

# Exploratory Data Analysis

## on

# Trending YouTube Videos

By

Venugopal Katragadda(A20544776)

Sannihitha Gudimalla (A20560248)

Sruthi Kondapalli (A20554780)

Bezawada Sai Sravanya (A20561552)



# Executive Summary



## Key research issue

This research aims to explore the dynamics behind YouTube trending videos and pinpoint the critical factors that drive their success. Through a comprehensive analysis, we investigated various aspects of trending content to gain meaningful insights into what contributes to a video's popularity on the platform.



## Findings

Our analysis provided important insights into several factors influencing trending videos, including optimal category choices, ideal title lengths, and the influence of sentiment in descriptions on viewership. We also compared our results with existing trends to offer a broader perspective on the data. However, our efforts to apply cluster analysis for deeper insights were less successful than anticipated.



## Future Work

Future work could involve a deeper examination of the identified clusters and the development of predictive models to forecast a video's likelihood of trending. Such models would leverage key input features to not only predict trending potential but also provide strategic recommendations for optimizing content to enhance its chances of achieving trending status.

# Overview of project

## Methodology



DATA COLLECTION



DATA PREPROCESSING

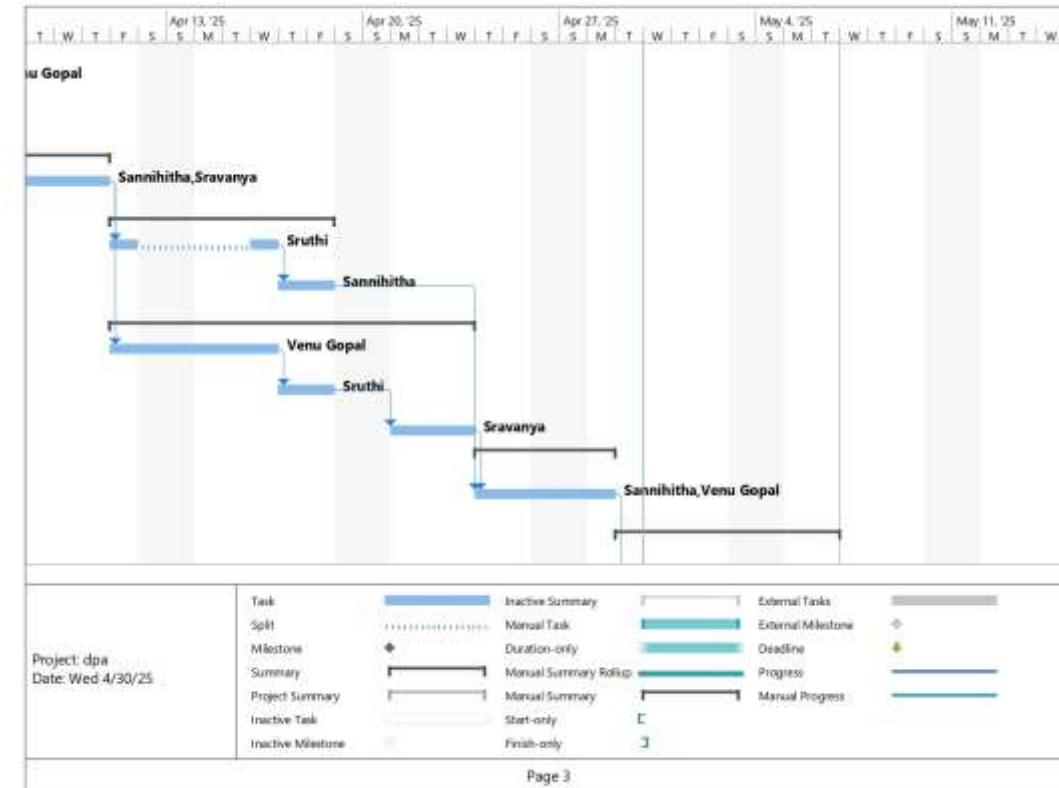
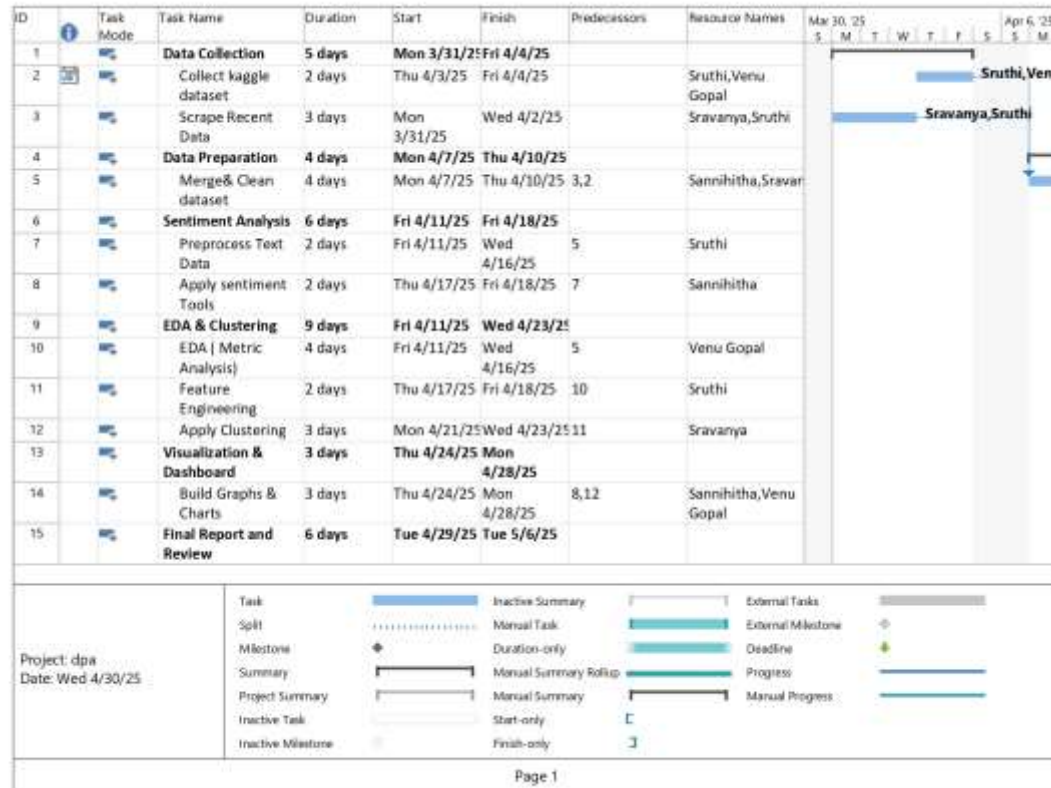


