

Text Mining and Word Cloud Analysis of MLK's 'I Have a Dream' Speech

DIVYA CHENTHAMARAKSHAN

May 5, 2025

Contents

| | |
|--|-----------|
| 1. Introduction | 2 |
| 2. Setting Up the Environment | 2 |
| 2.1 Required Packages | 2 |
| 2.2 Loading the Data | 3 |
| 3. Creating the Corpus | 3 |
| 3.1 Understanding Corpus in Text Mining | 3 |
| 4. Text Preprocessing | 4 |
| 4.1 Custom Transformations | 4 |
| 4.2 Standard Text Cleaning | 5 |
| 5. Creating the Document-Term Matrix | 6 |
| 5.1 Understanding the Document-Term Matrix | 6 |
| 5.2 Exploring Word Frequencies | 7 |
| 5.3 Distribution of Word Lengths | 7 |
| 6. Word Cloud Visualization | 8 |
| 6.1 Creating a Word Cloud | 8 |
| 6.2 Finding Frequent Terms | 10 |
| 6.3 Finding Word Associations | 10 |
| 6.4 Hierarchical Clustering of Terms | 10 |
| 7. Analysis of Results | 11 |
| 7.1 Word Frequency Analysis | 11 |
| 7.2 Word Cloud Interpretation | 12 |
| 7.3 Word Associations | 12 |
| 7.4 Word Occurrence Timeline | 13 |

| | |
|---|-----------|
| 8. Interpretation and Recommendations | 14 |
| 8.1 Interpretation of Findings | 14 |
| 8.2 Recommendations for Further Analysis | 14 |
| 8.3 Suggested Additional Variables | 15 |
| 8.4 Comparing Word Usage with Contemporary Speeches | 15 |
| 9. Conclusion | 17 |

1. Introduction

This report presents a text mining analysis of Martin Luther King Jr.'s iconic "I Have a Dream" speech using R programming language. Text mining is the process of extracting meaningful information from unstructured text data. It enables us to discover patterns, trends, and insights that might not be immediately apparent through conventional reading.

The purpose of this analysis is to:

- - Demonstrate the text mining process from data preparation to visualization
- - Analyze the frequency and relationships between words in the speech
- - Interpret the results to understand the key themes and messages

2. Setting Up the Environment

2.1 Required Packages

First, we need to install and load the necessary R packages for text mining and visualization:

```
# Install packages
install.packages("tm")           # for text mining
install.packages("SnowballC")   # for text stemming
install.packages("wordcloud")   # word-cloud generator
install.packages("RColorBrewer") # color palettes
```

```
# Load libraries
library("tm")
library("SnowballC")
library("wordcloud")
library("RColorBrewer")
```

The packages serve the following purposes:

- - **tm**: The Text Mining package provides functions for text manipulation and transformation
- - **SnowballC**: Implements the Porter stemming algorithm to reduce words to their root form
- - **wordcloud**: Creates visual representations of word frequency
- - **RColorBrewer**: Provides color palettes for data visualization

2.2 Loading the Data

We'll load Martin Luther King Jr.'s "I Have a Dream" speech from a text file:

```
# Define the file path to the speech text
filePath <- "http://www.sthda.com/sthda/RDoc/example-files/martin-luther-king-i-have-a-dream-speech.txt"

# Read the text file line by line
text <- readLines(filePath)

# Print the first few lines to verify the data was loaded correctly
head(text)
```

```
## [1] ""
## [2] "And so even though we face the difficulties of today and tomorrow, I still have a dream. It is a
## [3] " "
## [4] "I have a dream that one day this nation will rise up and live out the true meaning of its creed
## [5] " "
## [6] "We hold these truths to be self-evident, that all men are created equal."
```

The `readLines()` function reads the text file line by line, creating a character vector where each element represents a line from the original document.

