

Regulatory Documents via Neural Nets - Documentation

Sebastian Knigge

18 9 2019

Contents

- 1 Setup
- 2 Import data
- 3 Apply ANN
 - 3.1 Splitting and Fitting
 - 3.2 Different Architecture
- 4 Training of the most Promising Structure

1 Setup

Following libraries are used in the code:

```
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 like for the LDA part. This first part - again - focusses on reading in the pdf documents.

```
# 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")
```

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

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

When cleaning up data, we have to take into account certain circumstances of the regulatory documents. For example, there are many formulas and technical abbreviations in the documents. Every variable, every estimator, and every index is included as a single word in the bag of words. These terms sometimes have a big influence on the documents, because they are very specific for individual documents and occur quite often. To avoid this, we exclude all mixed words with characters and numeric values, as well as all terms with special characters (e.g. Greek letters).

```
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 stopwords 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:

```
# initial 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 Apply ANN

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

We will reduce the dimensionality again (i.e. performing embedding).

```
dtm_embedding2 <- dtm %>%  
  dtm_remove_lowfreq(minfreq = 2)  
# ordering  
dtm_embedding2 <- dtm_embedding2[order(as.numeric(rownames(dtm_embedding2))),]  
  
dtm_embedding_all <- dtm %>%  
  dtm_remove_lowfreq(minfreq=0)  
# ordering  
dtm_embedding_all <- dtm_embedding_all[order(as.numeric(rownames(dtm_embedding_all))),]  
  
dtm_embedding_tfidf <- dtm %>%  
  dtm_remove_tfidf(top=M*0.5)  
# ordering  
dtm_embedding_tfidf <- dtm_embedding_tfidf[order(as.numeric(rownames(dtm_embedding_tfidf))),]  
  
# convert x matrix into a form such that it can be used for tensorflow  
adjust_tensor_format <- function(classes, dtm){  
  x_chapters <- apply(dtm, 1, function(x) as.matrix(x)) %>% t()  
  # one hot encoding for the Groups (y)  
  topics_categorical <- classes$Group %>% -1 %>%  
    to_categorical()  
  ret <- list(  
    x=x_chapters,  
    y=topics_categorical,  
    topics=classes$Group  
  )  
  return(ret)  
}  
  
dtm_tensor_2 <- adjust_tensor_format(classes, dtm_embedding2)  
dtm_tensor_all <- adjust_tensor_format(classes, dtm_embedding_all)  
dtm_tensor_tfidf <- adjust_tensor_format(classes, dtm_embedding_tfidf)
```

3.1 Splitting and Fitting

The following function splits the sample in an manner, such that each cluster is equally to its size represented in the test data and the validation data.

```
sample_cluster_wise <- function(data, test_ratio=0.2, val_ratio=0.2, seed=1234){  
  X <- data$x; y <- data$y  
  cluster=data$topics
```

```

set.seed(seed)
{
  # setting the absolute number of observations for the sample of each cluster
  n_test <- (table(cluster)*test_ratio) %>% floor()
  n_val <- (table(cluster)*val_ratio) %>% floor()
  # function to get the correct sample indices for validation and test sample
  samp_ind <- function(i, n_list) which(cluster==i) %>% sample(n_list[i])
  test_indices <- unique(cluster) %>% sort() %>%
    sapply(function(i) samp_ind(i, n_test)) %>%
    unlist()
  val_indices <- unique(cluster) %>% sort() %>%
    sapply(function(i) samp_ind(i, n_val)) %>%
    unlist()
}
ret <- list(partial_x_train=X[-c(val_indices,test_indices),],
           partial_y_train=y[-c(val_indices,test_indices),],
           x_val=X[val_indices,], y_val = y[val_indices,],

           x_test=X[test_indices,], y_test = y[test_indices,])
return(ret)
}

```

The following two functions are setting up the model and evaluate the goodness of fit.

```

# The whole model is set up and trained within this function
set_up_n_fit <- function(split, books_n=n_books, epochs=5, batchsize=2^9){
  # starting with 64 neurons and scaling it down to 46 in the
  # mid layer turned out to be a well predicting model
  model <- keras_model_sequential() %>%
    layer_dense(units=64, activation="relu", input_shape=ncol(split$partial_x_train)) %>%
    layer_dense(units=46, activation="relu") %>%
    # we want to classify for 7 categories
    layer_dense(units=7, activation="softmax")

  model %>% compile(
    optimizer="rmsprop",
    loss="categorical_crossentropy",
    metrics=c("accuracy"))

  history <- model %>% fit(
    split$partial_x_train,
    split$partial_y_train,
    # from experience the model tends to
    # overfit for more than 5 epochs
    epochs=epochs,
    batch_size=batchsize,
    validation_data=list(split$x_val,split$y_val)
  )
  return(
    list(history=history,
         model=model))
}

```