EXPLORATORY DATA  
ANALYSIS (EDA)



CLUSTER ANALYSIS

# Project Plan



## Updated Plan:

1. We made a lot more graphs and charts than we first thought. These visuals helped us see trends in the YouTube data more clearly.
2. We faced unexpected problems with the YouTube data scrapping. It took us longer to get the data we needed, but we found solutions and kept going with our analysis.

# Data Sources

## Kaggle Dataset

For the analysis we have chosen a data set which has the following characteristics.

**Dataset Type:** Multivariate

**Dataset Size:** 10 documents categorized by country, each having 16 columns and around 20000 records.

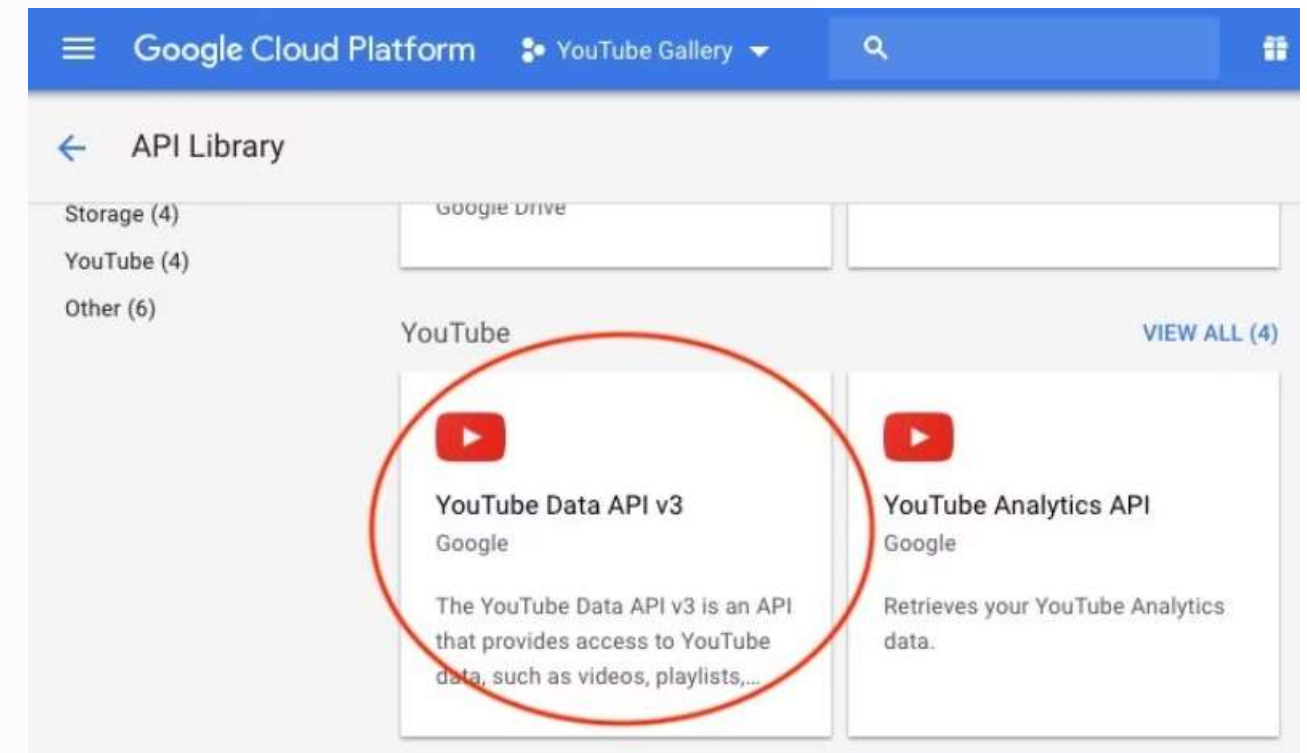
**Countries chosen:** Canada, USA, Great Britain, France, Germany, Russia, Mexico, South Korea, India, Japan

**Missing Values:** Yes

**Feature Descriptions:**

1. **video\_id:** A unique identifier for each video on YouTube. (String)
2. **title:** The title of the video. (String)
3. **trending\_date:** Date when the video appeared on the trending list. (String)
4. **channel\_title:** Name of the YouTube channel that uploaded the video. (String)
5. **category\_id:** Category to which the video belongs. (Integer)
6. **published\_time:** Date and time when the video was published on YouTube. (DateTime)
7. **tags:** Keywords or phrases that describe the content of the video. (String)
8. **views:** Number of views. (Integer)
9. **likes:** Number of likes. (Integer)
10. **dislikes:** Number of dislikes. (Integer)
11. **comment\_count:** Number of comments. (Integer)
12. **thumbnail\_link:** Link to the thumbnail image that represents the video. (String)
13. **comments\_disable:** Whether comments have been disabled for the video. (Boolean)
14. **ratings\_disable:** Whether ratings have been disabled for the video. (Boolean)
15. **video\_error\_or\_removed:** Whether there was an error with the video or if the video has been removed from YouTube. (Boolean)
16. **description:** Description of the video. (String)

## YouTube API



# Data Preprocessing

1

Converting the categorical variable `category_ID` and incorporating a new `country` variable.

