

Test

Andreas Methling

26/11/2021

```
library(readr)
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5      v dplyr   1.0.7
## v tibble  3.1.4      v stringr 1.4.0
## v tidyr   1.1.3      v forcats 0.5.1
## v purrr   0.3.4
```

```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
#data_start <- read_csv("C:/Users/Simon ik mig/Downloads/lyrics-data.csv.zip") #Simon
#artists_data <- read_csv("C:/Users/Simon ik mig/Downloads/artists-data (1).csv")# Simon
data_start <- read_csv("C:/Users/andre/Desktop/lyrics-data.csv")
```

```
## Rows: 209522 Columns: 5
```

```
## -- Column specification -----
## Delimiter: ","
## chr (5): ALink, SName, SLink, Lyric, Idiom
```

```
##
## 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.
```

```
artists_data <- read_csv("C:/Users/andre/Downloads/artists-data.csv")
```

```
## Rows: 3242 Columns: 6
```

```
## -- Column specification -----
## Delimiter: ","
## chr (4): Artist, Link, Genre, Genres
## dbl (2): Songs, Popularity
```

```
##
## 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.
```

Making the artist data ready for import

```
artists = artists_data %>%
  group_by(Artist) %>%
  count(Genre) %>%
  pivot_wider(names_from = Genre, values_from = n) %>%
  replace_na(list(Pop = 0, "Hip Hop" = 0, Rock = 0, "Funk Carioca" = 0, "Sertanejo" = 0, Samba = 0)) %>%
  ungroup() %>%
  left_join(artists_data, by = c("Artist")) %>%
  select(-c(Genre, Genres, Popularity)) %>%
  distinct()
```

Combining the data and only keeping a column with the lyrics and a combined band name song name column.

```
glimpse(data_start)
```

```
## Rows: 209,522
## Columns: 5
## $ ALink <chr> "/10000-maniacs/", "/10000-maniacs/", "/10000-maniacs/", "/10000~
## $ SName <chr> "More Than This", "Because The Night", "These Are Days", "A Camp~
## $ SLink <chr> "/10000-maniacs/more-than-this.html", "/10000-maniacs/because-th~
## $ Lyric <chr> "I could feel at the time. There was no way of knowing. Fallen l~
## $ Idiom <chr> "ENGLISH", "ENGLISH", "ENGLISH", "ENGLISH", "ENGLISH", "ENGLISH"~
```

```
data = data_start %>%
  filter(Idiom == "ENGLISH") %>%
  rename("Link" = "ALink") %>%
  inner_join(artists, by = c("Link")) %>%
  distinct() %>%
  mutate(name = paste(Artist, SName)) %>%
  rename(text=Lyric) %>%
  filter(Rock==1 & Pop==1) %>%
  select(name, text) %>%
  distinct(name, .keep_all = T)
```

```
data %>%
  count(name, sort = T)
```

```
## # A tibble: 3,348 x 2
##   name                                     n
##   <chr>                                <int>
## 1 10000 Maniacs A Campfire Song         1
## 2 10000 Maniacs A Room For Everything   1
## 3 10000 Maniacs Across The Fields      1
## 4 10000 Maniacs All That Never Happens 1
## 5 10000 Maniacs Among The Americans    1
## 6 10000 Maniacs Angels, From The Realms Of Glory 1
```

```
## 7 10000 Maniacs Anthem For Doomed Youth 1
## 8 10000 Maniacs Arbor Day 1
## 9 10000 Maniacs Back O' The Moon 1
## 10 10000 Maniacs Because The Night 1
## # ... with 3,338 more rows
```

Make labels

Our data set didnt contain labels so we made them ourselves by first tokenizing

```
library(tidytext)
text_tidy = data %>% unnest_tokens(word, text, token = "words")

head(text_tidy)
```

```
## # A tibble: 6 x 2
##   name                word
##   <chr>              <chr>
## 1 10000 Maniacs More Than This i
## 2 10000 Maniacs More Than This could
## 3 10000 Maniacs More Than This feel
## 4 10000 Maniacs More Than This at
## 5 10000 Maniacs More Than This the
## 6 10000 Maniacs More Than This time
```

We remove stopwords and words less than two words.

```
text_tidy %<%>%
  filter(str_length(word) > 2 ) %>%
  group_by(word) %>%
  ungroup() %>%
  anti_join(stop_words, by = 'word')
```

Then we stem our words.

```
library(hunspell)
text_tidy %>%
  mutate(stem = hunspell_stem(word)) %>%
  unnest(stem) %>%
  count(stem, sort = TRUE)
```

```
## # A tibble: 8,047 x 2
##   stem      n
##   <chr> <int>
## 1 love   5789
## 2 time   3466
## 3 feel   3005
## 4 yeah   2512
## 5 baby   2215
## 6 gonna  2166
```

```
## 7 day      2042
## 8 wanna    1944
## 9 life     1769
## 10 heart   1730
## # ... with 8,037 more rows
```

```
text_tidy %<>%
  mutate(stem = hunspell_stem(word)) %>%
  unnest(stem) %>%
  select(-word) %>%
  rename(word = stem)
```