```

# making a prediction on the test data and calculating the
# misspecification rate; we also want to save the true categories and the predicted ones
evaluate_model <- function(model_fit, split) {
  y=split$y_test
  x=split$x_test
  prediction <- model_fit %>% predict(x)
  pred <- apply(prediction, 1, which.max)
  true_value <- apply(y, 1, which.max)
  misspecified <- sum(!pred==true_value)/length(pred)
  ret <- list(misspecified=misspecified,
             # the function also dicloses the true and the predicted
             # values for exact evaluation, if needed
             pred=pred, true_value=true_value)
  return(ret)
}

```

Now we can start applying the model. Here we use a split of 20% validation-, 17% test- and 60% training-set. All 3 embeddings are evaluated. We use 49 simulations to test each embedding.

```

tim1 <- Sys.time()
n <- 49
results <- rep(NA,n)
for(i in 1:n){
  split <- sample_cluster_wise(dtm_tensor_2, seed=i*201, test_ratio = 0.17)
  results[i] <- set_up_n_fit(split) %>% .$model %>%
    evaluate_model(split=split) %>% .$misspecified
}
tim2 <- Sys.time()

```

```

u_2 <- tim2-tim1
results_2 <- results
misc.rate_2 <- results_2 %>% mean
var_2 <- results_2 %>% var

```

```

tim1 <- Sys.time()
n <- 49
results <- rep(NA,n)
for(i in 1:n){
  split <- sample_cluster_wise(dtm_tensor_all, seed=i*201, test_ratio = 0.17)
  results[i] <- set_up_n_fit(split) %>% .$model %>%
    evaluate_model(split=split) %>% .$misspecified
}
tim2 <- Sys.time()

```

```

u_all <- tim2-tim1
results_all <- results
misc.rate_all <- results_all %>% mean
var_all <- results_all %>% var

```

```

tim1 <- Sys.time()
n <- 49
results <- rep(NA,n)
for(i in 1:n){
  split <- sample_cluster_wise(dtm_tensor_tfidf, seed=i*201, test_ratio = 0.17)
  results[i] <- set_up_n_fit(split) %>% .$model %>%

```

```

    evaluate_model(split=split) %>% .$misspecified
}
tim2 <- Sys.time()

u_tfidf <- tim2-tim1
results_tfidf <- results
misc.rate_tfidf <- results_tfidf %>% mean
var_tfidf <- results_tfidf %>% var

# overview
performance_matrix <- data.frame(freq2.embedding=c(misc.rate_2, var_2, u_2),
                                all.embedding=c(misc.rate_all, var_all, u_all),
                                tfidf=c(misc.rate_tfidf, var_tfidf, u_tfidf))
rownames(performance_matrix) <- c("missc. rate", "variance", "time")
performance_matrix %>% stargazer(summary=FALSE, header=F, title = "ANN for Reg. Docs.")

```

Table 1: ANN for Reg. Docs.

	freq2.embedding	all.embedding	tfidf
missc. rate	0.418	0.418	0.367
variance	0.108	0.149	0.133
time	2.666	7.668	11.792

3.2 Different Architecture

In this section we want to try a different architecture. Using batch size 2^5 instead of 2^9 In Table 2 the results of this architecture are shown, including the variances of the missclassification rates over all 49 simulations.

```

n <- 49
results_2_d <- rep(NA,n)
for(i in 1:n){
  split <- sample_cluster_wise(dtm_tensor_2, seed=i*201, test_ratio = 0.17)
  results_2_d[i] <- set_up_n_fit(split, batchsize = 2^5) %>% .$model %>%
    evaluate_model(split=split) %>% .$misspecified
}

n <- 49
results_all_d <- rep(NA,n)
for(i in 1:n){
  split <- sample_cluster_wise(dtm_tensor_all, seed=i*201, test_ratio = 0.17)
  results_all_d[i] <- set_up_n_fit(split, batchsize = 2^5) %>% .$model %>%
    evaluate_model(split=split) %>% .$misspecified
}

n <- 49
results_tfidf_d <- rep(NA,n)
for(i in 1:n){
  split <- sample_cluster_wise(dtm_tensor_tfidf, seed=i*201, test_ratio = 0.17)
  results_tfidf_d[i] <- set_up_n_fit(split, batchsize = 2^5) %>% .$model %>%
    evaluate_model(split=split) %>% .$misspecified
}

```

```
# overview
performance_matrix <- data.frame(freq2.embedding=c(results_2_d %>% mean, results_2_d %>% var),
                                all.embedding=c(results_all_d %>% mean, results_all_d %>% var),
                                tfidf=c(results_tfidf_d %>% mean, results %>% var))
rownames(performance_matrix) <- c("misc. rate", "variance")
performance_matrix %>% stargazer(summary=FALSE, header=F, title = "ANN for Reg. Docs. Different Archite
```

Table 2: ANN for Reg. Docs. Different Architecture (7 epochs)

	freq2.embedding	all.embedding	tfidf
misc. rate	0.418	0.418	0.388
variance	0.118	0.118	0.125

4 Training of the most Promising Structure

In the course of training and optimizing, many hundret models were fitted. In the following I want to show the training-plot of the most promising approach of the ANN for the regulatory documents.