2

Handling missing values and Duplicates

3

Convert the 'trending\_date' and 'published\_time' columns to datetime format

4

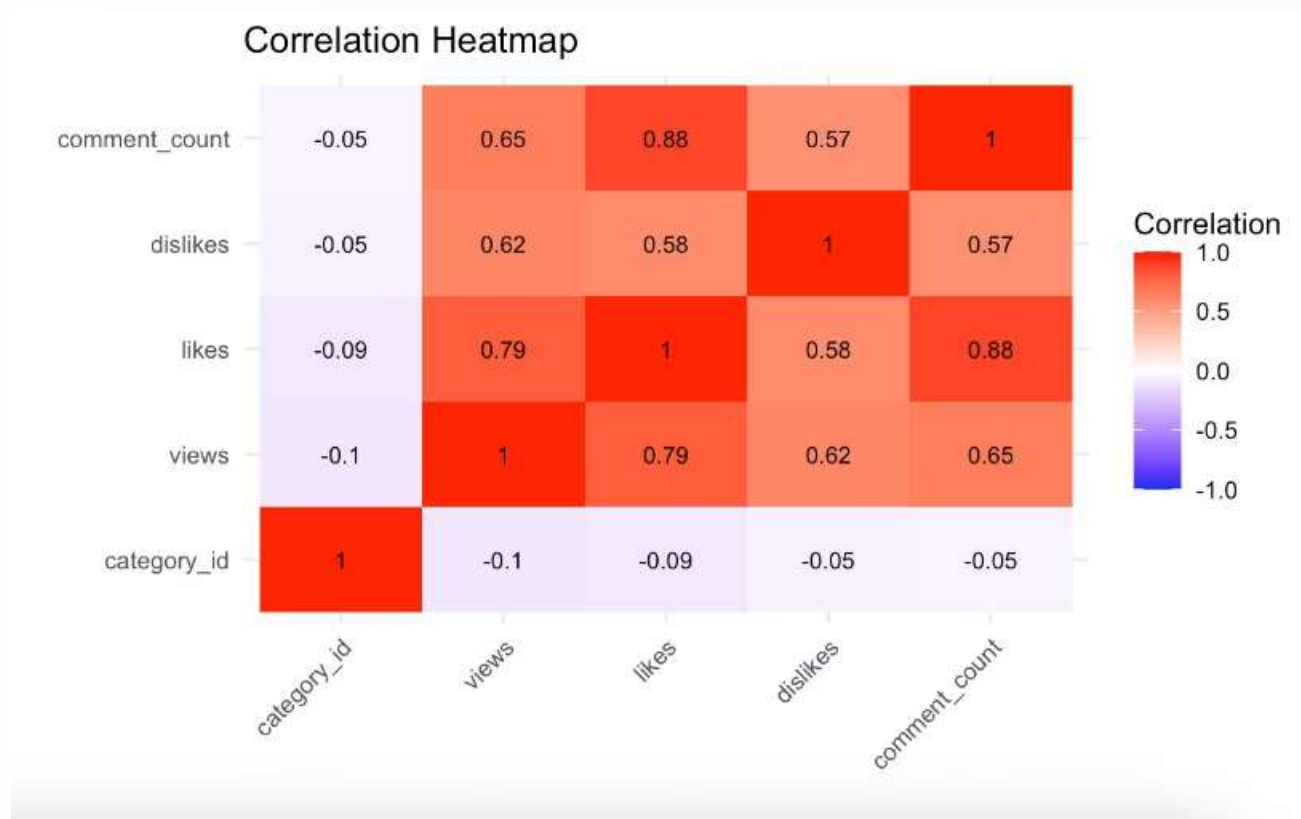
Extracting from 'published\_time' such as the day of the week or hour of the day the video was published.