Then we take the 10000 top words, but our data set after preprocessing only contains 8047 words so we move forward with them

```
top_10000_words=text_tidy %>%
  count(word,sort = T) %>%
  head(10000) %>%
  select(word)

data_top_10000=top_10000_words %>%
  left_join(text_tidy, by= c("word"))
```

Bing

Then we make our sentiment labels using the bing dictionary by giving every word a sentiment of either 0 for negative or 1 for positive and then we group by song and summarise the sentiment of every word in a song to get the mean, which then becomes the label of the song.

```
sentiment_bing= data_top_10000 %>%
  inner_join(get_sentiments("bing")) %>%
  mutate(sentiment= ifelse(sentiment == "positive", 1,0))
```

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

```
sentiment_bing %<>%
  group_by(name) %>%
  summarise(mean= mean(sentiment))%>%
  mutate(label= ifelse(mean>=0.5, 1,0))
```

Afinn

We have also made labels for the data using the Afinn dictionary, which give the words a value from -5 to 5 so the words can now also be very positive or negative. Then we do the same and find the mean of every song to get the labels.

```
sentiment_afinn= data_top_10000 %>%
  inner_join(get_sentiments("afinn"))
```

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

```
sentiment_afinn %<>%  
  group_by(name) %>%  
  summarise(mean= mean(value))%>%  
  mutate(label= ifelse(mean>=0, 1,0))
```

Data

```
data_bing= sentiment_bing %>%  
  inner_join(data)%>%  
  select(text, label, name)
```

```
## Joining, by = "name"
```

```
data_afinn= sentiment_afinn %>%  
  inner_join(data)%>%  
  select(text, label, name)
```

```
## Joining, by = "name"
```

Neural Network Bing

Using our bing data set

```
data_bing_n= data_bing %>%  
  rename( y = label )
```

We start by creating test and training data

```
library(rsample)  
  
split= initial_split(data_bing_n, prop = 0.75)  
  
train_data= training(split)  
test_data= testing(split)
```

```
x_train_data= train_data %>% pull(text)  
y_train_data= train_data %>% pull(y)  
  
x_test_data= test_data %>% pull(text)  
y_test_data= test_data %>% pull(y)
```

Now it is time to load keras and make some adjustments to the data. The data are lyrics so not a lot of special characters are used but we still remove them just to be sure. And then we tokenize our data as we know from basic machine learning to get like a bag of words from our song lyrics and lastly we create a list where every song has a vector which includes the words as a numerical character if the words contained in the tweets are among the 100000 most used words in the data set.

```
library(keras)

#for training data
tokenizer_train <- text_tokenizer(num_words = 5000,
                                   filters = "!\"#$%&()*+,-./:;<=>?@[\\]^_`{|}~\t\n" ) %>%
  fit_text_tokenizer(x_train_data)

sequences_train = texts_to_sequences(tokenizer_train, x_train_data)

#For test data
tokenizer_test <- text_tokenizer(num_words = 5000,
                                  filters = "!\"#$%&()*+,-./:;<=>?@[\\]^_`{|}~\t\n" ) %>%
  fit_text_tokenizer(x_test_data)

sequences_test = texts_to_sequences(tokenizer_test, x_test_data)
```

Baseline model

One-hot encoding

we use this function Daniel made to vectorize the sequences :

```
vectorize_sequences <- function(sequences, dimension) {
  results <- matrix(0, nrow = length(sequences), ncol = dimension)
  for(i in 1:length(sequences)){
    results[i, sequences[[i]]] <- 1
  }
  return(results)
}
```

we use it on the training and test data

```
x_train <- sequences_train %>% vectorize_sequences(dimension = 5000)
x_test <- sequences_test %>% vectorize_sequences(dimension = 5000)

str(x_train[1,])
```

```
## num [1:5000] 1 1 1 1 1 0 0 1 1 1 ...
```

What the above has done to the data is, that every tweet now is a row and every feature/word now is a column and then if the tweets has e.g. word 1 then it would have the value 1 otherwise zero. So we basically now have a matrix of size [2488x5000] [number of song in training set x number of words].

