

# Topic Modeling

Haoran Cui

## Packages

```
library(tidyverse)
```

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
v ggplot2    3.5.1      v tibble     3.2.1
v lubridate  1.9.3      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
```

```
library(tm)
```

Warning: package 'tm' was built under R version 4.4.2

Loading required package: NLP

Attaching package: 'NLP'

The following object is masked from 'package:ggplot2':

annotate

```
library(topicmodels)
```

Warning: package 'topicmodels' was built under R version 4.4.2

```
library(ldatuning)
```

Warning: package 'ldatuning' was built under R version 4.4.2

```
library(tidytext)
```

Warning: package 'tidytext' was built under R version 4.4.2

```
library(Rtsne)
```

Warning: package 'Rtsne' was built under R version 4.4.2

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

Warning: package 'wordcloud' was built under R version 4.4.2

Loading required package: RColorBrewer

```
library(RColorBrewer)
```

## Data Cleaning

```
movie_data = read.csv("movie_plots.csv")  
view(movie_data)
```

```
plots_by_word <- movie_data %>% unnest_tokens(word, Plot)  
plot_word_counts <- plots_by_word %>%  
  anti_join(stop_words, by = join_by(word)) %>%  
  count("Movie Name", word, sort = TRUE)
```

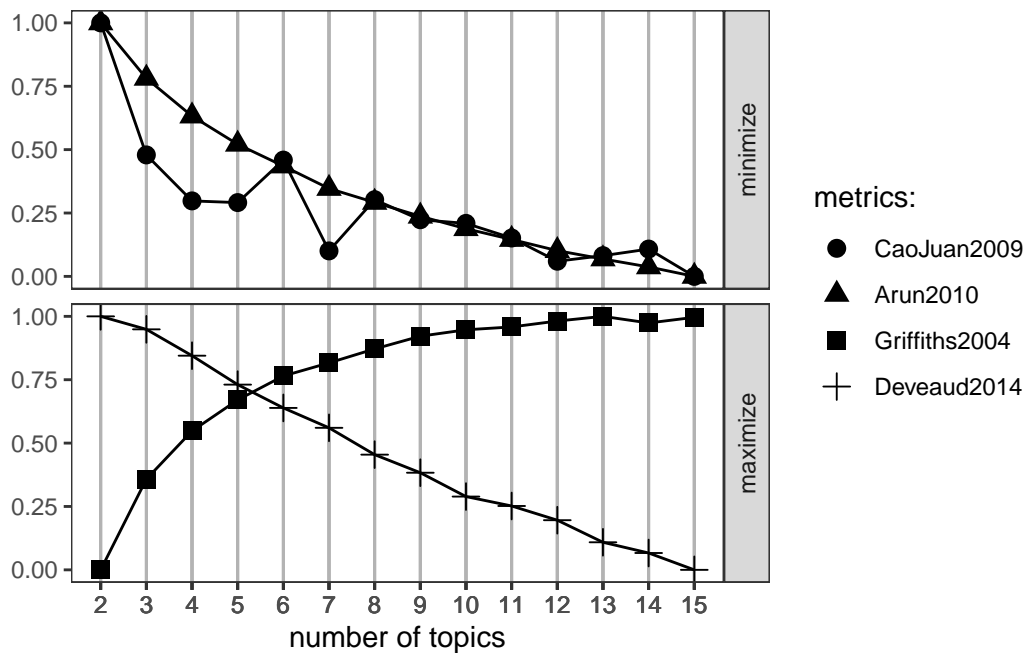


```
FindTopicsNumber_plot(result)
```

Warning: The ``scale`` argument of ``guides()`` cannot be ``FALSE``. Use "none" instead as of ggplot2 3.3.4.

i The deprecated feature was likely used in the ldatuning package.

Please report the issue at <https://github.com/nikita-moor/ldatuning/issues>.



- **Setting Up Topic Modeling Parameters:** Uses `FindTopicsNumber` to test different topic numbers (from 2 to 15) on the document-term matrix (DTM) with specified evaluation metrics.
- **Defining the Method and Controls:** Sets the topic modeling method to “Gibbs” sampling, uses a random seed for reproducibility, and specifies the number of processing cores.
- **Visualizing Optimal Topics:** Calls `FindTopicsNumber_plot(result)` to plot the results and help identify the optimal number of topics based on the evaluation metrics.

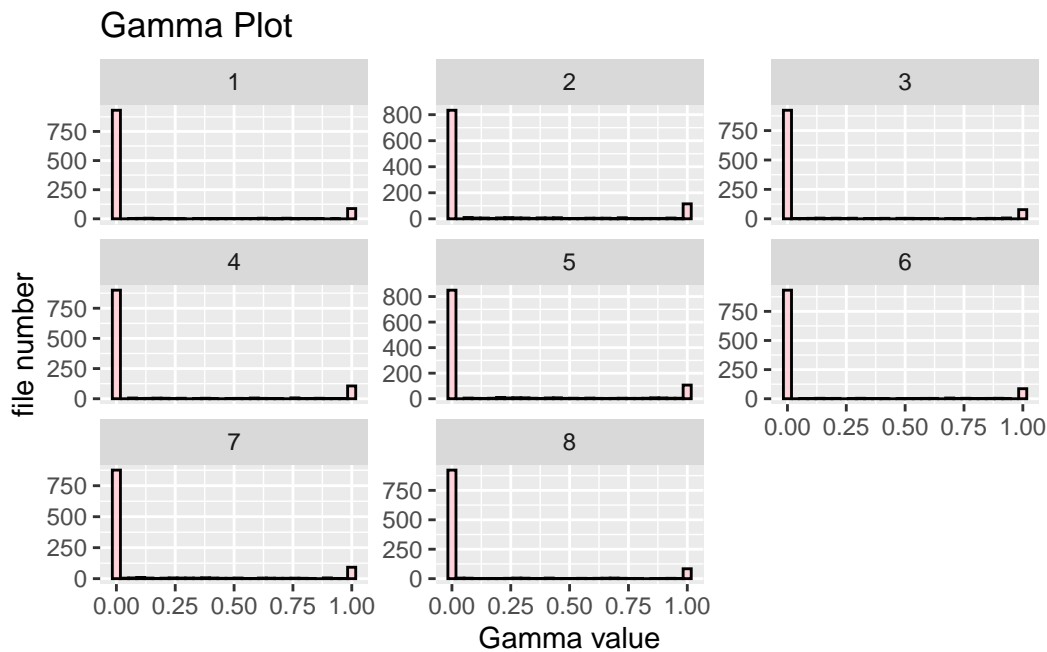
```
#set k to 8
k <- 8
lda_model <- LDA(dtm, k = k, control = list(seed = 321))
```



cluster by topic, illustrating the distinctiveness and relationships between topics by showing how closely documents with similar themes group together.