5

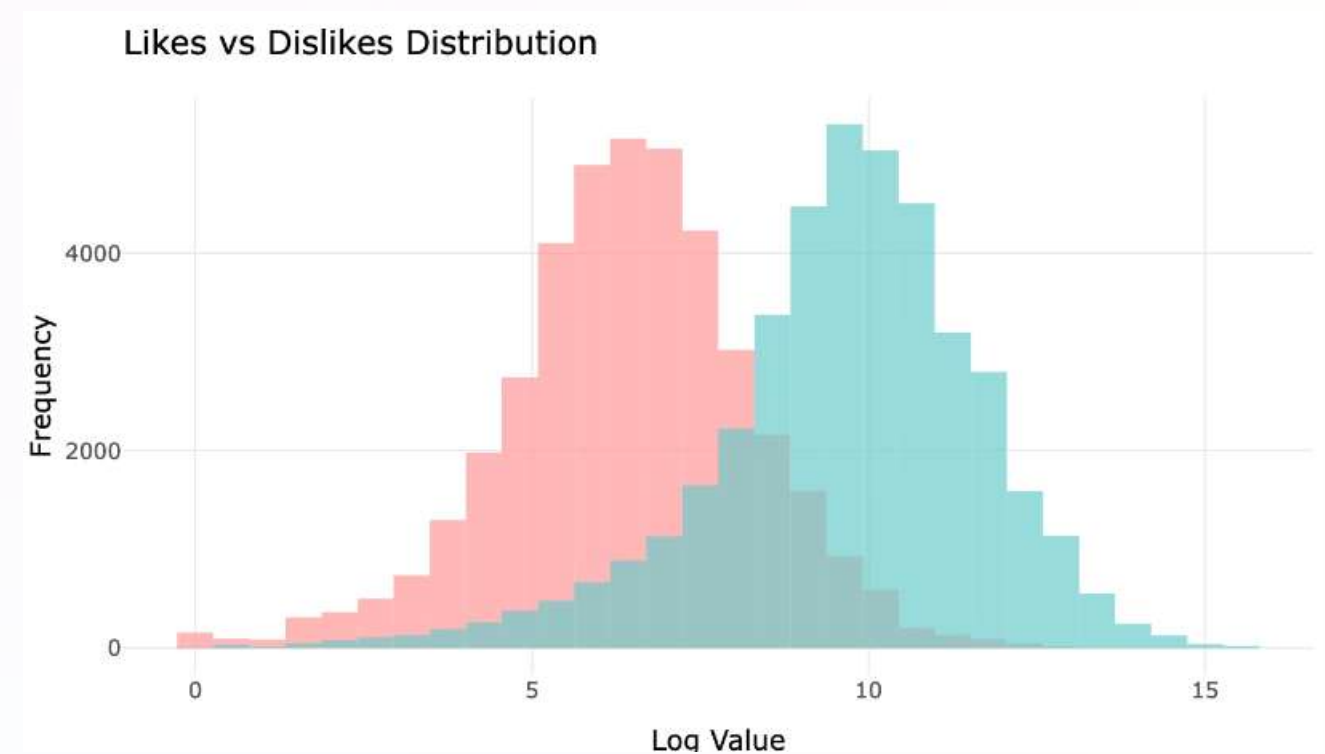
Text pre - processing for Sentiment analysis