The model

The above data is then used in our baseline model with an input shape of 5000 because that is the size of our input. Then we run it through two dense “relu” layers which is normal procedure for a baseline model. Lastly we have a dense layer with the output which is of unit 1 and is a “sigmoid” layer which means it returns a value between 0 and 1 as we want, as we wanna figure out if a song is positive or negative.

```
model_keras <- keras_model_sequential()

model <- model_keras %>%
  layer_dense(units = 128, activation = "relu", input_shape = c(5000)) %>%
  layer_dense(units = 128, activation = "relu") %>%
  layer_dense(units = 1, activation = "sigmoid")
```

We use baseline model compiling with optimizer “adam”, loss “binary” as we are dealing with a binary case and the metric we wanna maximize is accuracy.

```
model %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy",
  metrics = "accuracy"
)
```

Here the structure of the model can be viewed, where it can be seen that the model has 656769 tunable parameters, so not the biggest of models but not the smallest either.

```
summary(model)
```

```
## Model: "sequential"
## -----
## Layer (type)                Output Shape          Param #
## -----
## dense_2 (Dense)             (None, 128)           640128
## -----
## dense_1 (Dense)             (None, 128)           16512
## -----
## dense (Dense)                (None, 1)              129
## -----
## Total params: 656,769
## Trainable params: 656,769
## Non-trainable params: 0
## -----
```

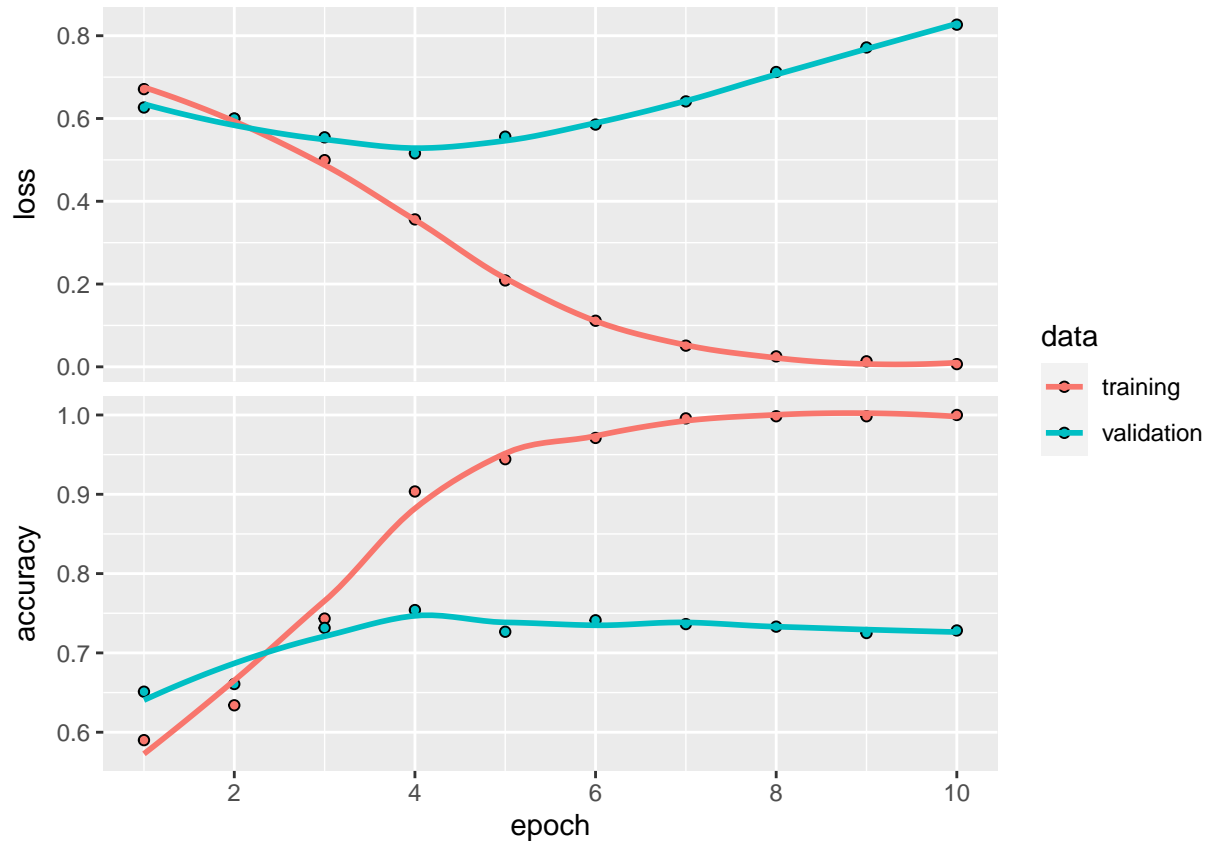
And now the model is run 10 times with a batch size of 256

```
set.seed(12345)
history_ann <- model %>% fit(
  x_train,
  y_train_data,
  epochs = 10,
  batch_size = 256,
  validation_split = 0.25
)
```

We then plot the result of the model

```
plot(history_ann)
```

```
## 'geom_smooth()' using formula 'y ~ x'
```



The top graph shows the loss of the model, where the blue line is the loss of the validation set, which initially falls a bit but then it rises back up. The lower graph shows the same, that the moment the loss rises in the validation set the accuracy falls again. The training set does a lot better than the validation set, which is an indicator of, that our model is over fitted.

