

# Statistical Satire: A Data Analytics Approach to Stand-Up Comedy

Chris Heffernan

Sunday 10<sup>th</sup> December, 2023

# Contents

Abstract and Introduction	3
Data Description and Exploratory Data Analysis	4
Analysis	5
Model Development and Application of Model(s)	9
Conclusions and Discussion	14
Bibliography	15
Code Appendix	16

# Abstract and Introduction

With the modern intersection of stand-up comedy and digital mediums, stand-up comedy is undergoing its own renaissance. Initially experiencing stand-up comedy feels like witnessing magic unfold on stage. However, delving behind the curtain reveals an appreciation for the artistry involved, showing how comedians skillfully construct sets – molding words and timing to craft moments that elicit genuine laughter from their audience. This project leverages analytical tools and concepts to unravel the underlying patterns and trends within comedic content.

Motivated by the desire to understand the dynamics of humor and audience engagement, the report poses questions about prevalent humor styles, linguistic nuances, and the impact of audience reactions. Leveraging sentiment analysis, word frequency analysis, and clustering techniques, the analysis unveils distinct humor clusters, identifies frequently used comedic elements, and explores the interplay between language and laughter.

Results reinforce the complexity of analyzing stand-up comedy by its words. Sentiment analysis exposes the predominantly positive nature of comedic language, while word frequency analysis highlights the significance of specific words in crafting relatable and resonant humor. Clustering techniques unveil the nuanced topics discussed in stand-up specials, emphasizing the individuality of each comedian’s style.

This project not only pays tribute to the field of data analytics but also provides insights for working comedians. The findings showcase the potential of data-driven approaches in decoding the intricate art of stand-up comedy and in doing so, provide a level of understanding that makes the seemingly intimidating nature of the craft into an accessible endeavor.

# Data Description and Exploratory Data Analytics

While non-verbal cues offer great value in eliciting laughter, the essence of humor lies fundamentally in a comedian’s choice of words. A stand-up special serves as a compilation of a comedian’s most successful jokes within a specified time frame, so the dataset is exclusively centered around the transcript of stand-up specials to conduct essential statistical analyses. Renowned comedians traditionally release their specials on platforms like Netflix, HBO, and Comedy Central, attracted by the substantial payout. However, a notable shift has occurred, where comedians are increasingly opting for YouTube as a platform for special releases, providing cost-free access and in turn reaching a broader audience. Exploiting YouTube’s open-source API capabilities, this project used closed-captioning functionalities to extract content from various comedy sets. This extraction facilitated the content and timing analyses.

The overall dataset was comprised of 50 specials featuring 50 different comedians and nearly 50 hours of content. Each individual special had its own csv file containing text information, the time with which it took to deliver those words, and its timestamp in the video. These columns were labeled *text*, *start*, and *duration* as directly sourced from the transcript API. The text column not only featured the comic’s words, but also behaviors from the special including applause, music, and most importantly laughter. This facilitated manipulating individual specials and creating dataframes in R.

Initial EDA was performed on individual comedy sets. This process involved an in-depth examination of each set’s unique characteristics, including audience behavior, sentiment analysis, word frequency patterns, and cadence. Building on the insights gained from individual set analyses, a comprehensive EDA was conducted on all comedy sets combined. This allowed for the aggregation of patterns and trends across the entire dataset. Preprocessing the text consisted of removing any unrecognizable characters and handling YouTube’s approach to cuss words. YouTube transcripts don’t provide text for certain profanities and instead replace the word with [\_\_] in its closed captions. Consequently, this symbol was separated from the text, and allowed for some interesting analysis. Some datasets also did not feature behaviors like laughter, applause, and music, so this was taken into account when developing some of the graphics and statistics. Another area for preprocessing was that the captions sometimes confused laughter with music. Given the unlikely event for a stand-up special to feature music in the set anywhere besides the beginning and end, music detected outside of the beginning and end was assumed to be laughter.

# Analysis

Sentiment analysis was conducted in the EDA to understand the tone and mood expressed in the text. The sentiment analysis was conducted after removing explicit. This was because there was no way to automate replacing the profanity symbol with the actual word. However, stand-up often undergoes controversy as some comics can be characterized as vulgar. In removing the explicit and focusing on all other words, the sampled comedy sets proved to be much showcase little negative content. The SentimentAnalysis library in R served as the primary tool for sentiment analysis. This library provides functionalities to assess sentiment in textual data and is particularly well-suited for applications in the field of natural language processing.

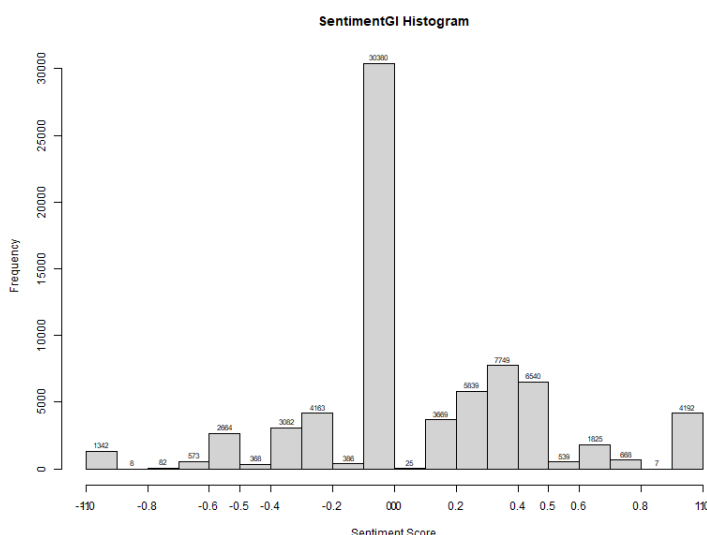


Figure 1: Histogram of Sentiment Analysis

The language seems to be remarkably neutral, but if one were to polarize between positive and negative, the histogram showcases that the majority of words chosen across all specials are positive. The sentiment analysis in this context is more concerned with capturing the humorous and positive aspects rather than negative sentiments. Profanities are often associated with negative emotions, but in the realm of comedy, they can be used playfully or satirically to enhance the humor without necessarily conveying a negative sentiment. The special that featured the most profanity was Andrew Schulz's special titled "Infamous" with 169 documented cases of profanity. Ali Siddiq's special "The Domino Effect" showed zero cases of profanity, and their histograms prove how similar their sentiments are.