# Exploratory Data Analysis

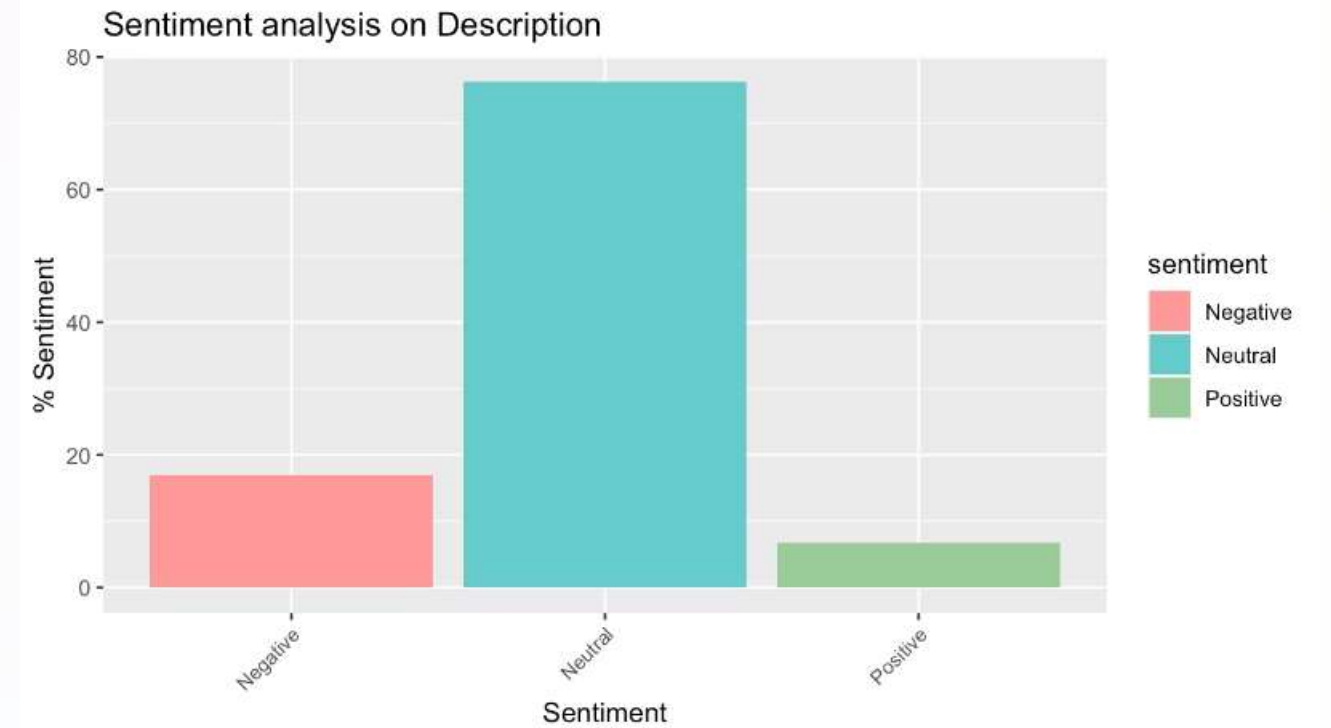
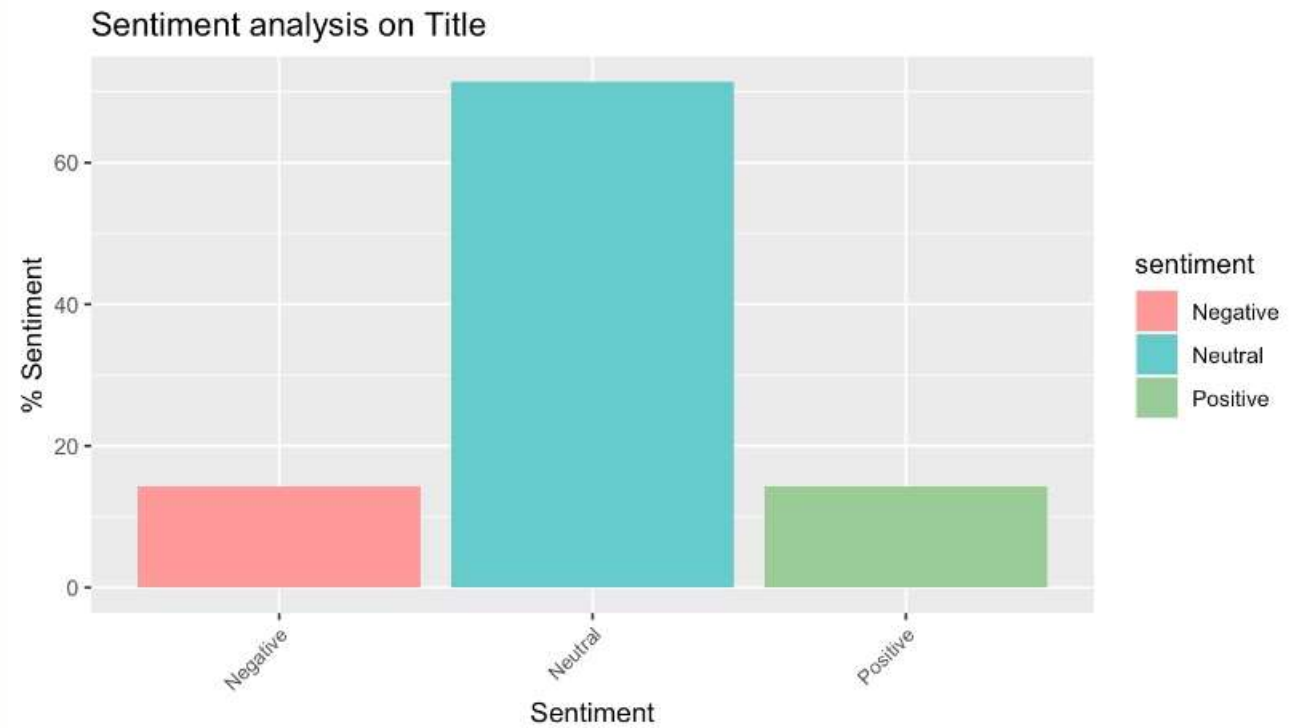
## CORRELATION MATRIX