Running our model on our test data also shows a bad result

```
metrics = model %>% evaluate(x_test, y_test_data); metrics
```

```
##      loss  accuracy
## 1.8711323 0.5421687
```

We will now try to tune the model to get a better result and prevent the over fitting.

Model tuning

We introduce some dropout layers and reduce the weights of each layer to minimize the number of parameters in the model to prevent the overfitting we saw above.

```
model_keras <- keras_model_sequential()

model2 <- model_keras %>%
  layer_dense(units = 16, activation = "relu", input_shape = c(5000)) %>%
  layer_dropout(rate = 0.1) %>%
  layer_dense(units = 16, activation = "relu") %>%
```



```
layer_dropout(rate = 0.1) %>%
layer_dense(units = 1, activation = "sigmoid")
```

We use baseline model compiling with optimizer “adam”, loss “binary” as we are dealing with a binary case and the metric we wanna maximize is accuracy.

```
model2 %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy",
  metrics = "accuracy"
)
```

Here the structure of the model can be viewed, where it can be seen that the model has 656769 tunable parameters, so not the biggest of models but not the smallest either.

```
summary(model2)
```

```
## Model: "sequential_1"
## -----
## Layer (type)                Output Shape          Param #
## -----
## dense_5 (Dense)              (None, 16)            80016
## -----
## dropout_1 (Dropout)          (None, 16)            0
## -----
## dense_4 (Dense)              (None, 16)            272
## -----
## dropout (Dropout)            (None, 16)            0
## -----
## dense_3 (Dense)              (None, 1)             17
## -----
## Total params: 80,305
## Trainable params: 80,305
## Non-trainable params: 0
## -----
```

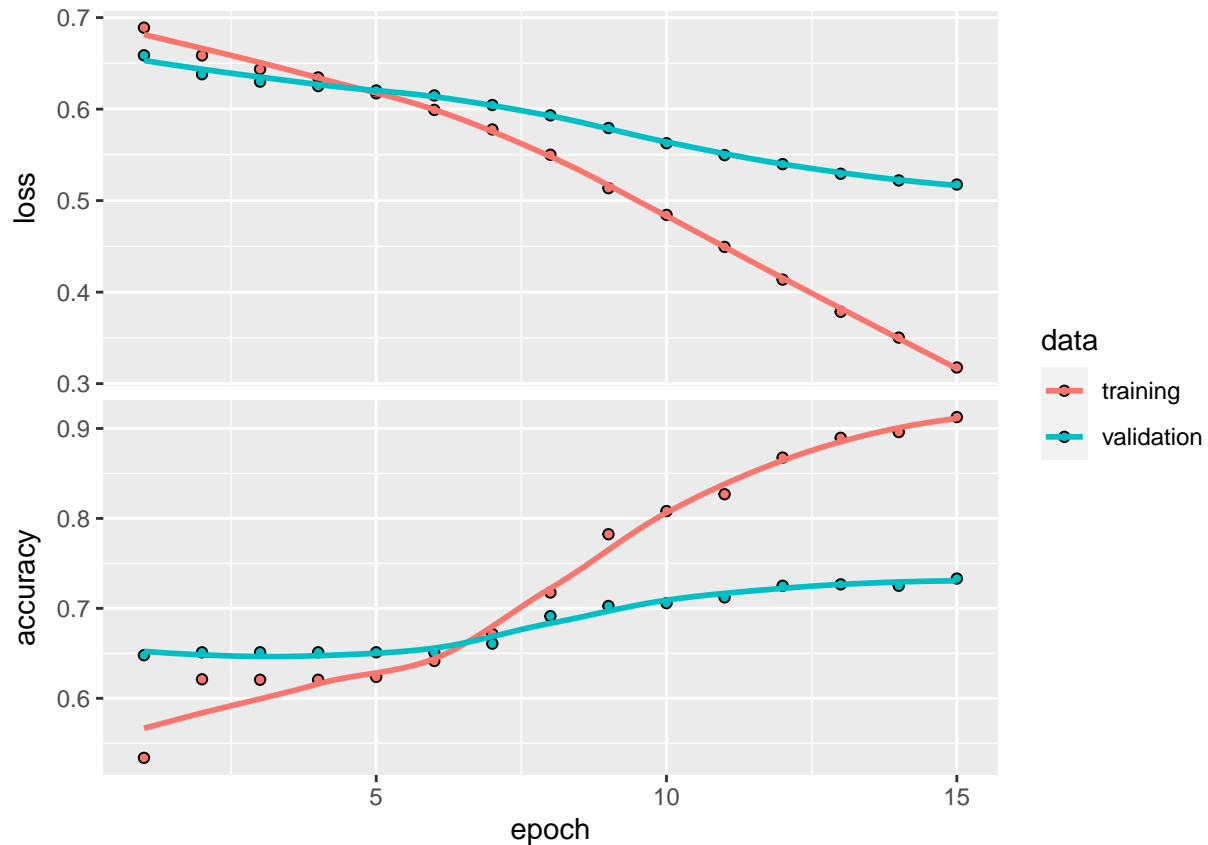
And now the model is run 10 times with a batch size of 512, so a bigger batch size than the baseline model.

```
set.seed(12345)
history_ann2 <- model2 %>% fit(
  x_train,
  y_train_data,
  epochs = 15,
  batch_size = 512,
  validation_split = 0.25
)
```