```
document_topics <- tidy(lda_model, matrix = "gamma")

# Gamma plot
ggplot(document_topics, aes(x = gamma)) +
  geom_histogram(bins = 30, fill = "pink", color = "black", alpha = 0.7) +
  facet_wrap(~ topic, scales = "free_y") +
  labs(title = "Gamma Plot", x = "Gamma value", y = "file number")
```



The Gamma Plot visualizes the distribution of topic probabilities (gamma values) for each document across the eight topics in the LDA model. Each subplot shows the gamma values for a specific topic, revealing that most documents have low gamma values for most topics but a high gamma value for one topic, indicating strong association with a single dominant topic. This pattern confirms that the model effectively assigns documents to distinct topics, as expected in well-separated topic modeling, where each document is primarily linked to one topic rather than spread across multiple topics.

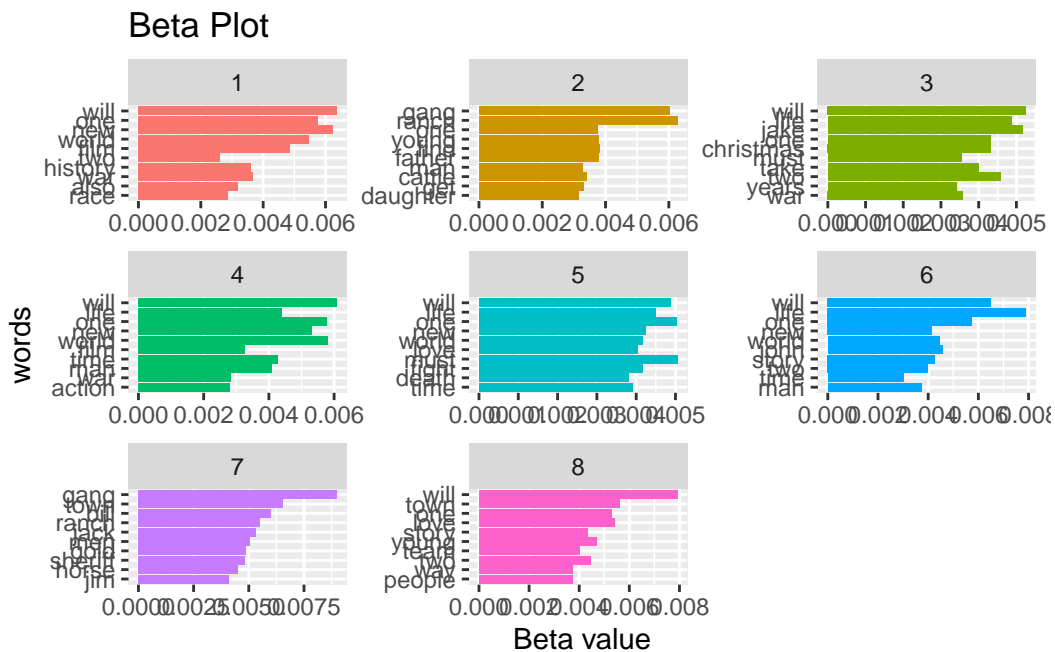
```
topic_terms <- tidy(lda_model, matrix = "beta")
top_terms <- topic_terms %>%
  group_by(topic) %>%
```

```

slice_max(beta, n = 10) %>%
ungroup() %>%
arrange(topic, -beta)

ggplot(top_terms, aes(x = reorder(term, beta), y = beta, fill = factor(topic))) +
  geom_bar(stat = "identity", show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  coord_flip() +
  labs(title = "Beta Plot", x = "words", y = "Beta value")

```



The Beta Plot displays the top 10 terms for each of the eight topics in the LDA model, with each term's importance measured by its beta value. Beta values represent the probability of a word being associated with a particular topic, so higher beta values indicate words that are more representative of the topic. Each subplot corresponds to one topic, showing the most significant words in descending order of their beta values. This plot provides insight into the defining terms of each topic, making it easier to interpret and label the topics based on the prominent words associated with them. This visualization helps understand the thematic structure of each topic in the model.

```

library(RColorBrewer)

all_terms <- top_terms %>%

```