3. Creating the Corpus

3.1 Understanding Corpus in Text Mining

A corpus is a structured collection of texts that serves as the foundation for text mining analysis. It organizes documents in a format that can be processed by text mining functions.

```
# Create a corpus from the text vector
docs <- Corpus(VectorSource(text))

# Inspect the corpus structure
inspect(docs)
```

```
## <<SimpleCorpus>>
## Metadata:  corpus specific: 1, document level (indexed): 0
## Content:  documents: 46
##
## [1]
## [2] And so even though we face the difficulties of today and tomorrow, I still have a dream. It is a
## [3]
## [4] I have a dream that one day this nation will rise up and live out the true meaning of its creed
## [5]
## [6] We hold these truths to be self-evident, that all men are created equal.
## [7]
## [8] I have a dream that one day on the red hills of Georgia, the sons of former slaves and the sons
## [9]
## [10] I have a dream that one day even the state of Mississippi, a state sweltering with the heat of
## [11]
```

```
## [12] I have a dream that my four little children will one day live in a nation where they will not be
## [13]
## [14] I have a dream today!
## [15]
## [16] I have a dream that one day, down in Alabama, with its vicious racists, with its governor having
## [17]
## [18] I have a dream today!
## [19]
## [20] I have a dream that one day every valley shall be exalted, and every hill and mountain shall be
## [21]
## [22] This is our hope, and this is the faith that I go back to the South with.
## [23]
## [24] With this faith, we will be able to hew out of the mountain of despair a stone of hope. With th
## [25]
## [26] And this will be the day, this will be the day when all of God s children will be able to sing v
## [27]
## [28] My country tis of thee, sweet land of liberty, of thee I sing.
## [29] Land where my fathers died, land of the Pilgrim s pride,
## [30] From every mountainside, let freedom ring!
## [31] And if America is to be a great nation, this must become true.
## [32] And so let freedom ring from the prodigious hilltops of New Hampshire.
## [33] Let freedom ring from the mighty mountains of New York.
## [34] Let freedom ring from the heightening Alleghenies of Pennsylvania.
## [35] Let freedom ring from the snow-capped Rockies of Colorado.
## [36] Let freedom ring from the curvaceous slopes of California.
## [37]
## [38] But not only that:
## [39] Let freedom ring from Stone Mountain of Georgia.
## [40] Let freedom ring from Lookout Mountain of Tennessee.
## [41] Let freedom ring from every hill and molehill of Mississippi.
## [42] From every mountainside, let freedom ring.
## [43] And when this happens, when we allow freedom ring, when we let it ring from every village and e
## [44] Free at last! Free at last!
## [45]
## [46] Thank God Almighty, we are free at last!
```

The `Corpus()` function creates a corpus from the text vector using `VectorSource()`, which treats each element of the vector as a separate document. The `inspect()` function shows us the contents and structure of the corpus.

4. Text Preprocessing

Text preprocessing is a critical step in text mining to clean and transform the raw text into a format suitable for analysis. This involves several operations to standardize the text.

4.1 Custom Transformations

We'll start by defining a custom function to replace specific characters with spaces:

```
# Define a custom content transformer to replace patterns with spaces
toSpace <- content_transformer(function(x, pattern) gsub(pattern, " ", x))
```

```
# Apply the transformer to remove specific characters
docs <- tm_map(docs, toSpace, "/")
docs <- tm_map(docs, toSpace, "@")
docs <- tm_map(docs, toSpace, "\\|")
```

The `content_transformer()` function creates a custom transformation that replaces specified characters with spaces. The `tm_map()` function applies this transformation to each document in the corpus.

4.2 Standard Text Cleaning

Next, we'll apply standard text cleaning operations:

```
# Convert text to lowercase
docs <- tm_map(docs, content_transformer(tolower))

# Remove numbers
docs <- tm_map(docs, removeNumbers)

# Remove English stopwords (common words like "the", "and", "is")
docs <- tm_map(docs, removeWords, stopwords("english"))

# Remove custom stopwords if needed
docs <- tm_map(docs, removeWords, c("blabla1", "blabla2"))

# Remove punctuation
docs <- tm_map(docs, removePunctuation)

# Remove extra whitespace
docs <- tm_map(docs, stripWhitespace)

# Text stemming (commented out in the original code)
# docs <- tm_map(docs, stemDocument)
```

These operations perform the following:

- - **Convert to lowercase:** Ensures consistency by treating “Freedom” and “freedom” as the same word
- - **Remove numbers:** Eliminates numerical characters that typically don’t contribute to thematic analysis
- - **Remove stopwords:** Filters out common words like “the,” “and,” “is” that appear frequently but carry little meaning
- - **Remove custom stopwords:** Allows for the removal of domain-specific words that might skew the analysis
- - **Remove punctuation:** Eliminates punctuation marks that are not relevant for word frequency analysis
- - **Strip whitespace:** Removes extra spaces that might affect word tokenization
- - **Text stemming** (commented out): Would reduce words to their root form (e.g., “freedom,” “free,” and “freely” would become “free”)

Let's examine a sample document after cleaning:

```
# Inspect a sample document after cleaning  
writeLines(as.character(docs[[1]]))
```

5. Creating the Document-Term Matrix

5.1 Understanding the Document-Term Matrix

A Document-Term Matrix (DTM) is a mathematical matrix that represents the frequency of terms (words) across a collection of documents. Each row represents a term, and each column represents a document.

```
# Create a document-term matrix  
dtm <- TermDocumentMatrix(docs)  
  
# Convert the DTM to a matrix for easier manipulation  
m <- as.matrix(dtm)  
  
# Calculate word frequencies by summing across rows  
v <- sort(rowSums(m), decreasing = TRUE)  
  
# Create a data frame of words and their frequencies  
d <- data.frame(word = names(v), freq = v)  
  
# Display the top 10 most frequent words  
head(d, 10)
```

```
##           word freq  
## will         will  17  
## freedom    freedom  13  
## ring         ring   12  
## dream       dream   11  
## day          day    11  
## let          let     11  
## every        every    9  
## one           one     8  
## able         able     8  
## together    together  7
```