We then plot the result of the tunned model

```
plot(history_ann2)
```

```
## 'geom_smooth()' using formula 'y ~ x'
```



The model we have tried to tune seems better as we can see the loss function both of the validation data and the training data continues to go down after every epoch. The accuracy also increases in every period until it stalls a bit but still this model looks better than the baseline model.

```
metrics2 = model2 %>% evaluate(x_test, y_test_data); metrics2
```

```
##      loss  accuracy
## 0.7530288 0.5409638
```

The accuracy of this model is better than the baseline model, but with an accuracy of 59% it is still not any good.

Rnn model with padded data

Padding

In our first baseline model, we used a document-term matrix as inputs for training, with one-hot-encodings (= dummy variables) for the 10.000 most popular terms. This has a couple of disadvantages. Besides being a very large and sparse vector for every review, as a “bag-of-words”, it did not take the word-order (sequence) into account.

This time, we use a different approach, therefore also need a different input data-structure. We now use `pad_sequences()` to create a integer tensor of shape (samples, word_indices). However, song vary in length, which is a problem since Keras requires the inputs to have the same shape across the whole sample. Therefore, we use the `maxlen = 300` argument, to restrict ourselves to the first 300 words in every song.

The data is padded

```
x_train_pad <- sequences_train %>% pad_sequences(maxlen=300)
x_test_pad <- sequences_test %>% pad_sequences(maxlen=300)
```

```
glimpse(x_train_pad)
```

```
## num [1:2488, 1:300] 0 0 0 0 590 0 0 5 0 174 ...
```

Now if the value in e.g. the first column of the first tweet is 0 it means that the first word in the first tweet is not one of the 100000 most used words and there for our model has no integer for it. If there is an integer e.g. “386” it means that that the 386 most commonly used word is the first word in the tweet.

The model setting up the model we will first use a layer_embedding to compress our initial one-hot-encoding vector of length 5000 to a “meaning-vector” (=embedding) of the lower dimensionality of 32. Then we add a layer_simple_rnn on top, and finally a layer_dense for the binary prediction of review sentiment.

```
model_keras2 <- keras_model_sequential()

model_rnn <- model_keras2 %>%
  layer_embedding(input_dim = 5000, output_dim = 32) %>%
  layer_simple_rnn(units = 32, activation = "tanh") %>%
  layer_dense(units = 1, activation = "sigmoid")
```

Here the structure of the model can be seen

```
summary(model_rnn)
```

```
## Model: "sequential_2"
## -----
## Layer (type)                Output Shape          Param #
## -----
## embedding (Embedding)       (None, None, 32)      160000
## -----
## simple_rnn (SimpleRNN)      (None, 32)            2080
## -----
## dense_6 (Dense)             (None, 1)              33
## -----
## Total params: 162,113
## Trainable params: 162,113
## Non-trainable params: 0
## -----
```

Again we use a basic setup for binary prediction.