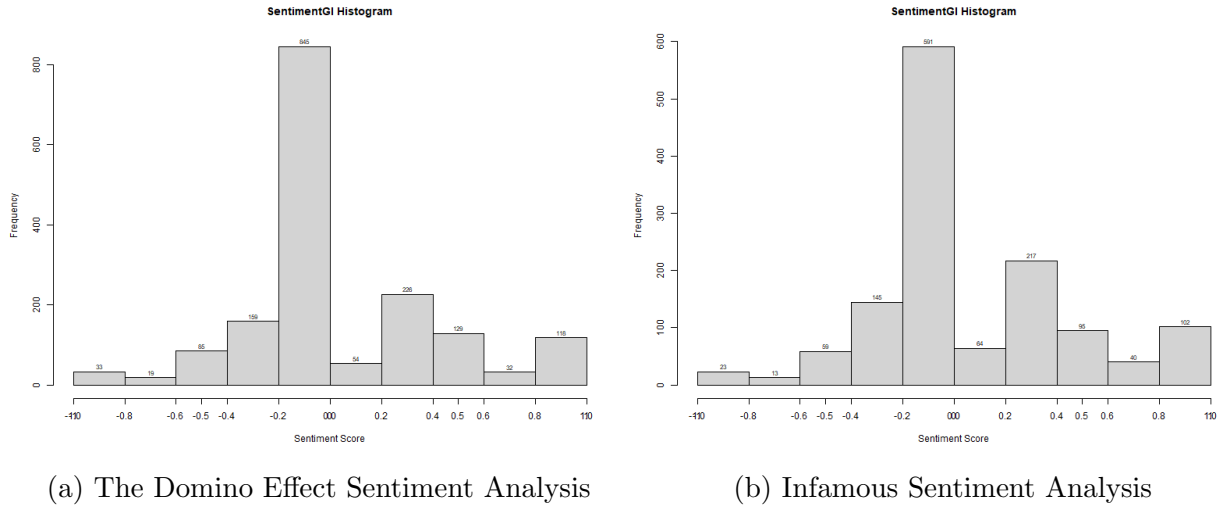


Figure 2: Comparative Sentiment Analysis

Next, analysis was conducted to determine which words were most prevalent in a stand-up comedy special. Identifying the most common words helps to gain a quick overview of the main themes or topics present in stand-up specials. These words are likely to be key elements in understanding a comedian’s style, recurring themes, or the overall tone of the performance.

### Word Frequencies Without Stop Words

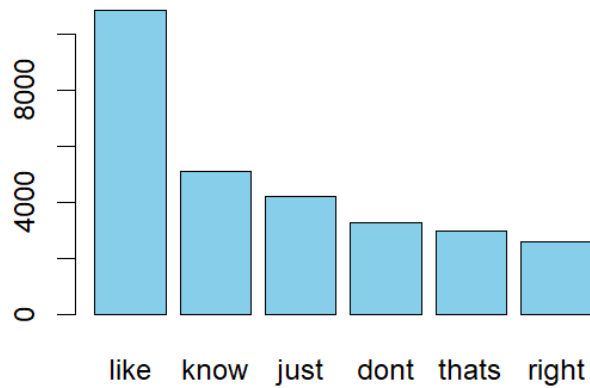


Figure 3: Word Frequency

The significance of "like" being the most common word in stand-up specials could be attributed to a couple of different factors. Stand-up comedians often use comparisons and similes as a humorous device. The word "like" is commonly used to introduce such comparisons, making it a frequent occurrence. In spoken language, especially in casual and

conversational settings like stand-up comedy, people tend to use filler words or phrases. "Like" is one such filler word that is often used to provide emphasis or convey a sense of approximation. If "like" is a frequently used word, it suggests that the audience may respond positively to comparisons and relatable content, adding a layer to why comics might choose to incorporate this in their act.

Next, the EDA considered a comic's timing and how cadence impacts the nature of stand-up. Timing and cadence play a pivotal role in stand-up comedy, as the way a comedian paces their speech, the emphasis on certain words, or the use of pauses can significantly contribute to the overall comedic effect. The transcript data was initially cleaned to ensure uniformity and eliminate potential discrepancies. Speech rate, measured in words per second, was calculated for each segment of the stand-up specials. This metric serves as an indicator of how quickly or slowly a comedian delivers their lines. The cubic spline interpolation technique was applied to individual stand-up specials. Pacing throughout the sets was sporadic and almost sinusoidal, so linear interpolation was not effective in capturing cadence, but the locally estimated scatterplot smoothing (LOESS) method was employed to also offer a simplified visualization of speech rate changes over time. Cubic spline interpolation helped in smoothing out the speech rate fluctuations, providing a clearer view of the underlying trends. The individual spline points from all stand-up specials were then aggregated into one comprehensive dataset. Here is a sample of the cadence featured in Jim Bruer's "Country Boy Will Survive":

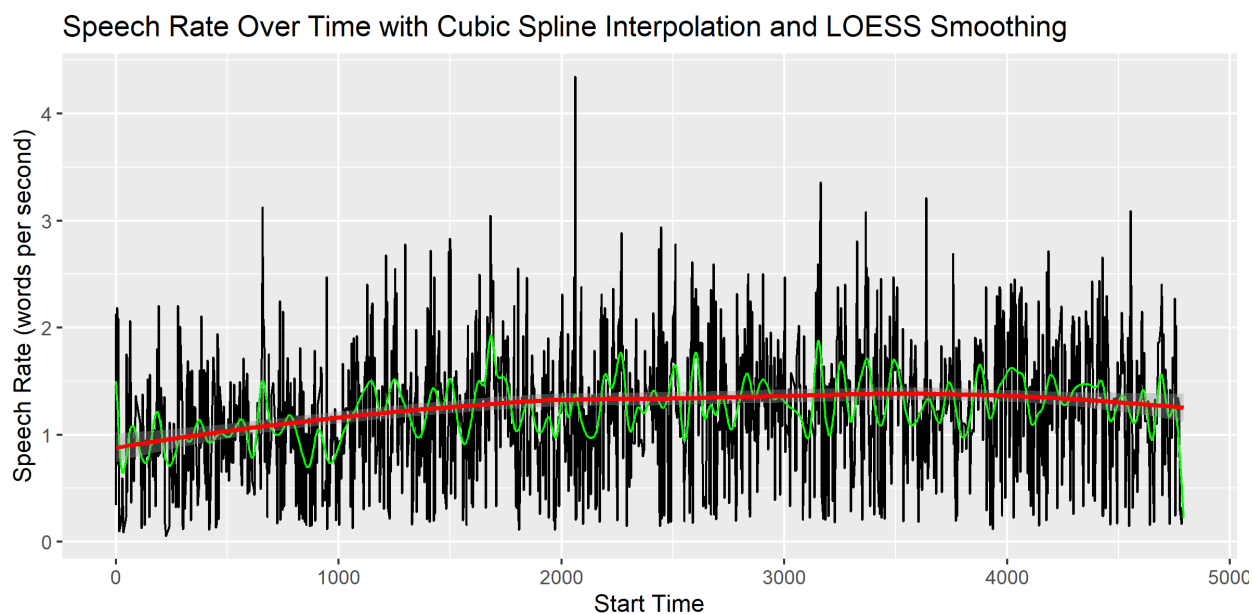


Figure 4: "Country Boy Will Survive" Cadence

The aggregate cadence interpolation shows how comics tend to oscillate around 2 words per second, but that rate accounts for the audience’s laughing and applause. The end of the special intuitively shows how speech rate decreases, but it appears that speech rate increases around 75% of the way through the special.

