

topic_mod

Group 10: Bill Gao, Jiun Lee, Priam Vyas, Danya Zhang

2022-11-16

##Import Data

```
imdb <- read_csv ("/Users/jiunlee/MSSP22/Fidelity/IMDB Dataset.csv",  
                  col_names = TRUE,  
                  show_col_types = FALSE)  
##remove sentiment column since we don't use it.  
imdb <- imdb %>% select(-sentiment)  
  
##remove duplicates of review  
imdb <- unique(imdb)  
  
##sample down for 100reviews.  
set.seed(456)  
review_index <- 1:dim(imdb)[1]  
text_df <- cbind(review_index,imdb)  
text_df <- text_df %>% slice_sample(n = 100, replace = FALSE)  
rm(review_index)
```

Cleaning: Tokenizing, Removing stopwords,tf-idf

##The imdb dataset has a lot of stopwords and meaningless words. #We will remove stopwords and words unnecessary for guessing movie genre.

```
library(stringr)  
library(tidytext)  
## tokenizing, count the number of words within each review.  
token <- text_df %>%  
  unnest_tokens(word, review) %>%  
  count(review_index, word, sort=TRUE) %>%  
  rename(count=n)  
  
##There's a lot of stopwords. Let's remove them.  
  
## create a stop word vector  
stop <- unlist(stop_words[,1])  
## drop the attribute  
stop <- StripAttr(stop)  
##restore tokens dataset to check their  
check <- token  
## check words against stop word lists
```

```

remove <- check$word %in% stop
## to make it easier to see create a data frame
d <- cbind(token,remove)
## create an index of words(not stopwords)
f <- which(d$remove == FALSE)
##clean tokens that has no stopwords
clean_token <- d %>% slice(f) %>% select(-remove)

##Let's subset the data frame that has only meaningful words

##vector that has meaningless words
strings <- c("br","movie","film", "scene", "character","story","bit","lot","bad","act","hard","awful",")
##detect numbers of rows that has meaningless words
meaningless <- str_detect(clean_token$word, paste(strings, collapse = "|"))
##detect numbers of meaningful words
has_meaning <- which(meaningless==F)
##subset: tokens without meaningless
clean_token <- clean_token %>% slice(has_meaning)

##remove redundant datas and values
rm(d,check,f,meaningless,has_meaning,strings,stop,remove,token)

```

```

#create lda model
library(topicmodels)

#convert sample token tibble to document term matrix for lda
clean_token_dmat <- clean_token %>%
  cast_dtm(review_index, word, count)

#select k=6 because 6 general film genres
imdb_lda <- LDA(clean_token_dmat, k = 6, control = list(seed = 1234))
imdb_topics <- tidy(imdb_lda, matrix = "beta")

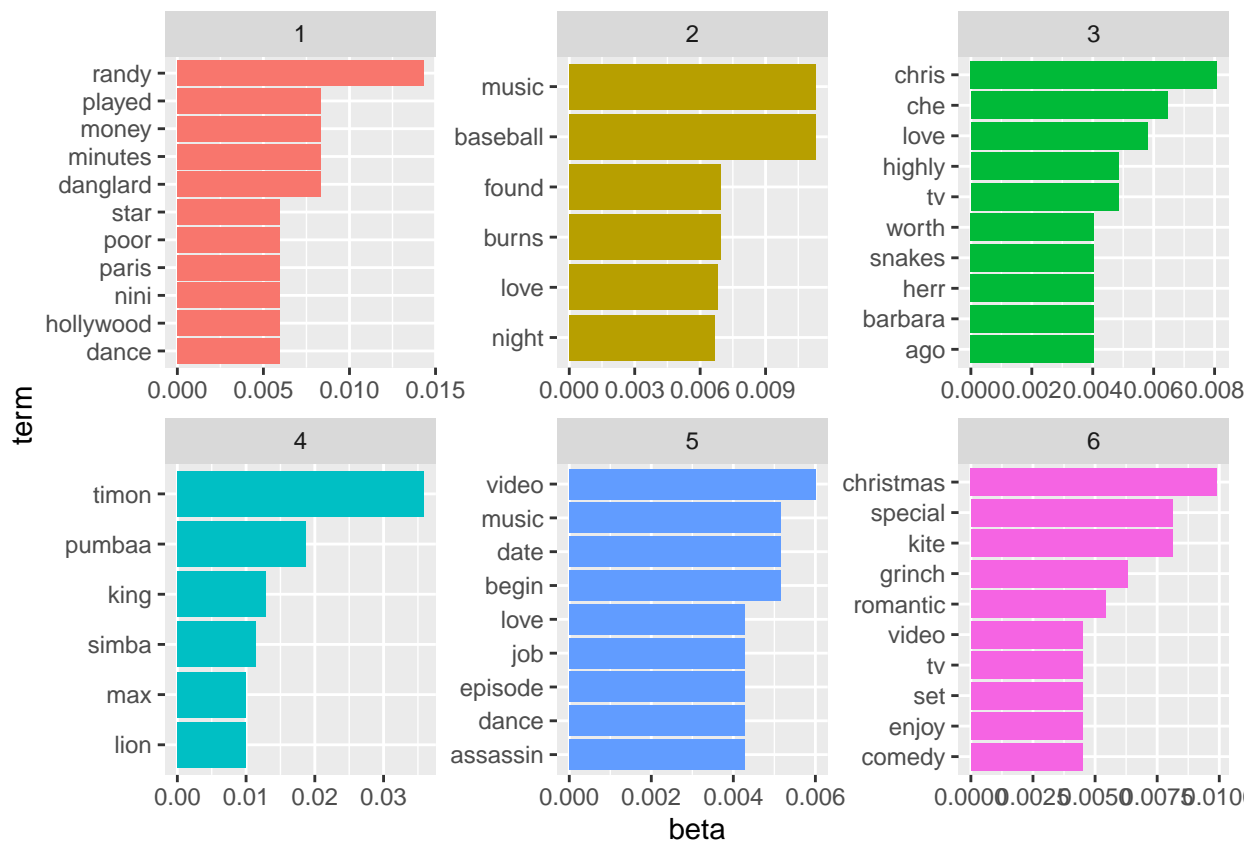
```

```

library(ggplot2)
imdb_top_terms <- imdb_topics %>%
  group_by(topic) %>%
  slice_max(beta, n = 6) %>%
  ungroup() %>%
  arrange(topic, -beta)

imdb_top_terms %>%
  mutate(term = reorder_within(term, beta, topic)) %>%
  ggplot(aes(beta, term, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  scale_y_reordered()

