

# Topic Modeling

Gary Wang

2024-11-01

```
library(lexicon)
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.4      v readr      2.1.4
## v forcats    1.0.0      v stringr   1.5.0
## v ggplot2    3.5.1      v tibble    3.2.1
## v lubridate  1.9.2      v tidyr     1.3.1
## v purrr      1.0.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```
library(topicmodels)
library(tidytext)
library(factoextra)
```

```
## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

```
library(dplyr)
library(ggplot2)
library(wordcloud)
```

```
## Loading required package: RColorBrewer
```

```
library(RColorBrewer)
```

```
movies <- read.csv("movie_plots.csv")
```

Group by Genre, summarize common words, and find genre frequency

```
# Extract genres based on common keywords found in the plot descriptions
movies <- movies %>%
  mutate(Genre = case_when(
    str_detect(Plot, "(?i)science|space|experiment|future|alien|
                    robot|technology|planet|human|new world|earth") ~ "Sci-Fi",
    str_detect(Plot, "(?i)love|romance|relationship|affair|
```

```

        wedding|couple|meets|girl") ~ "Romance",
str_detect(Plot, "(?i)war|battle|army|soldier|conflict|
    military|enemy|tank|battlefield|dead") ~ "War",
str_detect(Plot, "(?i)ghost|haunt|horror|fear|terror|
    scary|supernatural|creepy") ~ "Horror",
str_detect(Plot, "(?i)crime|detective|murder|investigate|
    thriller|mafia|heist|mystery") ~ "Crime",
str_detect(Plot, "(?i)action|fight|fighting|adventure|hero|explosion|
    battle|rescue") ~ "Action",
str_detect(Plot, "(?i)comedy|funny|humor|laugh|joke|
    satire|parody") ~ "Comedy",
str_detect(Plot, "(?i)history|historical|biography|true story|
    period drama|century|ancient") ~ "History",
str_detect(Plot, "(?i)fantasy|magic|myth|legend|superhero|
    kingdom|evil") ~ "Fantasy",
str_detect(Plot, "(?i)western|cowboy|wild west|sheriff|ranch|
    town|outlaw") ~ "Western",
str_detect(Plot, "(?i)documentary|docu|true events|reality|
    biopic") ~ "Documentary",
str_detect(Plot, "(?i)sport|game|team|match|championship|
    wrestling") ~ "Sport",
str_detect(Plot, "(?i)home|people|brother|daughter|brothers|
    friend|wife|son|father|mother") ~ "Family",
TRUE ~ "Other"
))

# Calculate the frequency of each genre
genre_frequency <- movies %>%
  count(Genre, name = "Frequency")

# Tokenize the plots and remove stop words
plot_words <- movies %>%
  unnest_tokens(word, Plot) %>%
  anti_join(get_stopwords()) %>%
  count(Genre, word, sort = TRUE)

```

## Joining with 'by = join\_by(word)'

```

# Group by Genre, summarize common words, and find genre frequency
nested_data <- plot_words %>%
  group_by(Genre) %>%
  summarize(
    Words = paste(unique(word), collapse = ", ")
  ) %>%
  left_join(genre_frequency, by = "Genre")

view(nested_data)

```

Create a Document Term Matrix

```

dtm <- plot_words %>%
  cast_dtm(Genre, word, n)

```

```
dtm
```

```
## <<DocumentTermMatrix (documents: 14, terms: 14842)>>
## Non-/sparse entries: 30812/176976
## Sparsity          : 85%
## Maximal term length: 17
## Weighting          : term frequency (tf)
```

Fir the LDA Model

```
# Since we have 14 genres, we try to set k close to the number of genres.
k <- 20
```

```
# Fit the LDA model with k = 14
lda_model <- LDA(dtm, k = k, control = list(seed = 999))
```