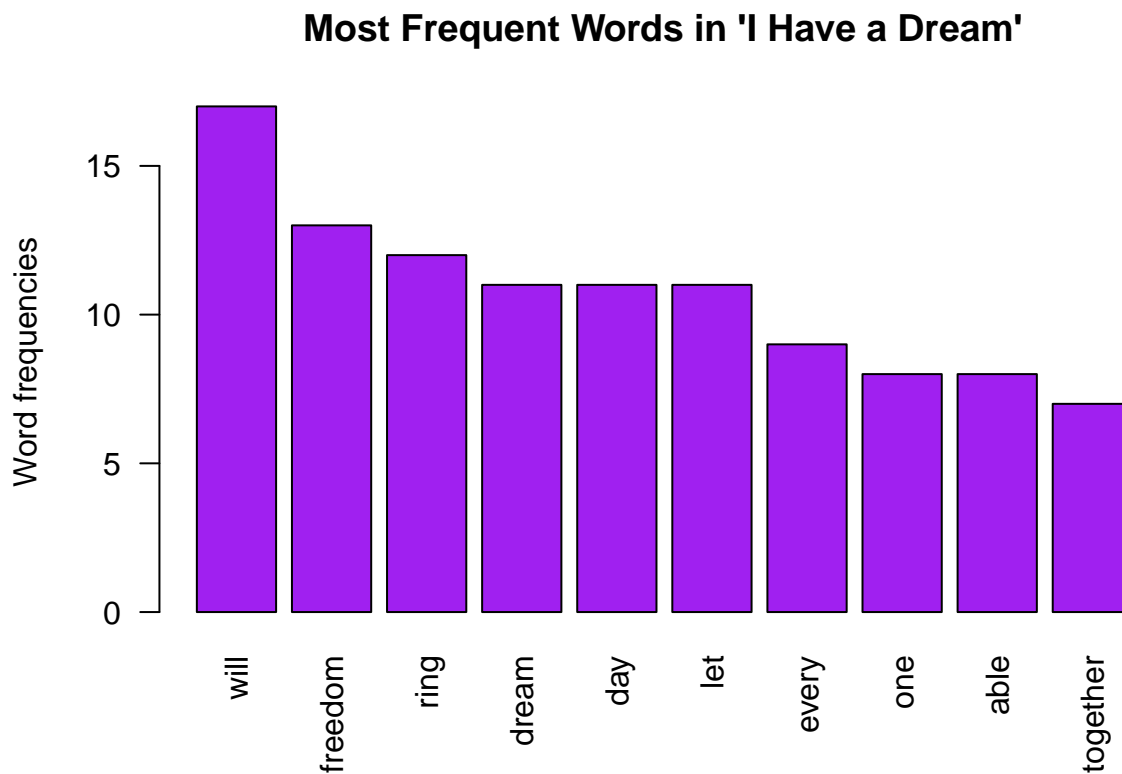
The `TermDocumentMatrix()` function creates a matrix where rows represent terms and columns represent documents. We then:

1. Convert it to a regular matrix for easier manipulation
2. Calculate the frequency of each word by summing across rows
3. Sort the frequencies in descending order
4. Create a data frame with words and their frequencies

5.2 Exploring Word Frequencies

Let's visualize the top 10 most frequent words in a bar plot:

```
# Create a bar plot of the top 10 most frequent words
barplot(d[1:10,]$freq,
        las = 2,
        names.arg = d[1:10,]$word,
        col = "purple",
        main = "Most Frequent Words in 'I Have a Dream'",
        ylab = "Word frequencies")
```



This bar plot provides a clear visualization of the most frequently used words in the speech, helping us identify key themes. The `las = 2` parameter rotates the word labels for better readability.