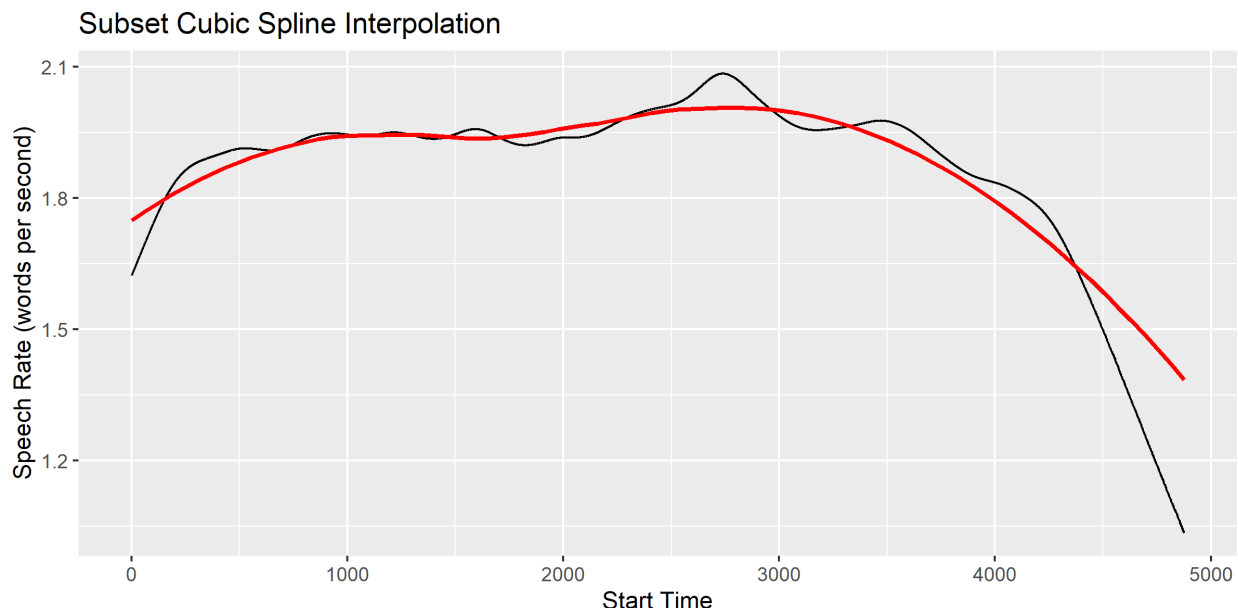


Figure 5: Aggregate Cadence Plot

The final component of the EDA was dedicated to determining the topics discussed in stand-up sets. In order to do so, the text information was input to a topic modeling method known as Latent Dirichlet Allocation (LDA). LDA is a type of statistical model used for discovering abstract topics within a collection of documents. The LDA model assumes that each document is a mixture of a small number of topics and that each word in the document is attributable to one of the document’s topics, which seemed perfect for the various stand-up sets.

Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6	Topic 7	Topic 8	Topic 9	Topic 10
"right"	"just"	"like"	"like"	"thats"	" "	"know"	"applause"	"people"	"get"
"got"	"make"	"youre"	"hes"	"one"	"going"	"dont"	"never"	"think"	"yeah"
"now"	"around"	"gonna"	"well"	"time"	"man"	"want"	"ive"	"theyre"	"can"
"back"	"cant"	"guys"	"guy"	"good"	"audience"	"didnt"	"music"	"theres"	"see"
"little"	"trying"	"much"	"shes"	"thing"	"said"	"say"	"ever"	"love"	"ill"
"come"	"still"	"way"	"look"	"first"	"laughing"	"even"	"two"	"really"	"take"
"show"	"saying"	"god"	"feel"	"day"	"dude"	"mean"	"years"	"lot"	"tell"
"okay"	"put"	"thank"	"kind"	"every"	"something"	"need"	"life"	"always"	"okay"
"new"	"give"	"getting"	"also"	"went"	"hey"	"talk"	"thought"	"white"	"laugh"
"talking"	"away"	"coming"	"pretty"	"whole"	"whats"	"doesnt"	"old"	"women"	"let"

Surprisingly, the LDA did not reveal any distinct and easily interpretable topics. The lack of coherent topics indicates that each of the specials makes a unique imprint. Analysis on individual specials could provide more valuable insights into topics discussed, but that would not contribute to the current mission to understand the comedy landscape as a whole.



## Model Development and Application of Model(s)

The 2 models constructed for this project are used to get a more robust understanding of the topics discussed in a comedy set and the cadence. Initial LDA analysis for the all-encompassing dataset proved to be ineffective, but through clustering techniques, it's possible to how topics are interconnected in individual sets and the stand-up landscape as a whole. The model built to understand cadence utilizes regression to correlate the duration of laughter and the sentiment of the words leading up to the laughter. With this we could potentially identify if people tend to laugh harder at raunchy humor or not.

For clustering topics using Latent Dirichlet Allocation results, hierarchical clustering is the most suitable method. It creates a hierarchy of clusters, providing a visual representation of topic relationships. The code takes a list of LDA models where one document represents one comedy set, and performs hierarchical clustering on the topics using the beta matrix of each LDA model. The beta matrix represents the probability distribution of terms in each topic. Each row in the beta matrix corresponds to a topic, and each column corresponds to a unique term in the vocabulary. The values in the matrix represent the probability of a term occurring in a particular topic. The beta matrix is particularly useful for clustering topics because it captures the underlying structure of topics in terms of their word distributions. By using the beta matrix for hierarchical clustering, the code is essentially comparing topics based on how similar their word distributions are. This can help identify topics that share common terms or have similar themes. The clustering is executed using Ward's method (`method = "ward.D2"`) based on the Euclidean distance between topics. Ward's method is known for minimizing the total within-cluster variance during the merging of clusters. In the context of stand-up comedy transcripts, where topics may represent common themes or recurring jokes, minimizing within-cluster variance is desirable.

The distance matrix quantifies the level of similarity between the various topics in the dataset. When initially building the model, distance matrix was calculated using the default Euclidean distance which measures a straight line between topics in the multidimensional space. The dendrogram resulting from hierarchical clustering shows the hierarchy of topic similarities. The branches of the dendrogram represent the clusters formed at different levels of similarity. Similar topics are joined together at lower levels of the dendrogram, while more dissimilar topics are joined at higher levels. Featured here is the dendrogram produced from Anjelah Johnson-Reyes' special "Say I Won't":

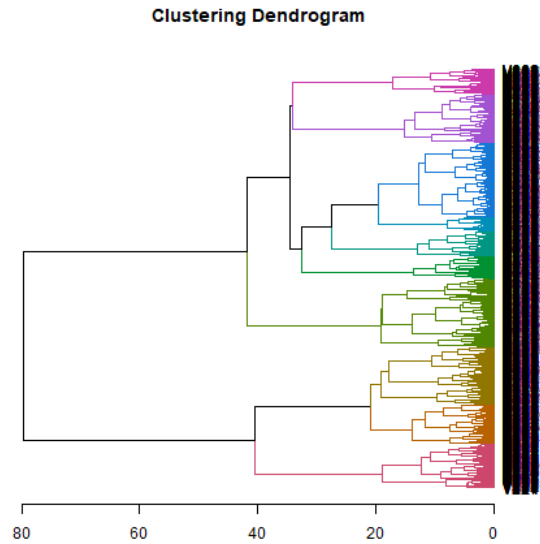


Figure 6: "Say I Won't" Dendrogram

There are 10 distinct colors specifying the partition of the set into 10 clusters based on the results. There are so many words to choose from and the variability in topics can be so diverse. Comics like Mark Normand deliver one-line jokes where each joke is completely unrelated to the next, but a comic like Mike Birbiglia is notable for building his set around an overarching theme. The dataset included a transcript from Mark Normand's special "Don't Be Yourself", which showed interesting results when compared to Anjelah Johnson's special.