```
# Extract the topic-term matrix
topics <- tidy(lda_model, matrix = "beta")
```

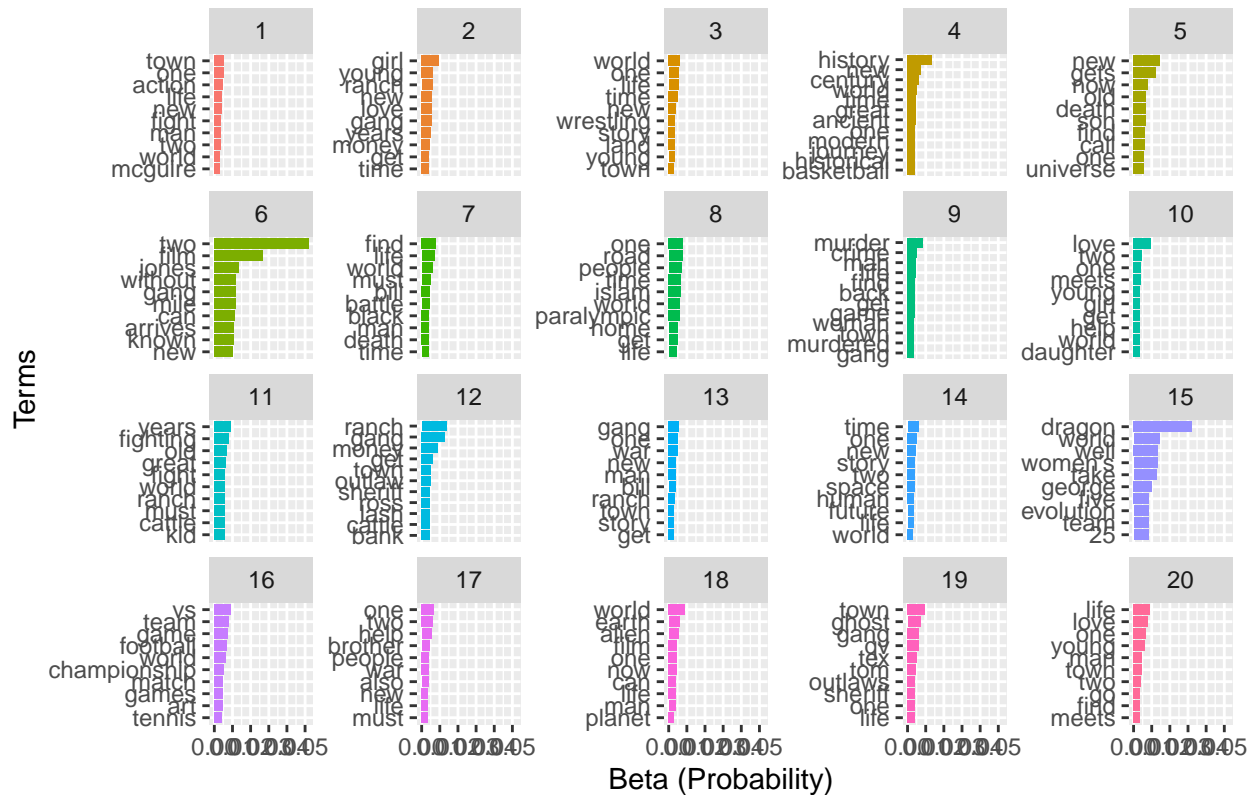
```
# View the top 10 terms for each topic
top_terms <- topics %>%
  group_by(topic) %>%
  slice_max(beta, n = 10) %>%
  ungroup() %>%
  arrange(topic, -beta)
```

```
print(top_terms)
```

```
## # A tibble: 205 x 3
##   topic term      beta
##   <int> <chr>    <dbl>
## 1     1 town    0.00552
## 2     1 one     0.00545
## 3     1 action 0.00474
## 4     1 life   0.00427
## 5     1 new    0.00387
## 6     1 fight  0.00378
## 7     1 man    0.00347
## 8     1 two    0.00339
## 9     1 world  0.00300
## 10    1 mcguire 0.00296
## # i 195 more rows
```