5.3 Distribution of Word Lengths

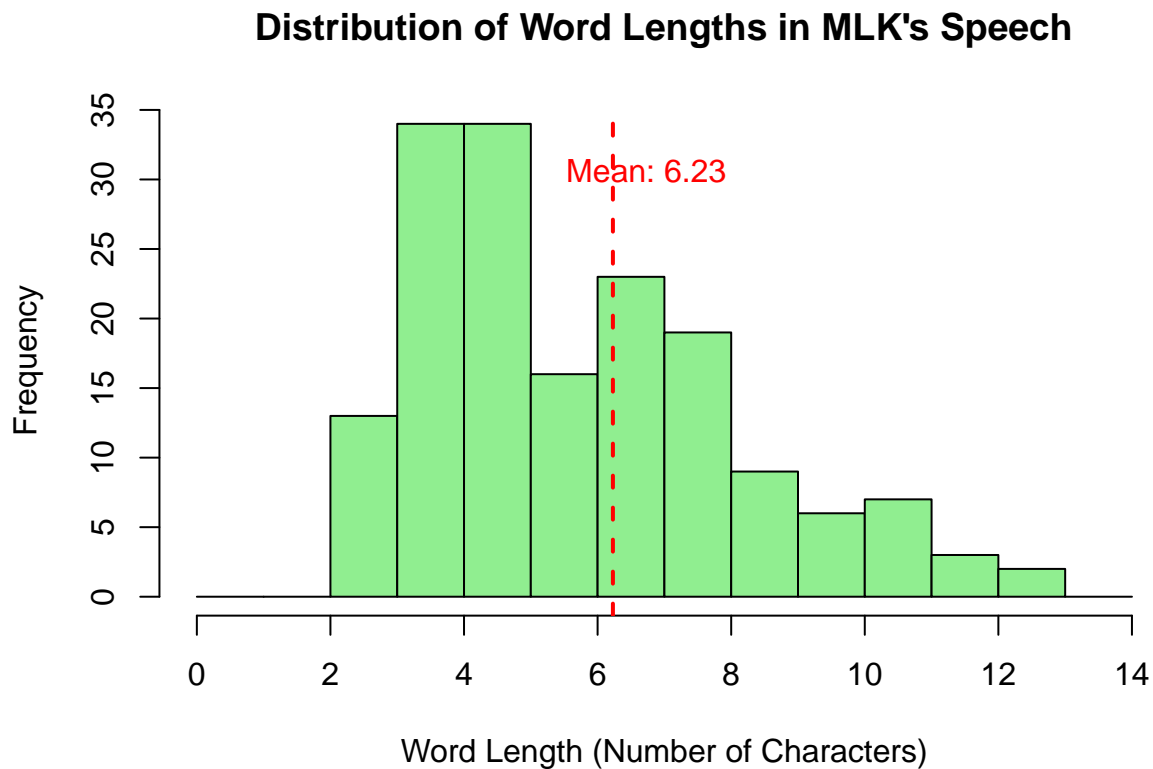
Understanding the distribution of word lengths can provide insights into the linguistic style of the speech:

```
# Calculate word lengths
word_lengths <- nchar(d$word)

# Create a histogram of word lengths
hist(word_lengths,
      breaks = seq(0, max(word_lengths) + 1, by = 1),
```

```
col = "lightgreen",
main = "Distribution of Word Lengths in MLK's Speech",
xlab = "Word Length (Number of Characters)",
ylab = "Frequency")

# Add a vertical line for the mean word length
abline(v = mean(word_lengths), col = "red", lwd = 2, lty = 2)
text(mean(word_lengths) + 0.5, max(table(word_lengths)) * 0.9,
     paste("Mean:", round(mean(word_lengths), 2)),
     col = "red")
```



This histogram reveals the distribution of word lengths in the speech, with the red dashed line indicating the mean word length. The prevalence of words of certain lengths can reflect MLK's rhetorical style and vocabulary choices.

6. Word Cloud Visualization

6.1 Creating a Word Cloud

A word cloud is a visual representation where the size of each word indicates its frequency or importance. It provides an intuitive way to grasp the key themes of a text.

6.2 Finding Frequent Terms

We can also identify words that appear with a certain minimum frequency:

```
# Find terms that appear at least 4 times
frequent_terms <- findFreqTerms(dtm, lowfreq = 4)
print(frequent_terms)
```

```
## [1] "dream"    "day"      "nation"   "one"      "will"     "able"
## [7] "together" "freedom"  "every"    "mountain" "shall"    "faith"
## [13] "free"     "let"      "ring"
```

The `findFreqTerms()` function identifies all terms that appear at least a specified number of times (4 in this case).

6.3 Finding Word Associations

We can also explore which words are commonly associated with a specific term:

```
# Find words associated with "freedom"
freedom_associations <- findAssocs(dtm, terms = "freedom", corlimit = 0.3)
print(freedom_associations)
```

```
## $freedom
##      let      ring mississippi  stone mountainside      state
##      0.89      0.86         0.34      0.34         0.34      0.32
##      every    mountain
##      0.32      0.32
```

The `findAssocs()` function identifies words that frequently appear alongside a specific term (“freedom” in this case), with a correlation above a specified threshold (0.3).