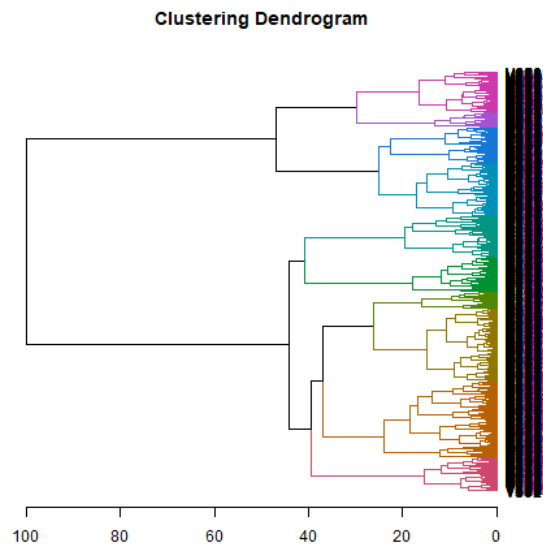


Figure 7: "Don't Be Yourself" Dendrogram

	Say I Won't	Don't Be Yourself
merge	2642	3236
height	1321	1618
order	1322	1619
labels	1322	1619

Using the summary function, the observed difference in the number of merges between 'Say I Won't' and 'Don't Be Yourself' suggests intriguing distinctions in the structure of these comedy sets. The higher number of merges in 'Don't Be Yourself,' particularly in comparison to 'Say I Won't,' implies a more intricate hierarchy of topics. This metric was discovered unintentionally, derived from the hierarchical clustering process, effectively quantifying the range and diversity of topics discussed in a stand-up set. In essence, a higher number of merges reflects a broader spectrum of themes and subjects covered within the same time duration. This novel metric serves as a valuable quantitative measure for assessing the richness and variety of content in stand-up comedy specials, offering a nuanced perspective beyond traditional topic modeling techniques. As comedy specials typically share a similar time frame, this metric becomes especially pertinent in gauging the efficiency with which comedians navigate through diverse topics, providing a unique lens through which to evaluate the distinct comedic styles of performers like Mark Normand and Anjelah Johnson-Reyes.

When generating the cluster for all combined sets, each bta matrix was stacked into 1 matrix and transposed which showcased the total different types of topics discussed based on the similarity in word choices. The intention was to determine if there were some overarching theme similarities across all the data. Finding a massive similarity would expose that comedians tend to drift to certain topics when producing specials as they could be classified as the funniest. As it turns out, and an homage to the art, there are extremely low similarities in the topics discussed in comedy specials. This indicates that comedians don't rely on the subject matter in order to produce humor. This information does not make it clear as to what a comic needs to say in order to be funny, but in fact eliminates any misconception that what a comedian talks about needs to be inherently funny.

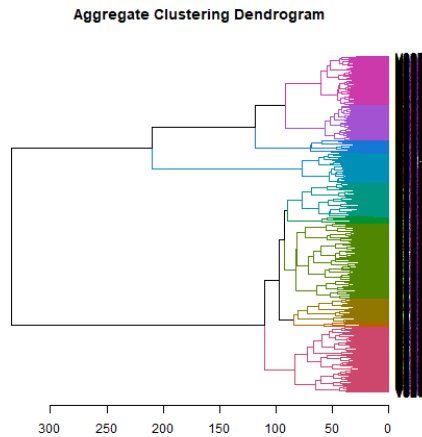


Figure 8: Aggregated Specials Dendrogram

The next model utilizes regression to create different plots showcasing the relationship between the sentiment of the vernacular leading up to a laugh and the duration of laughter. This relationship will in theory demonstrate if the sentiment of language can dictate how hard the audience will laugh. This regression will also showcase the sentiment of the text leading up to applause. Comedians have expressed their disdain for applause during their sets as they believe the audience is compensating for their lack of laughter. Laughter can be characterized as an involuntary and uncontrollable reaction to a given situation, whereas applause tends to be a conscious decision. Of course, it is not appropriate to generalize applause as a negative, but it begs the question, do certain reactions from the audience stem from the sentiment of the jokes being told?

The stereotypical structure of a joke is set-up and punchline; however, the structure of jokes vary wildly for different comedians. Comedians may even employ various styles within their own sets. This paper unfortunately does not suggest any groundbreaking techniques that can digitize a special into jokes. Instead, the regression will extract the text before the audience reaction and split it in half denoting the set-up and punchline respectively.

The regression analysis revealed intriguing insights. For laughter, the model indicates a minimal effect of LaughTime on the Average Sentiment, with a non-significant coefficient ( $t$  value = -0.144,  $p$ -value = 0.886). The overall model's explanatory power is low, as indicated by the low Multiple R-squared (5.918e-05) and Adjusted R-squared (-0.002814). This suggests that the sentiment of the language leading up to a laugh might not strongly predict the intensity or duration of laughter.

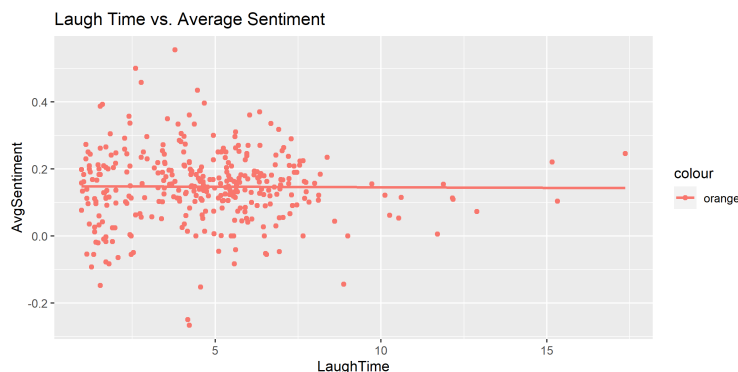


Figure 9: Aggregate Laughter Regression

Similarly, the regression for applause showed that the ApplauseTime has a non-significant effect on the Average Sentiment (coefficient = 0.002023,  $t$  value = 1.238,  $p$ -value = 0.216). The Multiple R-squared (0.001569) and Adjusted R-squared (0.0005447) also indicate a low overall explanatory power of the model. These findings suggest that the sentiment of the text before applause might not significantly influence the audience's conscious decision to applaud.