```
# Plot the top terms for visualization
ggplot(top_terms, aes(x = reorder_within(term, beta, topic), y = beta,
  fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free_y") +
  coord_flip() +
  scale_x_reordered() +
  labs(title = "Top 10 Terms for Each Topic", x = "Terms",
    y = "Beta (Probability)")
```

## Top 10 Terms for Each Topic



## Cluster Plot and Gamma Plot

```
# Filter the topic-term matrix to include only the top 50 most frequent terms
filtered_topics <- topics %>%
  group_by(term) %>%
  summarize(total_beta = sum(beta)) %>%
  top_n(50, total_beta)

topic_term_matrix <- topics %>%
  filter(term %in% filtered_topics$term) %>%
  spread(term, beta) %>%
  column_to_rownames(var = "topic")

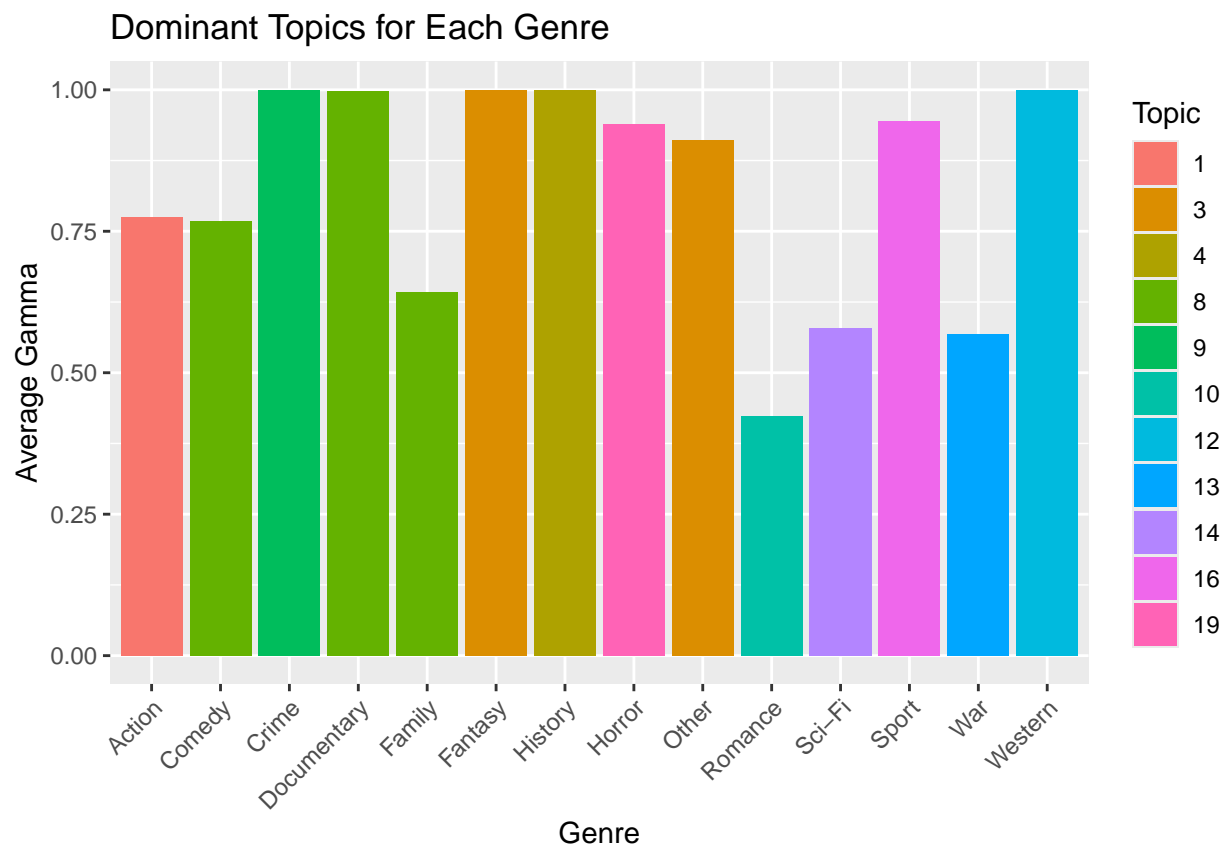
# Perform PCA on the filtered topic-term matrix
pca_result <- prcomp(topic_term_matrix, scale. = TRUE)