```



The plot above shows the top 6 words for each topic in our LDA. We've split up our LDA into 6 genres, which represents the number of topics we have. The topic 4 looks like a famous Disney movie 'Lion King'.

```
#gamma: per-document-per-topic probabilities
imdb_documents <- tidy(imdb_lda, matrix = "gamma")
imdb_documents
```

```
## # A tibble: 600 x 3
##   document topic    gamma
##   <chr>    <int>    <dbl>
## 1 35484      1 0.0000498
## 2 29214      1 0.999
## 3 12301      1 0.000217
## 4 320        1 0.000112
## 5 19086      1 0.0000577
## 6 8473       1 0.000135
## 7 21531      1 0.000137
## 8 5044       1 0.0000968
## 9 37209      1 1.00
## 10 6706      1 0.000117
## # ... with 590 more rows
## # i Use 'print(n = ...)' to see more rows
```

```
#most common words in document
tidy(clean_token_dmat) %>%
  filter(document == 6) %>%
  arrange(desc(count))
```

```
## # A tibble: 0 x 3
## # ... with 3 variables: document <chr>, term <chr>, count <dbl>
## # i Use 'colnames()' to see all variable names
```

```
assignments <- augment(imdb_lda, data = clean_token_dmat)
assignments
```

```
## # A tibble: 5,286 x 4
##   document term      count .topic
##   <chr>    <chr>    <dbl> <dbl>
## 1 35484    timon        25     4
## 2 35484    pumbaa       13     4
## 3 29214    randy        12     1
## 4 12301    christmas    11     6
## 5 320      baseball     10     2
## 6 38255    baseball      3     2
## 7 19086    kite         9      6
## 8 320      burns        8     2
## 9 8473     chris         8     3
## 10 14281   chris         2     3
## # ... with 5,276 more rows
## # i Use 'print(n = ...)' to see more rows
```

The assignments tibble above count up the words for each topic.