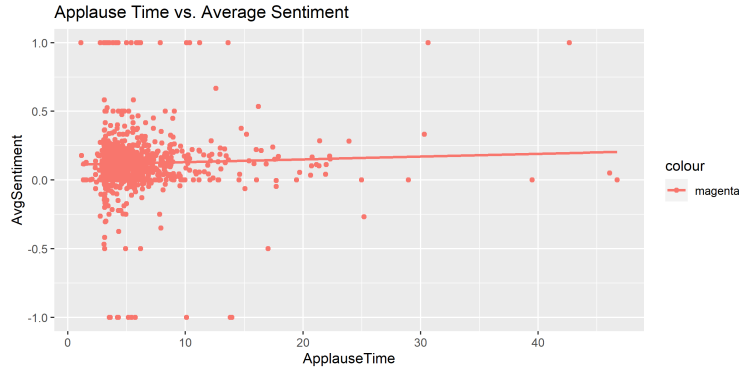


Figure 10: Aggregate Applause Regression

It's important to acknowledge the complexity of comedy dynamics since laughter is often considered an involuntary response, while applause is more deliberate. The stereotypical joke structure of set-up and punchline introduces another layer of nuance, although the regression approach simplifies this by extracting and analyzing the text before audience reactions. The outliers in the applause regression suggest a line different from laughter, but the data points show extremely similar values. This makes it complicated to distinguish between what types of jokes would lead to laughter as opposed to applause. This adds another layer of complexity to the art form and makes the art form that much more unique.

## Conclusions and Discussion

This project has shed light on various aspects of the comedic landscape. Through sentiment analysis, word frequency analysis, clustering techniques, and regression models, the report explored the dynamics of humor, linguistic patterns, and audience reactions across a diverse range of stand-up specials. Sentiment analysis revealed the predominantly positive nature of comedic language, emphasizing the role of humor in fostering a positive and enjoyable atmosphere. Word frequency analysis unveiled the significance of certain words, such as "like," providing insights into the stylistic choices and relatability that resonate with audiences. Hierarchical clustering and Latent Dirichlet Allocation (LDA) brought forth a nuanced understanding of the diverse topics covered by comedians. The lack of substantial similarity across topics emphasized the individuality and uniqueness of each comedian's comedic style, challenging preconceptions about the necessity of specific subjects for humor. The analysis of laughter and applause dynamics, while demonstrating limited predictive power, showcased the complexity of audience responses and the multifaceted nature of comedic delivery. Aspiring comedians can take solace in the data-driven revelation that there is no predetermined path to success in stand-up; instead, the diverse landscape of comedic expression welcomes experimentation and authenticity of any kind.

This project offered an interesting opportunity to explore an avenue that seems to be untouched by professionals. The analytical and language processing skills necessary to fully realize quantifiable models for stand-up are far above my ability. When you watch a lot of stand-up you start to recognize patterns, so I thought if I could as a human then surely I could train a computer to do so. Next time I would do my best to look more into language processing since that's what the paper was mostly concerned about. I was hoping to use OpenAI's ChatGPT to extract topics, and that seems like it would be a valuable tool for future work. This project was a motivating and exciting way to apply what we learned in class to a core passion of mine. In the time between graduation and finding a job, I plan to continue working on finding reliable models that can characterize stand-up.

# Bibliography

- [1] Google for Developers. “Channels: list — YouTube Data API.” Available: <https://developers-dot-devsite-v2-prod.appspot.com/youtube/v3/docs/channels/list>. [Accessed: Dec. 10, 2023]
- [2] Udacity Team. “Natural Language Processing With R.” Udacity, Oct. 08, 2020. Available: <https://www.udacity.com/blog/2020/10/natural-language-processing-with-r.html>
- [3] X. Wang, Y. Ming, T. Wu, H. Zeng, Y. Wang, and H. Qu. “DeHumor: Visual Analytics for Decomposing Humor.” Available: <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9488285>
- [4] R. Kulshrestha. “Latent Dirichlet Allocation(LDA).” Medium, Jul. 03, 2020. Available: <https://towardsdatascience.com/latent-dirichlet-allocation-lda-9d1cd064ffa2>
- [5] sicss.io. “Basic Text Analysis in R.” Available: [https://sicss.io/2020/materials/day3-text-analysis/basic-text-analysis/rmarkdown/Basic\\_Text\\_Analysis\\_in\\_R.html](https://sicss.io/2020/materials/day3-text-analysis/basic-text-analysis/rmarkdown/Basic_Text_Analysis_in_R.html). [Accessed: Dec. 10, 2023]
- [6] J. Depoix. “youtube-transcript-api: This is a Python API which allows you to get the transcripts/subtitles for a given YouTube video. It also works for automatically generated subtitles, supports translating subtitles and it does not require a headless browser, like other Selenium based solutions do!” PyPI. Available: <https://pypi.org/project/youtube-transcript-api/>
- [7] T. Munasinghe. “Labs: Trees, Hierarchical Clustering, Heatmaps.” Rensselaer Polytechnic Institute, Oct. 2023.

# Code Appendix

## Python Code: Transcript API

```
1 from youtube_transcript_api import YouTubeTranscriptApi
2 import csv
3
4 # READS THROUGH THE LIST OF COMICS, GETS THEIR YOUTUBE SPECIAL INFO
5 def getSpecials(id_list):
6     # Define the path to your CSV file
7     csv_file_path = 'ComicsList.csv'
8
9     # Specify the target column name (e.g., 'SPECIAL')
10    target_column_name = 'SPECIAL'
11
12    # Open the CSV file for reading
13    with open(csv_file_path, 'r', newline='') as csv_file:
14        # Create a CSV reader
15        csv_reader = csv.DictReader(csv_file)
16
17        # Iterate through the rows in the CSV file
18        for row in csv_reader:
19            # Get the content of the 'SPECIAL' column
20            specials_cell_content = row[target_column_name]
21            # Check if the 'SPECIAL' cell is not empty
22            if specials_cell_content:
23                # Parse the pseudo-JSON content within the 'SPECIAL' cell
24                # Extract the 'id' value using string manipulation
25                id_start = specials_cell_content.find('id:') # Find the position of 'id:'
26                if id_start != -1:
27                    id_start += len('id:') # Move the starting position after 'id:'
28                    id_end = specials_cell_content.find('}', id_start) # Find the end of the 'id' value
29                    if id_end != -1:
30                        id_value = specials_cell_content[id_start:id_end]
31                        id_list.append(id_value)
32
33    # PROCESS AND SAVE THE TRANSCRIPT TO A CSV WHERE THE FILENAME IS THE VIDEO ID
34    def saveTranscriptToCSV(id):
35        # Retrieve the available transcripts
36        transcript_list = YouTubeTranscriptApi.list_transcripts(id)
37
38        # Create a list to store the transcript data
39        transcript_data = []
40
41        # Iterate over all available transcripts
42        for transcript in transcript_list:
43            # Fetch the transcript data
44            transcript_text = transcript.fetch()
45
46            # Append each entry to the transcript_data list
47            transcript_data.extend(transcript_text)
48
49        # Define the CSV file name
50        csv_file_name = id+'.csv'
51
52        # Open the CSV file in write mode
53        with open(csv_file_name, 'w', newline='', encoding='utf-8') as csv_file:
54            # Create a CSV writer
55            csv_writer = csv.writer(csv_file)
56
57            # Write the header row
58            csv_writer.writerow(['text', 'start', 'duration'])
59
60            # Write data from the transcript data to the CSV
61            for entry in transcript_data:
62                csv_writer.writerow([entry['text'], entry['start'], entry['duration']])
63
64            print(f'CSV data has been saved to {csv_file_name}')
65
66    id_list = []
67    getSpecials(id_list)
68    for id in id_list:
69        saveTranscriptToCSV(id)
```