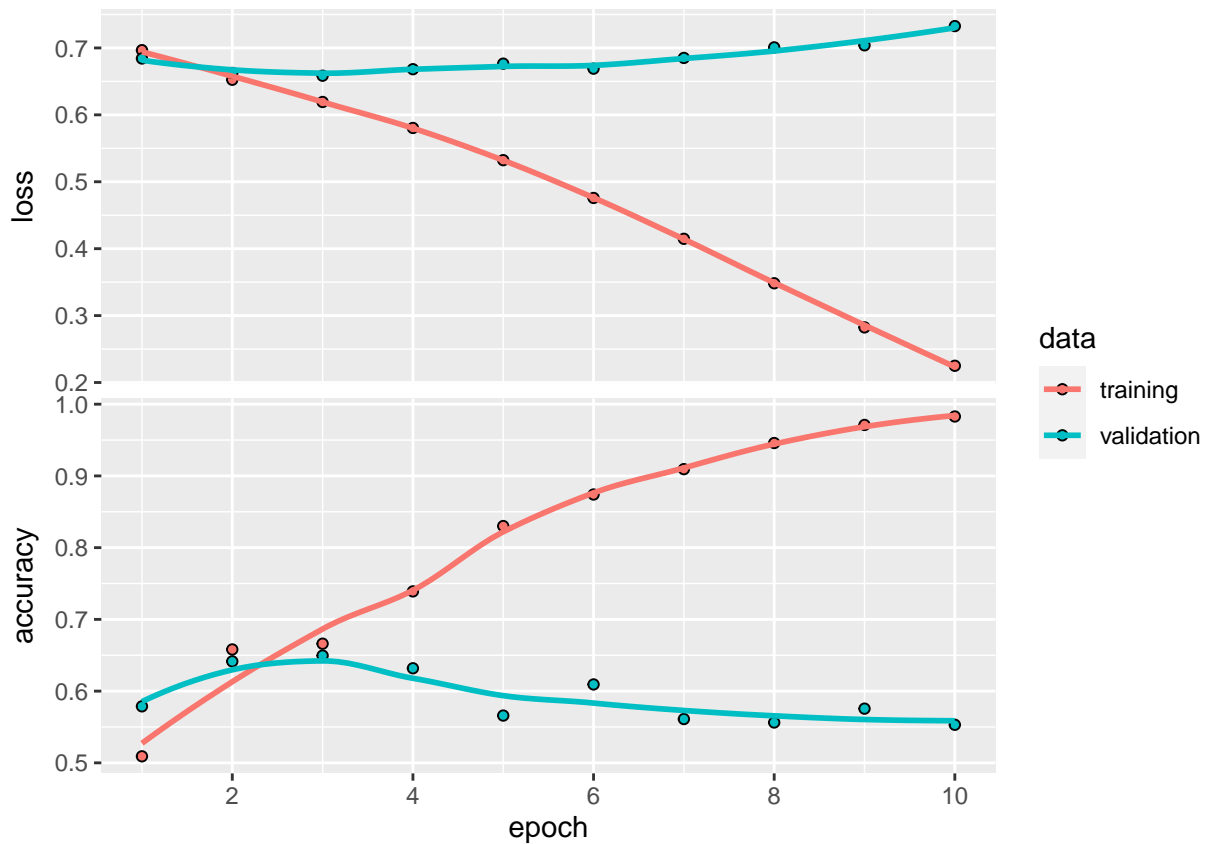
```
model_rnn %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy",
  metrics = "accuracy"
)
```

And run our model

```
set.seed(12345)
history_rnn <- model_rnn %>% fit(
  x_train_pad, y_train_data,
  epochs = 10,
  batch_size = 516,
  validation_split = 0.25
)
```

```
plot(history_rnn)
```

```
## 'geom_smooth()' using formula 'y ~ x'
```



Again the training set outperforms the validation set a lot, which shows our model is overfitted. Further we see that the loss of the validation set starts to climb after a couple of epochs and the accuracy to fall, so not a good model.

```
metrics3 = model_rnn %>% evaluate(x_test_pad, y_test_data); metrics3
```

```
##      loss  accuracy
## 0.7871893 0.5144578
```

Running our model on the test data shows an accuracy of 52%, but we will now try to fine tune it to make it better.

Tunning the model This time we again try to reduce the number of parameters to prevent over fitting. We also make another rnn layer and drop a fraction of the units with a drop_out input. Return_sequences = TRUE return the full state sequence of the first rnn layer so next rnn layer gets the full sequence of the input.xMethl

```
model_keras2 <- keras_model_sequential()

model_rnn2 <- model_keras2 %>%
  layer_embedding(input_dim = 5000, output_dim = 16) %>%
  layer_simple_rnn(units = 16, return_sequences = TRUE, activation = "tanh", recurrent_dropout=0.1) %>%
  layer_simple_rnn(units = 16, return_sequences = FALSE, activation = "tanh", recurrent_dropout=0.1) %>%
  layer_dense(units = 1, activation = "sigmoid")
```

Here the structure of the model can be seen

```
summary(model_rnn2)
```

```
## Model: "sequential_3"
## -----
## Layer (type)                Output Shape          Param #
## =====
## embedding_1 (Embedding)      (None, None, 16)      80000
## -----
## simple_rnn_2 (SimpleRNN)     (None, None, 16)      528
## -----
## simple_rnn_1 (SimpleRNN)     (None, 16)            528
## -----
## dense_7 (Dense)             (None, 1)             17
## =====
## Total params: 81,073
## Trainable params: 81,073
## Non-trainable params: 0
## -----
```

Again we use a basic setup for binary prediction.