##Let's look tf-idf to see what is the most important words in the whole reviews. ##tf-idf

```
review_tf_idf <- clean_token %>%
  bind_tf_idf(review_index, word, count)

##Look at terms with high tf-idf in reviews.
review_tf_idf <- review_tf_idf %>%
  arrange(desc(tf_idf))
##It looks like the high tf-idf's tf are mostly 1.
##For words that tf=1, it means those words are only contained on one review, and the tf-idf algorithm

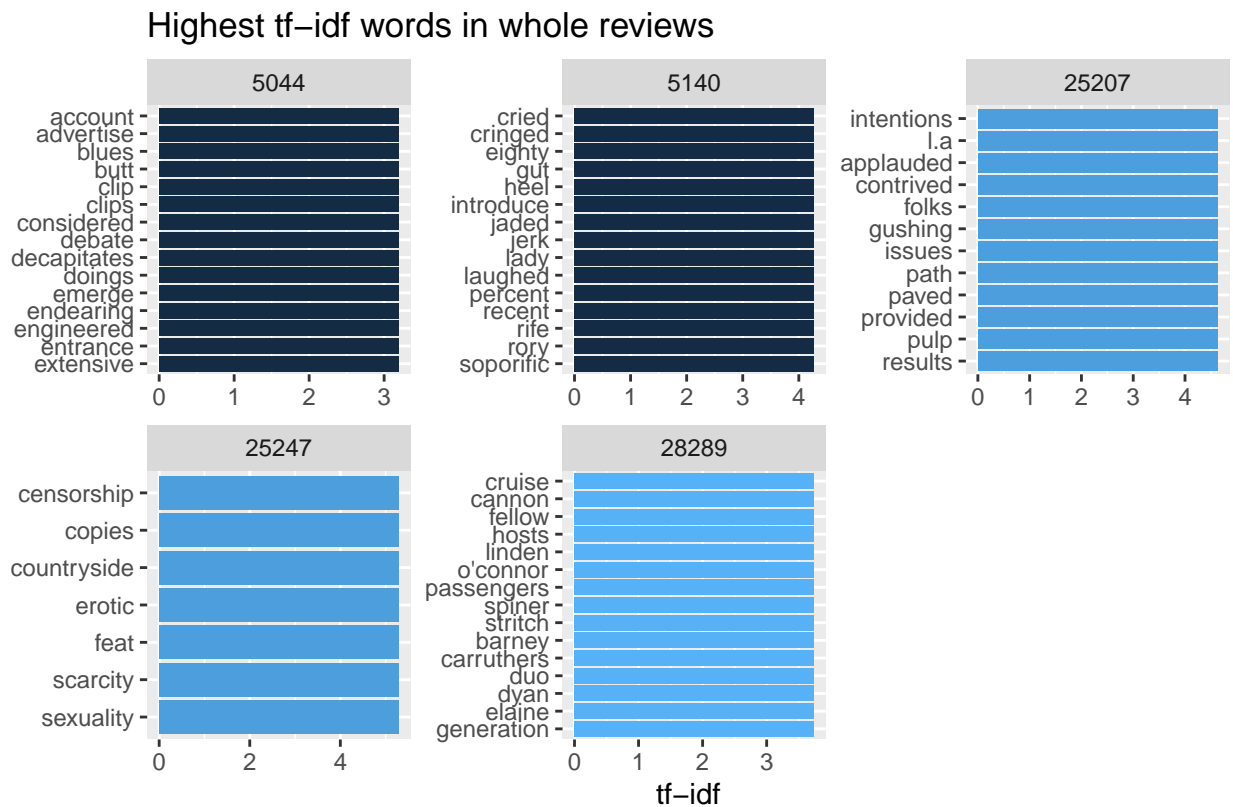
##So, remove all tf != 1.
tf_1 <- which(review_tf_idf$tf==1)
tf_idf_high <- review_tf_idf %>% slice(tf_1) %>% select(-count,-tf,-idf) #remove column 'count','tf','idf'

rm(review_tf_idf)

##review_index numbers in tf_idf_high
index <- unique(tf_idf_high$review_index)

##tf-idf plot
##Let's make the plots with only 6 review_index.
tf_idf_high %>%
  filter(review_index %in% c(25207,5044,5140,28289,25247)) %>%
  arrange(desc(tf_idf)) %>%
  group_by(review_index) %>%
```

```
distinct(word,review_index, .keep_all = TRUE) %>%
slice_max(tf_idf, n = 15, with_ties = FALSE) %>%
ungroup() %>%
mutate(word = factor(word, levels = rev(unique(word)))) %>%
ggplot(aes(tf_idf, word, fill = review_index)) +
geom_col(show.legend = FALSE) +
facet_wrap(~review_index, ncol = 3, scales = "free") +
labs(title = "Highest tf-idf words in whole reviews",
caption = "IMDB Dataset",
x = "tf-idf", y = NULL)
```



IMDB Dataset

- On each 6 plot, we can see top 15 words with high tf-idf.
- Among them, we can verify some meaningful words for checking their genres.
- For example, in review'25247', the words 'censorship','erotic','sexuality' imply that the review is about romance movie.