## R Code: Text EDA

```
1 # Load the ggplot2 library
2 library(ggplot2)
3 library(SentimentAnalysis)
4 library(tm)
5 library(NLP)
6 library(gutenbergr)
7 library(dplyr)
8 library(quanteda)
9 library(readr)
10 library(wordcloud)
11
```



```

12 fileNames <- list()
13 allData <- list()
14 wordDict <- c()
15 cuss <- "\\[_ _\\]"
16 soundEffectsAll <- c("[applause]"=0, "[music]"=0, "[laughter]"=0, "[cuss]"=0)
17
18 folder <- "C:/Users/heffec/Desktop/Semester_8/Data_Analytics/FinalProject"
19 # Get a list of CSV files in the directory
20 csv_files <- list.files(
21   path = "C:/Users/heffec/Desktop/Semester_8/Data_Analytics/FinalProject",
22   pattern = "\\*.csv$",
23   full.names = TRUE,
24 )
25 csv_files <- setdiff(csv_files, "C:/Users/heffec/Desktop/Semester_8/Data_Analytics/FinalProject/ComicsList.csv")
26 # ---- Data Preprocessing ---- #
27
28 # Iterate over each CSV file
29 for (file_path in csv_files) {
30   print(file_path)
31
32   soundEffects <- c("[applause]"=0, "[music]"=0, "[laughter]"=0, "[cuss]"=0)
33
34
35   # Read and clean CSV text
36   #print(file_path)
37   data <- read.csv(file_path)
38   data$text <- tolower(data$text)
39   data <- na.omit(data, cols = "text")
40
41   # Clean up file names to extract the specific dataset name
42   split_fileName <- strsplit(file_path, '/')
43   elements <- unlist(split_fileName)
44   last = length(elements)
45   modified_string <- substr(elements[last], 1, nchar(elements[last]) - 4)
46
47
48   # --- Audience Behavior ---
49   # This will record laughter, applause, and music from the set
50   # Saved on an individual basis and in aggregate
51   explicitUTFCode <- c(91,160,95,95,160,93) # [ _ _ ] expressed as explicit
52   for (line in data$text) {
53     for (words in line) {
54       for (words in strsplit(words, "_")){
55         for (word in words) {
56           if (word == "[applause]") {
57             soundEffects[word] = soundEffects[word] + 1
58             soundEffectsAll[word] = soundEffectsAll[word] + 1
59           }
60           # music is only present in the beginning and end of a set
61           # maybe once in between. Visual inspection of these sets proved that
62           # music is often confused with laughter in YouTubes captioning system
63           else if (word == "[music]") {
64             if (!is.na(soundEffects["music"]) && soundEffects["music"] > 3) {
65               soundEffectsAll["[laughter]"] <- soundEffectsAll["[laughter]"] + 1
66             }
67           }
68           else if (word == "[laughter]") {
69             soundEffects[word] = soundEffects[word] + 1
70             soundEffectsAll[word] = soundEffectsAll[word] + 1
71           }
72         }
73       }
74       else {
75         utf_codes <- utf8ToInt(word)
76         if (length(utf_codes) == 6 && all(utf_codes == explicitUTFCode)) {
77           soundEffects["[cuss]"] <- soundEffects["[cuss]"] + 1
78           soundEffectsAll["[cuss]"] <- soundEffectsAll["[cuss]"] + 1
79         }
80       }
81     }
82   }
83 }
84
85 print(soundEffects["[applause]"])
86 print(soundEffects["[laughter]"])
87 print(soundEffects["[cuss]"])
88 print("")
89
90 data$text <- gsub("[:punct:]", "", data$text)
91 data$text <- removeWords(data$text, stopwords("english"))
92 wordDict <- c(wordDict, data$text)
93
94 # --- Sentiment Analysis --- #
95 # Assuming file_path is a variable containing the file path
96 hist_path <- paste0(folder, modified_string, "_hist.png")
97 # Obtain SentimentGI scores
98 sentiment <- analyzeSentiment(data$text)
99 sentimentGI_scores <- sentiment$SentimentGI
100 # Determine breaks for the histogram
101 breaks <- "sturges"
102 # Set the width and height of the plotting device to make the image smaller
103 png(file = hist_path, width = 800, height = 600)

```

```

104 # Create a histogram using the SentimentGI scores
105 h <- hist(sentimentGI_scores,
106           breaks = breaks,
107           main = "SentimentGI_Histogram",
108           xlab = "Sentiment_Score",
109           ylab = "Frequency")
110
111 # Customize x-axis ticks and labels
112 axis(side = 1, at = pretty(sentimentGI_scores, n = 10), labels = pretty(sentimentGI_scores, n = 10))
113 # Adjust text labels to be smaller and inside the bars
114 text(h$mids, h$counts, labels = h$counts, adj = c(0.5, -0.5), cex = 0.7)
115 # Close the plotting device
116 dev.off()
117 }
118
119 allDtm <- DocumentTermMatrix(Corpus(VectorSource(wordDict)))
120 # Convert the document-term matrix to a data frame
121 allDtm_df <- as.data.frame(as.matrix(allDtm))
122
123 # Visualize word frequencies
124 word_freq_all <- colSums(allDtm_df)
125 word_freq_all <- sort(word_freq_all, decreasing = TRUE)
126
127 barplot(word_freq_all[1:5], col = "skyblue", main = "Word_Frequencies_Without_Stop_Words")
128 all_wordFreq_path = paste0(folder, "/all_wordFreq.png")
129 ggsave(all_wordFreq_path, plot = last_plot(), units = "px", device = "png", width = 800, height = 600)
130
131 # Generate a word cloud
132 wordcloud(words = names(word_freq_all), freq = word_freq_all, scale=c(3,0.5), min.freq = 1, colors = brewer.pal(8, "Dark2"))
133
134 # --- Overall Sentiment Analysis --- #
135 # Assuming file_path is a variable containing the file path
136 hist_path <- paste0(folder, "/all_hist.png")
137 # Obtain SentimentGI scores
138 sentiment <- analyzeSentiment(wordDict)
139 sentimentGI_scores <- sentiment$SentimentGI
140 # Determine breaks for the histogram
141 breaks <- "sturges"
142 # Set the width and height of the plotting device to make the image smaller
143 png(file = hist_path, width = 800, height = 600)
144 # Create a histogram using the SentimentGI scores
145 h <- hist(sentimentGI_scores,
146           breaks = breaks,
147           main = "SentimentGI_Histogram",
148           xlab = "Sentiment_Score",
149           ylab = "Frequency")
150
151 # Customize x-axis ticks and labels
152 axis(side = 1, at = pretty(sentimentGI_scores, n = 10), labels = pretty(sentimentGI_scores, n = 10))
153 # Adjust text labels to be smaller and inside the bars
154 text(h$mids, h$counts, labels = h$counts, adj = c(0.5, -0.5), cex = 0.7)
155 # Close the plotting device
156 dev.off()
157
158 # ---- Topic Modeling --- #
159 allDtm <- allDtm[rowSums(as.matrix(allDtm)) > 0, ]
160 lda <- LDA(allDtm, method = "Gibbs", k=10)
161 terms_df <- as.data.frame(terms(lda, 10))
162
163 # Save the list of data frames for further analysis if needed
164 save(allData, file = "allData.RData")
165
166 as.data.frame(wordDict)