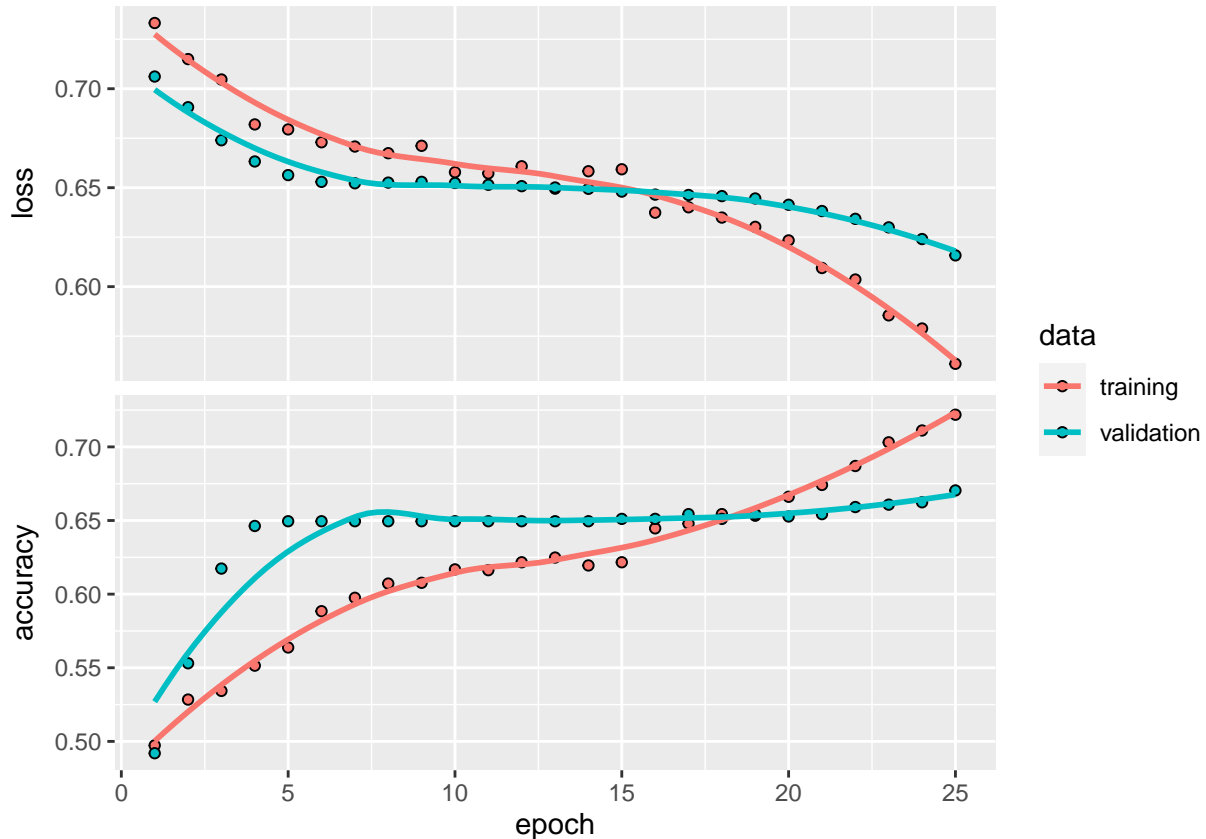
```
model_rnn2 %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy",
  metrics = "accuracy"
)
```

And run our model

```
set.seed(12345)
history_rnn2 <- model_rnn2 %>% fit(
  x_train_pad, y_train_data,
  epochs = 25,
  batch_size = 516,
  validation_split = 0.25
)
```

```
plot(history_rnn2)
```

```
## 'geom_smooth()' using formula 'y ~ x'
```



The loss of the training and validation data seems to follow each other quite well for the first couple of epochs. After that the validation set begins to vary a lot, flying up and down. The accuracy also follows each other a lot, but after around 12 epochs they cross and the training data runs off.

```
metrics4 = model_rnn2 %>% evaluate(x_test_pad, y_test_data); metrics4
```

```
##      loss  accuracy
## 0.7168422 0.5698795
```

Running our model on the test data shows an accuracy of 58%, this is better than the baseline model but not at all good.