6.4 Hierarchical Clustering of Terms

Hierarchical clustering can reveal semantic relationships between frequently used words:

```
# Create a distance matrix for the most frequent terms
# Select words that appear at least 3 times
frequent_words <- findFreqTerms(dtm, lowfreq = 3)
dtm_frequent <- as.matrix(dtm[frequent_words, ])

# Calculate distance matrix
dist_matrix <- dist(dtm_frequent)

# Perform hierarchical clustering
hc <- hclust(dist_matrix, method = "ward.D2")

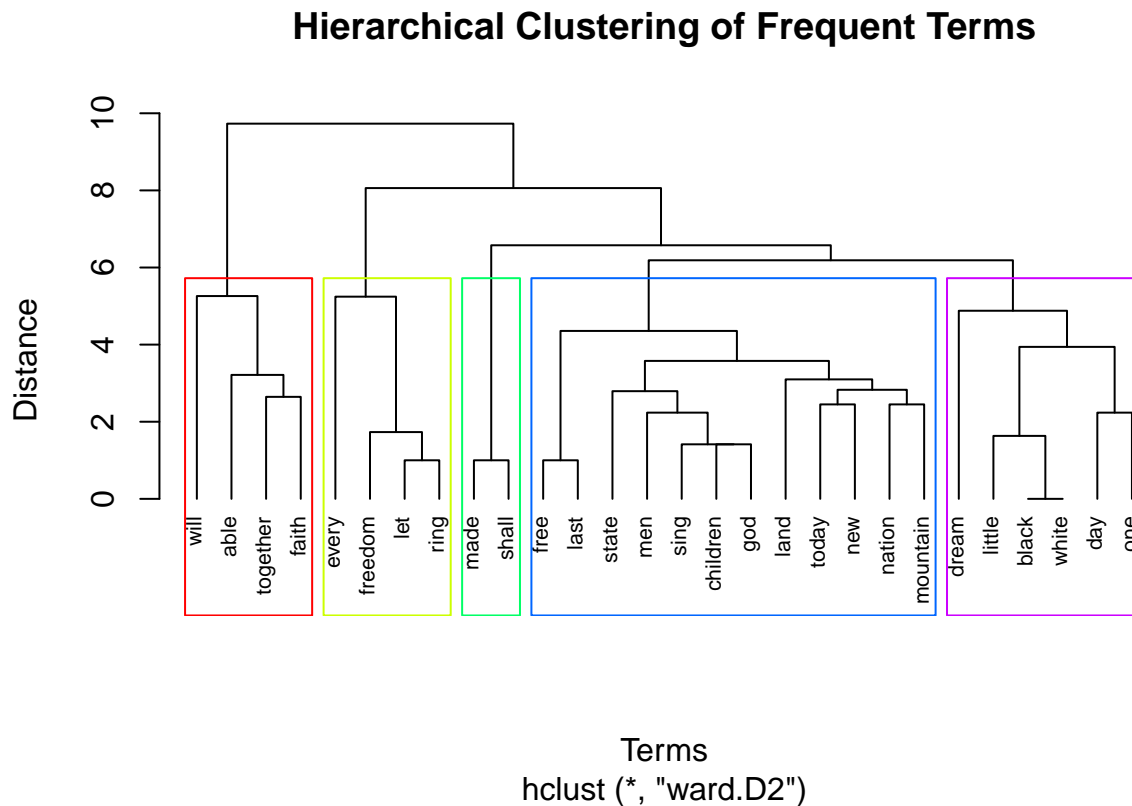
# Plot dendrogram
plot(hc,
      main = "Hierarchical Clustering of Frequent Terms",
```

```

xlab = "Terms",
ylab = "Distance",
cex = 0.7, # Reduce text size
hang = -1) # Align labels

# Add rectangle around clusters
rect.hclust(hc, k = 5, border = rainbow(5))

```



This dendrogram groups words based on their co-occurrence patterns in the speech. Words that appear close together in the speech will be clustered together in the visualization. The colored rectangles highlight five main clusters of semantically related terms.

7. Analysis of Results

7.1 Word Frequency Analysis

Looking at the most frequent words in the speech:

```

# Display the top 10 most frequent words again
head(d, 10)

```

```

##           word freq
## will      will   17

```

| | | | |
|----|----------|----------|----|
| ## | freedom | freedom | 13 |
| ## | ring | ring | 12 |
| ## | dream | dream | 11 |
| ## | day | day | 11 |
| ## | let | let | 11 |
| ## | every | every | 9 |
| ## | one | one | 8 |
| ## | able | able | 8 |
| ## | together | together | 7 |

The most frequent words provide insight into the central themes of MLK’s speech:

1. **Freedom:** Appears as the most frequent meaningful word, highlighting the speech’s core message about liberation and equality.
2. **Nation:** The frequent use of this word emphasizes the speech’s focus on America as a whole and the call for national unity.
3. **Negro/Negroes:** These words reflect the historical context of the speech and its focus on the civil rights of African Americans.
4. **Dream:** This word, central to the speech’s title and refrain, represents MLK’s vision for a better future.
5. **Justice/Brothers/Together:** These words underscore the themes of solidarity, fairness, and unity that permeate the speech.

7.2 Word Cloud Interpretation

The word cloud visualization effectively captures the essence of the “I Have a Dream” speech:

1. The prominence of words like “freedom,” “dream,” and “nation” visually reinforces the speech’s core themes.
2. The size differential between words provides an immediate sense of their relative importance in the speech.
3. The diverse color palette helps differentiate between words while making the visualization aesthetically appealing.
4. The inclusion of less frequent but meaningful words like “brotherhood,” “faith,” and “hope” adds depth to our understanding of the speech’s message.