# Plot the PCA result
fviz_pca_biplot(pca_result,
  repel = TRUE,
  col.var = "blue",
  col.ind = "red",
  pointsize = 3,
  alpha.ind = 0.6,
  title = "Cluster Plot of Topics")
```

```
## # A tibble: 14 x 3
##   Genre      topic avg_gamma
##   <chr>    <int>    <dbl>
## 1 Action      1    0.775
## 2 Comedy      8    0.769
## 3 Crime        9    1.00
## 4 Documentary  8    0.998
## 5 Family       8    0.643
```

```
## 6 Fantasy      3      1.00
## 7 History      4      1.00
## 8 Horror      19     0.939
## 9 Other        3     0.910
## 10 Romance    10     0.423
## 11 Sci-Fi     14     0.579
## 12 Sport      16     0.945
## 13 War        13     0.567
## 14 Western    12     1.00
```

```
ggplot(dominant_topics, aes(x = Genre, y = avg_gamma, fill = factor(topic))) +
  geom_col() +
  labs(title = "Dominant Topics for Each Genre",
       x = "Genre",
       y = "Average Gamma",
       fill = "Topic") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



Create a Word Cloud for each genre

```
# Set up color palette
palette <- brewer.pal(8, "Dark2")

# Create word clouds for each topic in a loop
k <- 14
for (i in 1:k) {
```

}



become brother daughter come york  
dead red life money knows  
work son old love set father outlaw  
like never three take sister also new ranch back even city  
men gun make help bill goes young girl years family soon  
use billy goes secret land takes town gang get falls arrives  
gold chief war killed west wife joins buck wild boy  
friend later try man





nypa g temple greatest women freedom  
tesla early part ancient era john letters  
college end journey great first many including  
following: irish jetsteam mentrue york trot wind  
long modern century europe dancers  
nemo just power story history evil stories  
see man big one new world c.s. winning  
video two life like time day basketball region river  
known film travel like historical dreams  
20th led kentucky dances return viewers  
breathhtaking

return episode christmas violent energy  
group people machine brothers  
series takes every two kill go mission sides  
like old death storm year  
son re now last religion  
maverick telling bring johnny get new call slaves time  
wife jack life map gets find father men  
many prison one man lies united  
players own self universe brings ody night



indians daughter across historical end  
lead forced takes death must goes behind way dan  
team black world film lost go  
dance man find story learns  
get can time save tom  
dead john life real events serial  
secret west new kill bill war  
many bob like battle two land army fight  
north old series city place

nhl team bob u.s heart goes  
trip back 2016 bill yea islam power  
field like rio people long  
world done order life religion  
son well road get brother  
look new time hz home exercises  
york city lost treasure

jan jack cattle young town come sam arrive murders holt  
night gold friend ranch find now men gets  
city way make burke becomes man life game card oneqv  
frank baby fight crime get mary just new  
run like story back time killer garrett  
takes two woman brother  
help murdered film detective

trially makes comes stiles end  
forced secret also like never  
way local killed gang world horse dead  
friend last time help get film meet  
falls another gold finds one takes family  
named leads take find john war jake york  
must leads take find john war jake york  
dr heart woman love father movie  
daughter two girl now go wants  
can young life man couple return  
outlaw comes soon story wife just real couple see  
line begins action high new friends



films like keep ex fight cattle son cr  
justice great way brothers  
back kid fight old must ned  
later team kill  
help billy years name  
action fighting young taken  
time can world young new men  
gang mine can world father joe find  
return west s



many way land men killed also toge  
streets learn ranch bill world order  
civil two story one town friend  
long john time gang must gets jim  
kill old sheriff war king get gold son  
wild line find first finds now new money ack title  
deadly can . . . night

