n_books){
  vec <- plotm[i,2:n_books] %>% unlist()
  barplot(vec/sum(vec), main=paste("Gutenberg ID", plotm[i,1]))
}
```