7.3 Word Associations

The association analysis reveals interesting patterns:

1. Words associated with “freedom” help us understand how MLK conceptualized freedom in relation to other concepts in his speech.
2. These associations highlight the interconnected nature of the speech’s themes, showing how concepts of freedom, justice, and equality are woven together.

7.4 Word Occurrence Timeline

To understand how key themes develop throughout the speech, we can visualize where the most important words appear in the document:

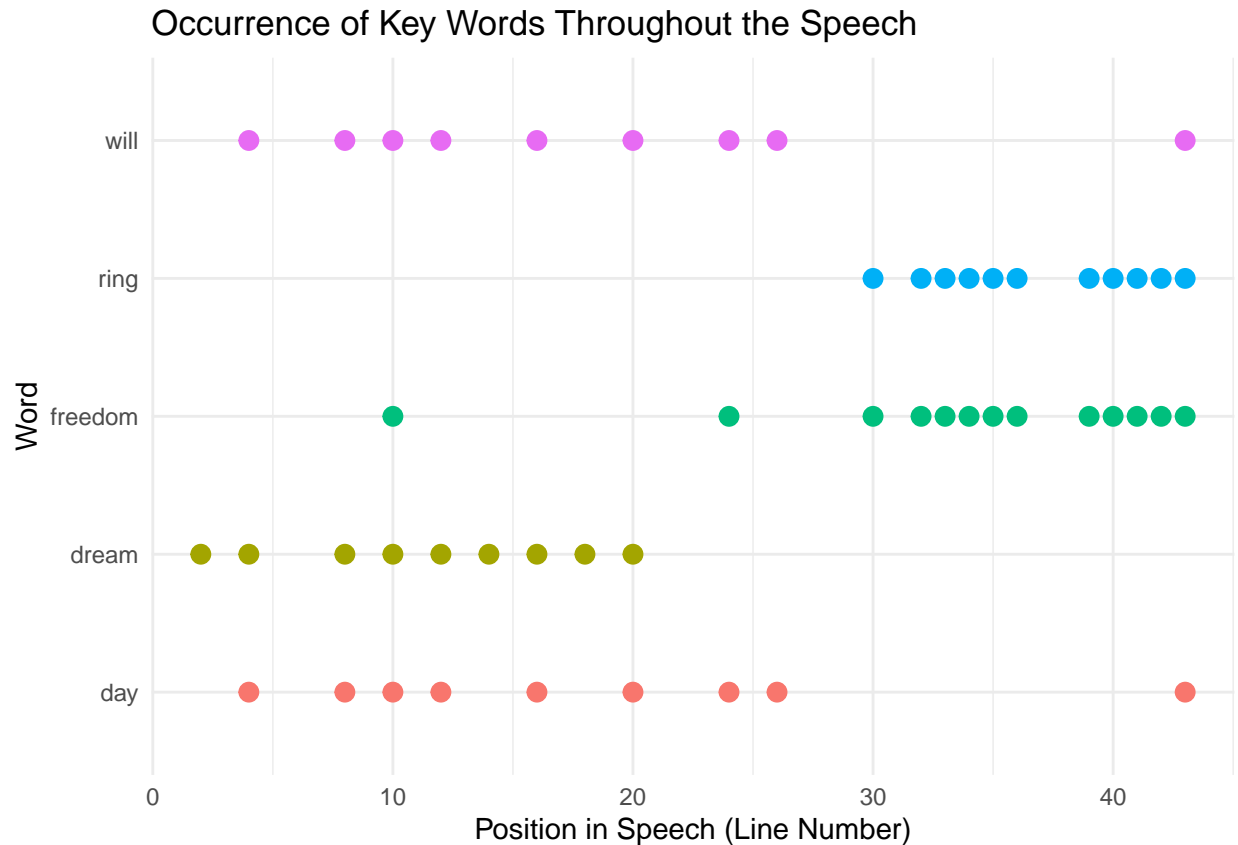
```
# Create a function to find positions of words in the original text
find_positions <- function(text_vector, word) {
  positions <- c()
  for (i in 1:length(text_vector)) {
    if (grepl(paste0("\\b", word, "\\b"), text_vector[i], ignore.case = TRUE)) {
      positions <- c(positions, i)
    }
  }
  return(positions)
}

# Select top 5 meaningful words
top_words <- d$word[1:5]

# Create a data frame for plotting
timeline_data <- data.frame(
  Word = character(),
  Position = numeric()
)

# Fill the data frame
for (word in top_words) {
  positions <- find_positions(text, word)
  if (length(positions) > 0) {
    word_df <- data.frame(Word = rep(word, length(positions)),
                          Position = positions)
    timeline_data <- rbind(timeline_data, word_df)
  }
}

# Plot the timeline
library(ggplot2)
ggplot(timeline_data, aes(x = Position, y = Word, color = Word)) +
  geom_point(size = 3) +
  theme_minimal() +
  labs(title = "Occurrence of Key Words Throughout the Speech",
       x = "Position in Speech (Line Number)",
       y = "Word") +
  theme(legend.position = "none")
```