```
{
set.seed(1234)

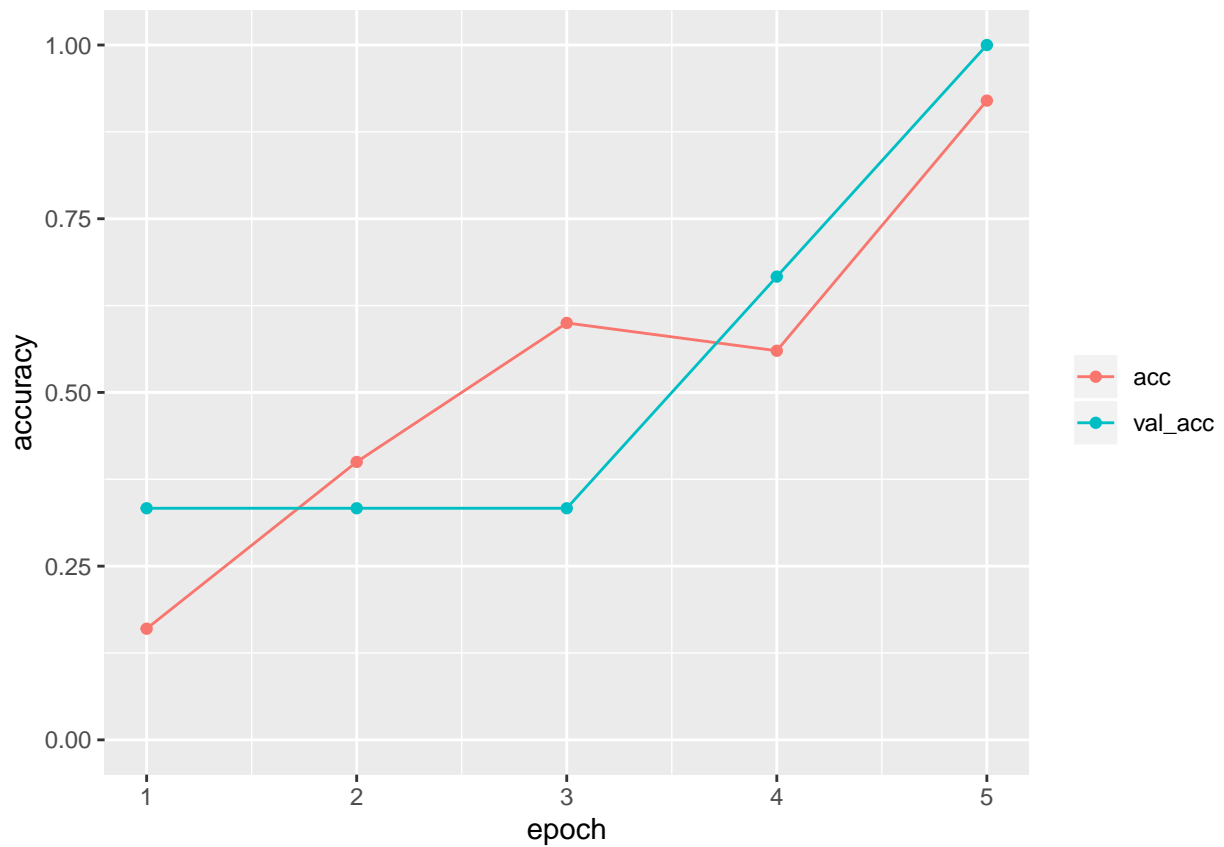
split <- sample_cluster_wise(dtm_tensor_2, test_ratio=0)
hist <- set_up_n_fit(split, batchsize = 2^9) %>% .$history

plot_df <- hist$metrics %>% do.call("cbind",.) %>% cbind(.,1:5) %>% as.data.frame()

plot_df <- plot_df %>% gather(value, key=measure, acc, loss, val_acc, val_loss)

}

ggplot(subset(plot_df, measure %in% c("acc", "val_acc")), aes(V5, value)) +
  geom_line(aes(color=factor(measure))) +
  geom_point(aes(color=factor(measure))) +
  ylim(c(0,1))+
  xlab("epoch")+
  ylab("accuracy")+
  theme(legend.title = element_blank())
```



```
ggplot(subset(plot_df, measure %in% c("loss", "val_loss")), aes(V5, value)) +  
  geom_line(aes(color=factor(measure))) +  
  geom_point(aes(color=factor(measure))) +  
  xlab("epoch") +  
  ylab("loss") +  
  theme(legend.title = element_blank())
```