LSTM

We will now try our data on a LSTM model, but here we are only running one model, since it takes a very long time to run it. We start by using an embedding layer and then we go to a LSTM layer with a unit size of our padded data, which has the size of 300 due to it being the 300 first words in every song. In the lstm layer, we freeze some of the input weights and also some of the state weights and since we only use one layer we set `sequence = false`, so it just compiles the input to a single output.

```
model_lstm <- keras_model_sequential() %>%
  layer_embedding(input_dim = 5000, output_dim = 32) %>%
  layer_lstm(units = 300, dropout = 0.25, recurrent_dropout = 0.25, return_sequences = FALSE) %>%
  layer_dense(units = 1, activation = "sigmoid")
```

We use a base compiling

```
model_lstm %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy",
  metrics = c("acc")
)
```

The model has a lot of parameters which makes it time consuming to run it.

```
summary(model_lstm)
```

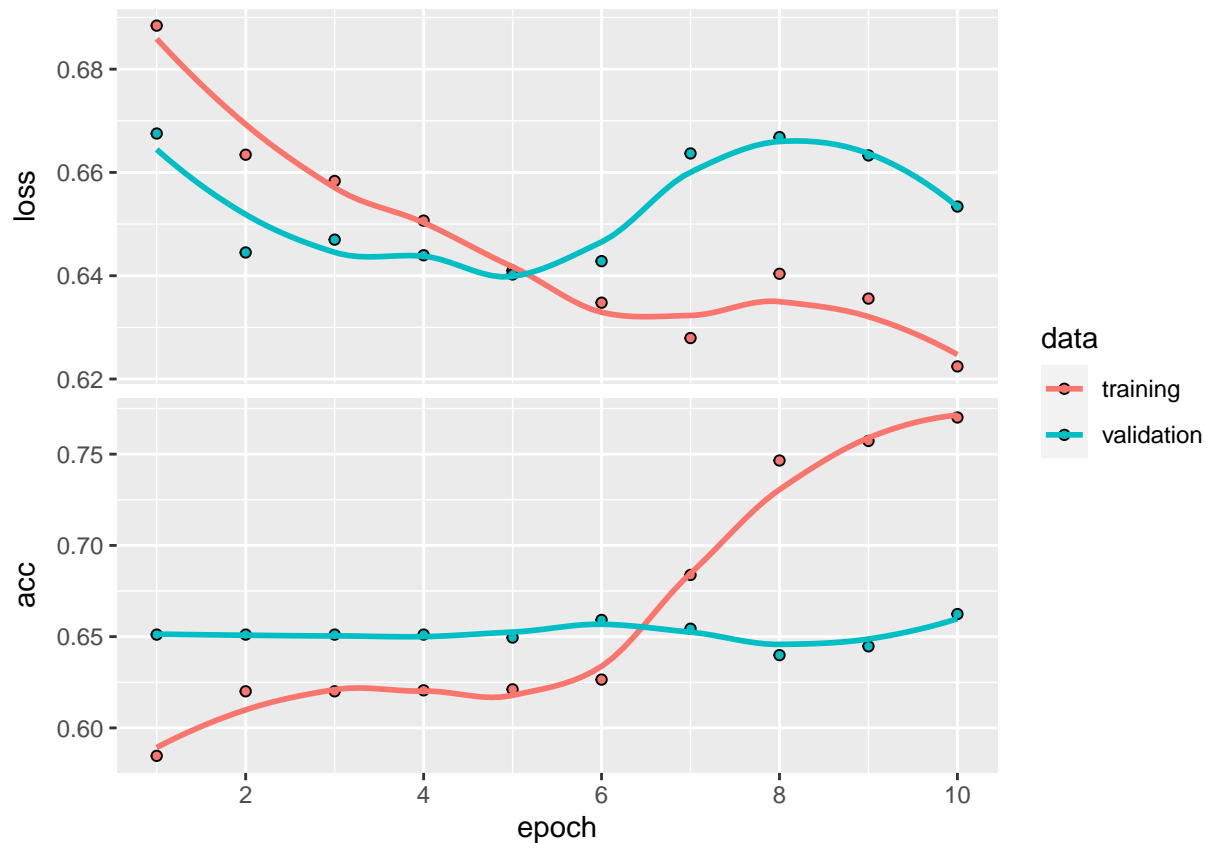
```
## Model: "sequential_4"
## -----
## Layer (type)                Output Shape          Param #
## =====
## embedding_2 (Embedding)      (None, None, 32)      160000
## -----
## lstm (LSTM)                  (None, 300)           399600
## -----
## dense_8 (Dense)              (None, 1)              301
## =====
## Total params: 559,901
## Trainable params: 559,901
## Non-trainable params: 0
## -----
```

Then we run the mode

```
history_lstm <- model_lstm %>% fit(
  x_train_pad, y_train_data,
  epochs = 10,
  batch_size = 512,
  validation_split = 0.25
)
```

```
plot(history_lstm)
```

```
## 'geom_smooth()' using formula 'y ~ x'
```



The model doesn't seem to perform any better than the previous ones, but the training and validation data did follow each other quite well for some epochs, but then at the end they split up due to the increasing loss value of the validation data.

```
metrics5 = model_lstm %>% evaluate(x_test_pad, y_test_data); metrics5
```

```
##      loss      acc
## 0.6860912 0.5445783
```

A really poor result to say the least with an accuracy of 47%