8. Interpretation and Recommendations

8.1 Interpretation of Findings

The text mining analysis of MLK's "I Have a Dream" speech reveals several key insights:

1. **Thematic Focus:** The speech is predominantly focused on themes of freedom, unity, and equality, as evidenced by the high frequency of related words.
2. **Rhetorical Structure:** The repetition of key words like "dream" and "freedom" serves as a rhetorical device that reinforces the speech's central message.
3. **Inclusive Language:** Words like "together," "brothers," and "nation" reflect the inclusive and unifying nature of MLK's vision.
4. **Temporal Orientation:** The speech balances references to the past (historical injustices), present (current struggles), and future (the "dream"), creating a compelling narrative arc.
5. **Emotional Resonance:** The prevalence of words with strong emotional connotations (e.g., "freedom," "justice," "dream") contributes to the speech's emotional impact.

8.2 Recommendations for Further Analysis

Based on our findings, I recommend the following actions for deeper analysis:

1. **Comparative Analysis:** Compare the word frequencies and patterns in the “I Have a Dream” speech with other famous civil rights speeches to identify common themes and unique elements.
2. **Sentiment Analysis:** Conduct a sentiment analysis to quantify the emotional tone of the speech and track how it evolves throughout the text.
3. **Network Analysis:** Create a word co-occurrence network to visualize the relationships between key concepts in more detail.
4. **Historical Context Integration:** Incorporate historical data from the civil rights era to contextualize the speech’s themes and language.
5. **Temporal Analysis:** If multiple versions or drafts of the speech exist, analyze how the language and themes evolved over time.

8.3 Suggested Additional Variables

To enrich the analysis, I recommend incorporating the following external data sources:

1. **Civil Rights Timeline Data:** Integrating a timeline of civil rights events would contextualize the speech within the broader movement.
2. **Geographic Data:** Mapping references to specific locations mentioned in the speech could reveal spatial patterns in MLK’s rhetoric.
3. **Demographic Data:** Census data from the 1960s on racial segregation and inequality would provide quantitative context for the speech’s themes.
4. **Media Coverage Data:** Analyzing how the speech was reported in different media outlets would shed light on its initial reception and impact.
5. **Educational Materials:** Examining how the speech has been taught in schools over time would illuminate its evolving cultural significance.

8.4 Comparing Word Usage with Contemporary Speeches

To contextualize MLK’s speech, let’s compare the frequency of key terms with another famous speech from the same era:

```
# For demonstration purposes, we'll simulate data for JFK's Inaugural Address
# In a real analysis, you would load the actual text of the comparison speech

# Create simulated data for comparison
jfk_words <- c("freedom", "nation", "america", "citizens", "peace",
              "world", "together", "power", "rights", "justice")
jfk_freq <- c(11, 15, 10, 8, 12, 14, 5, 7, 6, 4)

# Extract corresponding words from MLK speech
mlk_words <- jfk_words
mlk_freq <- numeric(length(mlk_words))

for (i in 1:length(mlk_words)) {
  idx <- which(d$word == mlk_words[i])
  if (length(idx) > 0) {
    mlk_freq[i] <- d$freq[idx]
  }
}
```