```

## R Code: Timing EDA

```

1 library(ggplot2)
2
3 # List to store individual plots
4 individual_plots <- list()
5
6 # Get a list of CSV files in the directory
7 csv_files <- list.files(
8   path = "C:/Users/heffec/Desktop/Semester_8/Data_Analytics/FinalProject",
9   pattern = "\\\\.csv$",
10  full.names = TRUE,
11 )
12 csv_files <- setdiff(csv_files, "C:/Users/heffec/Desktop/Semester_8/Data_Analytics/FinalProject/ComicsList.csv")
13
14 # List to store spline points and outlier points for each file
15 all_spline_points <- list()
16 outlier_indices_list <- list()
17 folder = "C:/Users/heffec/Desktop/Semester_8/Data_Analytics/FinalProject"
18
19 for (file_path in csv_files) {
20   print(file_path)
21
22   # Read and clean CSV text
23   df <- read.csv(file_path)

```

```

24 df$text <- tolower(df$text)
25
26 # Check for missing values
27 if (any(is.na(df$start)) || any(is.na(df$speech_rate))) {
28   next() # Skip to the next iteration
29 }
30
31 # Calculate speech rate (words per second)
32 df$word_count <- sapply(strsplit(df$text, "\\s+"), length)
33 df$speech_rate <- df$word_count / df$duration
34
35 # Weighted average speech rate based on duration
36 weighted_avg_rate <- weighted.mean(df$speech_rate, df$duration)
37
38 # Cubic spline interpolation for the cleaned data
39 spline_points <- smooth.spline(df$start, df$speech_rate)
40 spline_df <- data.frame(x = spline_points$x, y = spline_points$y)
41 all_spline_points[[file_path]] <- spline_df
42
43 # Create a time series plot with cubic spline interpolation and LOESS smoothing
44 plot <- ggplot(df, aes(x = start, y = speech_rate)) +
45   geom_line() +
46   geom_line(data = spline_df, aes(x, y), color = "green") +
47   geom_smooth(method = "loess", se = TRUE, color = "red") +
48   labs(x = "Start Time", y = "SpeechRate(words_per_second)", title = "SpeechRateOverTimewithCubicSpline
49         InterpolationandLOESSSmoothing") +
50   scale_color_manual(values = c("ActualData", "CubicSplineInterpolation", "LocallyEstimatedScatterplotSmoothing"))
51
52 # Save individual plot as PNG
53 individual_plot_path <- paste0(file_path, "_spline_plot.png")
54 ggsave(individual_plot_path, plot, width = 2400, height = 1200, units = "px")
55
56 individual_plots[[file_path]] <- plot
57 }
58
59 # Combine all individual spline points into one big spline
60 all_spline_points_df_t <- data.table::rbindlist(all_spline_points)
61
62 # Get the first 80% of the data
63 subset_index <- 1:round(0.8 * nrow(all_spline_points_df_t))
64 subset_spline_points_df_t <- all_spline_points_df_t[subset_index, ]
65
66 # Fit a smooth spline to the subset data
67 spline_points_big_subset <- smooth.spline(subset_spline_points_df_t$x, subset_spline_points_df_t$y, spar = 0.8)
68
69 # Create a data frame from the spline points
70 spline_df_big_subset <- data.frame(x = spline_points_big_subset$x, y = spline_points_big_subset$y)
71
72 # Create a big plot with the subset spline
73 subset_plot <- ggplot(data = spline_df_big_subset, aes(x = x, y = y)) +
74   geom_line() +
75   geom_smooth(method = "loess", se = FALSE, color = "red") +
76   labs(x = "Start Time", y = "SpeechRate(words_per_second)", title = "SubsetCubicSplineInterpolation")
77
78 # Save the subset spline plot as PNG
79 subset_plot_path <- paste0(folder, "/subset_splinePlot.png")
80 ggsave(subset_plot_path, subset_plot, width = 2400, height = 1200, units = "px")

```

## R Code: Clustering Model

```

1 library(topicmodels)
2 library(tm)
3 library(dendextend)
4 library(dplyr)
5 library(cluster)
6
7 explicitUTFCode <- c(91, 160, 95, 95, 160, 93) # [ _ ] expressed as explicit
8 folder <- "C:/Users/heffec/Desktop/Semester8/DataAnalytics/FinalProject"
9
10 # Function to clean text
11 clean_text <- function(text) {
12   # Remove profanity placeholders
13   text <- gsub(paste0("\\", paste(explicitUTFCode, collapse = "\\s*")), NA, text, fixed = TRUE)
14
15   # Remove audience behavior annotations
16   text <- gsub("\\[laughter\\]", NA, text)
17   text <- gsub("\\[applause\\]", NA, text)
18   text <- gsub("\\[music\\]", NA, text)
19
20   # Convert to lowercase
21   text <- tolower(text)
22
23   return(text)
24 }
25
26 perform_LDA <- function(text, num_topics = 5) {
27   dtm <- DocumentTermMatrix(Corpus(VectorSource(text)))
28   dtm_matrix <- as.matrix(dtm)
29   dtmCleaned <- dtm_matrix[rowSums(dtm_matrix) > 0, ]
30 }