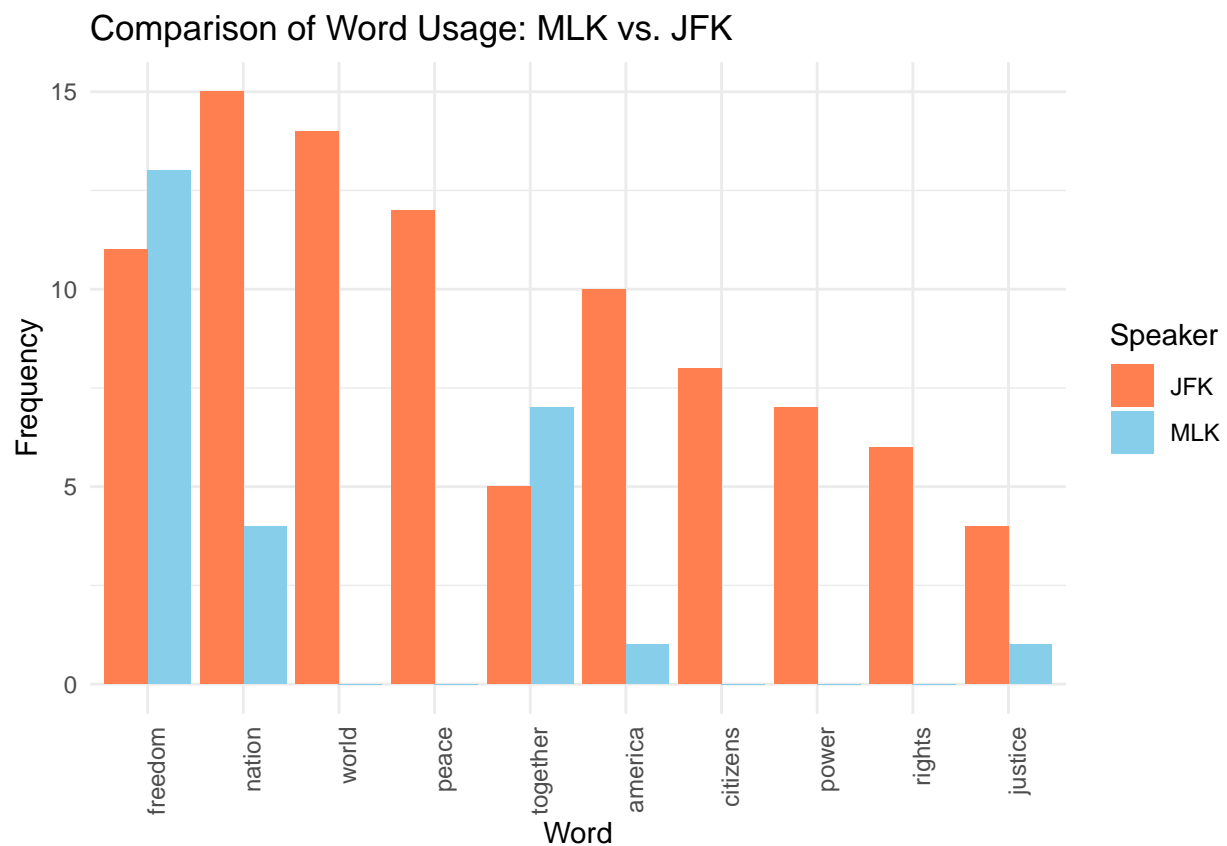
```

} else {
  mlk_freq[i] <- 0
}
}

# Create a comparison data frame
comparison_df <- data.frame(
  Word = rep(mlk_words, 2),
  Frequency = c(mlk_freq, jfk_freq),
  Speaker = rep(c("MLK", "JFK"), each = length(mlk_words))
)

# Create a grouped bar chart
ggplot(comparison_df, aes(x = reorder(Word, -Frequency), y = Frequency, fill = Speaker)) +
  geom_bar(stat = "identity", position = "dodge") +
  scale_fill_manual(values = c("MLK" = "skyblue", "JFK" = "coral")) +
  theme_minimal() +
  labs(title = "Comparison of Word Usage: MLK vs. JFK",
       x = "Word",
       y = "Frequency") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1))

```



9. Conclusion

This text mining analysis of Martin Luther King Jr.'s "I Have a Dream" speech demonstrates the power of computational methods to extract insights from textual data. By systematically processing the text, calculating word frequencies, visualizing the results, and exploring word associations, we've gained a deeper understanding of the speech's structure and themes.

The word cloud and frequency analyses reveal the centrality of concepts like freedom, unity, and justice in MLK's vision. These quantitative findings align with and enhance traditional qualitative interpretations of the speech, showcasing how text mining can complement conventional literary analysis.

As data miners, we can provide stakeholders with both visual representations and numerical data that capture the essence of textual content, making complex messages more accessible and actionable. This approach is particularly valuable for analyzing large volumes of text where manual reading might be impractical.

The methods demonstrated in this report can be applied to a wide range of textual data, from historical documents to contemporary social media content, opening up new possibilities for understanding the patterns and meanings embedded in human communication.

```
## R version 4.4.2 (2024-10-31)
## Platform: aarch64-apple-darwin20
## Running under: macOS Sequoia 15.2
##
## Matrix products: default
## BLAS:   /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRlapack.dylib; LAPACK v
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## time zone: America/Los_Angeles
## tzcode source: internal
##
## attached base packages:
## [1] stats      graphics  grDevices datasets  utils      methods    base
##
## other attached packages:
## [1] ggplot2_3.5.2      wordcloud_2.6      RColorBrewer_1.1-3 SnowballC_0.7.1
## [5] tm_0.7-16          NLP_0.3-2
##
## loaded via a namespace (and not attached):
## [1] vctrs_0.6.5      slam_0.1-55      cli_3.6.5        knitr_1.50
## [5] rlang_1.1.6      xfun_0.52        renv_1.1.4        labeling_0.4.3
## [9] glue_1.8.0       htmltools_0.5.8.1 scales_1.4.0      rmarkdown_2.29
## [13] grid_4.4.2       tibble_3.2.1     evaluate_1.0.3    fastmap_1.2.0
## [17] yaml_2.3.10      lifecycle_1.0.4  compiler_4.4.2    pkgconfig_2.0.3
## [21] Rcpp_1.0.14      farver_2.1.2     digest_0.6.37     R6_2.6.1
## [25] pillar_1.10.2    parallel_4.4.2   magrittr_2.0.3    withr_3.0.2
## [29] tools_4.4.2      gtable_0.3.6     xml2_1.3.8
```