```

```

31 if (nrow(dtmCleaned) == 0) {
32   warning("Empty document-term matrix after cleaning. Skipping LDA.")
33   return(NULL)
34 }
35
36 lda_model <- LDA(dtmCleaned, k = num_topics)
37 return(lda_model)
38 }
39
40 extract_topics <- function(lda_model, threshold = 0.1) {
41   topic_terms <- terms(lda_model, 10)
42   terms_df <- as.data.frame(topic_terms)
43   return(terms_df)
44 }
45
46 allTopics <- list()
47
48 # Iterate over each CSV file
49 for (file_path in csv_files) {
50   print(file_path)
51
52   data <- read.csv(file_path)
53   data$text <- clean_text(data$text)
54   data <- na.omit(data, cols = "text")
55   data$text <- removeWords(data$text, stopwords("english"))
56   data$text <- gsub("[:punct:]", "", data$text)
57   data <- na.omit(data, cols = "text")
58
59   individualLDA <- perform_LDA(data$text)
60
61   # Check if LDA was successful
62   if (!is.null(individualLDA)) {
63     extractedTopics <- extract_topics(individualLDA)
64     allTopics[[file_path]] <- individualLDA
65
66     # Transpose the matrix so that rows represent topics
67     topic_matrix_transposed <- t(as.data.frame(individualLDA@beta))
68
69     # Calculate distances
70     dist_matrix <- dist(topic_matrix_transposed)
71
72     # Perform hierarchical clustering
73     hclust_result <- hclust(dist_matrix, method = "ward.D2")
74     print(summary(hclust_result))
75
76     # Assign cluster labels to topics
77     cluster_labels <- cutree(hclust_result, k = 10)
78
79     # Set the output PNG file path
80     file_name <- basename(file_path)
81     png_file_path <- file.path(folder, paste0(file_name, "_dendrogramManhattan", ".png"))
82
83     # Open PNG file for plotting
84     png(png_file_path)
85
86     # Plot the dendrogram for each document
87     dendrogram_ind <- as.dendrogram(hclust_result)
88     dendrogram_ind <- color_branches(dendrogram_ind, k = length(unique(cluster_labels)))
89
90     # Plot the individual dendrogram
91     plot(dendrogram_ind, main = "Clustering Dendrogram", horiz = TRUE)
92
93     # Save the PNG file
94     dev.off()
95   }
96 }
97
98 # Combine topic matrices of all documents
99 all_topic_matrices <- lapply(allTopics, function(lda_model) as.data.frame(lda_model@beta))
100 combined_topic_matrix <- dplyr::bind_rows(all_topic_matrices)
101 combined_topic_matrix_transposed <- t(combined_topic_matrix)
102 dist_matrix <- dist(combined_topic_matrix_transposed)
103 hclust_result_all <- hclust(dist_matrix, method = "ward.D2")
104 cluster_labels_all <- cutree(hclust_result_all, k = 10)
105
106 # Plot the dendrogram for all documents
107 dendrogram_all <- as.dendrogram(hclust_result_all)
108 dendrogram_all <- color_branches(dendrogram_all, k = length(unique(cluster_labels_all)))
109
110 # Plot the combined dendrogram
111 png_file_path <- file.path(folder, paste0("all_dendrogram", ".png"))
112 png(png_file_path)
113 plot(dendrogram_all, main = "Aggregate Clustering Dendrogram", horiz = TRUE)
114 dev.off()
115

```

## R Code: Regression Model

```

1 # Load the required libraries
2 library(SentimentAnalysis)

```

```

3 library(ggplot2)
4
5 sets_with_laughter <- list()
6 sets_with_applause <- list()
7
8 # Iterate over each CSV file
9 for (file_path in csv_files) {
10
11   print(file_path)
12
13   soundEffects <- c("[applause]"=0, "[music]"=0, "[laughter]"=0, "[cuss]"=0)
14
15   # Read and clean CSV text
16   data <- read.csv(file_path)
17   data$text <- tolower(data$text)
18   explicitUTFCode <- c(91, 160, 95, 95, 160, 93) # [ __ ] expressed as explicit
19   data$text <- gsub(paste0("\\s+", paste(explicitUTFCode, collapse = "\\s*")), "", data$text, fixed = TRUE)
20   data$text <- removeWords(data$text, stopwords("english"))
21   data$text <- na.omit(data$text, cols = "text")
22
23   # Find indices where text is equal to "\\[[laughter\\]"
24   laughter_indices <- grep("\\[[laughter\\]", data$text)
25   applause_indices <- grep("\\[[applause\\]", data$text)
26
27   data$text <- gsub("[[:punct:]]", "", data$text)
28
29   # Check if there are more than 0 laughter indices
30   if (length(laughter_indices) > 1) {
31     prev <- 1 # Start from the first index
32     avg_sentiments <- numeric() # Vector to store average sentiments
33     laugh_times <- numeric() # Vector to store laugh times
34
35     for (index in laughter_indices) {
36       joke <- data$text[prev:(index - 1)]
37       sentiment <- analyzeSentiment(joke)
38       sentimentGI_scores <- sentiment$SentimentGI
39       sentimentGI_scores <- na.omit(sentimentGI_scores)
40       avg_sentiment <- mean(sentimentGI_scores)
41       avg_sentiments <- c(avg_sentiments, avg_sentiment)
42
43       laugh_time <- data$duration[index]
44       laugh_times <- c(laugh_times, laugh_time)
45       # Check if sentiment analysis result is valid
46       prev <- index + 1
47     }
48
49     # Create a data frame for each set with laughter
50     set_data <- data.frame(LaughTime = laugh_times, AvgSentiment = avg_sentiments)
51     sets_with_laughter[[file_path]] <- set_data
52   }
53
54   # Check if there are more than 0 laughter indices
55   if (length(applause_indices) > 1) {
56     prev <- 1 # Start from the first index
57     avg_sentiments <- numeric() # Vector to store average sentiments
58     applause_times <- numeric() # Vector to store laugh times
59
60     for (index in applause_indices) {
61       joke <- data$text[prev:(index - 1)]
62       sentiment <- analyzeSentiment(joke)
63       sentimentGI_scores <- sentiment$SentimentGI
64       sentimentGI_scores <- na.omit(sentimentGI_scores)
65       avg_sentiment <- mean(sentimentGI_scores)
66       avg_sentiments <- c(avg_sentiments, avg_sentiment)
67
68       applause_time <- data$duration[index]
69       applause_times <- c(applause_times, applause_time)
70       # Check if sentiment analysis result is valid
71       prev <- index + 1
72     }
73
74     # Create a data frame for each set with laughter
75     set_data <- data.frame(ApplauseTime = applause_times, AvgSentiment = avg_sentiments)
76     sets_with_applause[[file_path]] <- set_data
77   }
78 }
79
80 combined_data <- bind_rows(sets_with_laughter)
81 combined_data_applause <- bind_rows(sets_with_applause)
82
83 # Fit linear regression model - laugh
84 lm_model_laugh <- lm(AvgSentiment ~ LaughTime, data = combined_data)
85 print(summary(lm_model_laugh))
86 laugh_graph <- ggplot(combined_data, aes(x = LaughTime, y = AvgSentiment, color = "orange")) +
87   geom_point() +
88   geom_smooth(method = "lm", se = FALSE) +
89   labs(title = "LaughTime vs. Average Sentiment")
90
91 # Fit linear regression model - applause
92 lm_model_applause <- lm(AvgSentiment ~ ApplauseTime, data = combined_data_applause)
93 print(summary(lm_model_applause))
94 # Corrected applause regression plot

```

```
95 app_graph <- ggplot(combined_data_applause, aes(x = ApplauseTime, y = AvgSentiment, color = "magenta")) +
96   geom_point() +
97   geom_smooth(method = "lm", se = FALSE) +
98   labs(title = "ApplauseTime vs. AverageSentiment")
99
100 # Save the corrected plot
101 path_applause <- paste0(folder, "/all_applauseRegression.png")
102 ggsave(path_applause, app_graph, width = 2400, height = 1200, units = "px")
```