## DISTRIBUTION PLOT

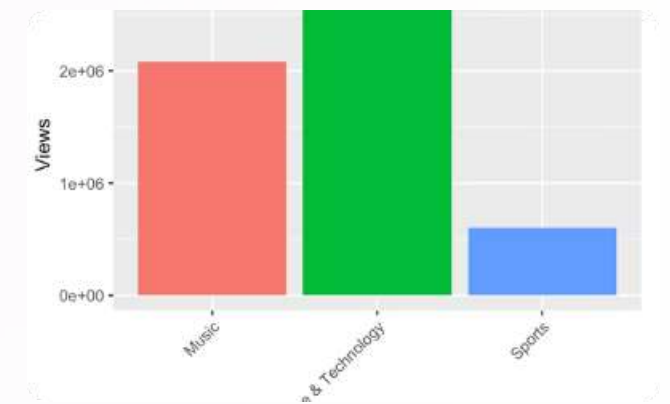
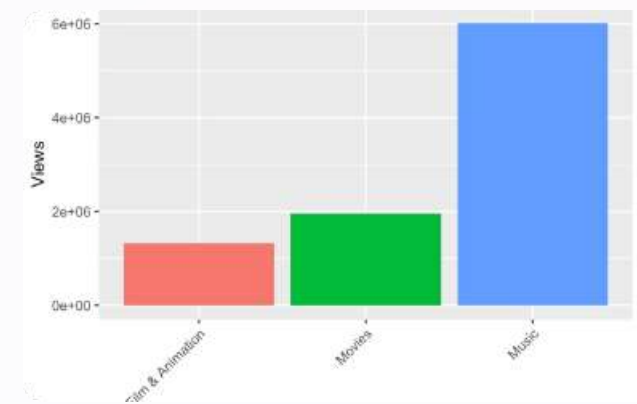
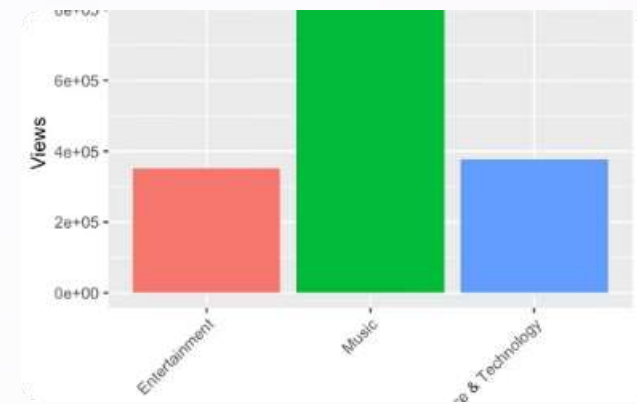
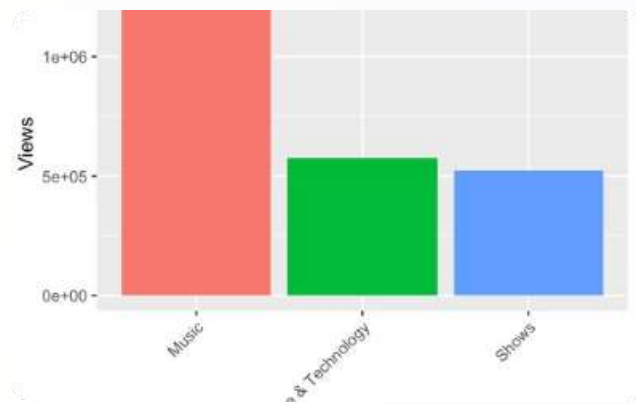


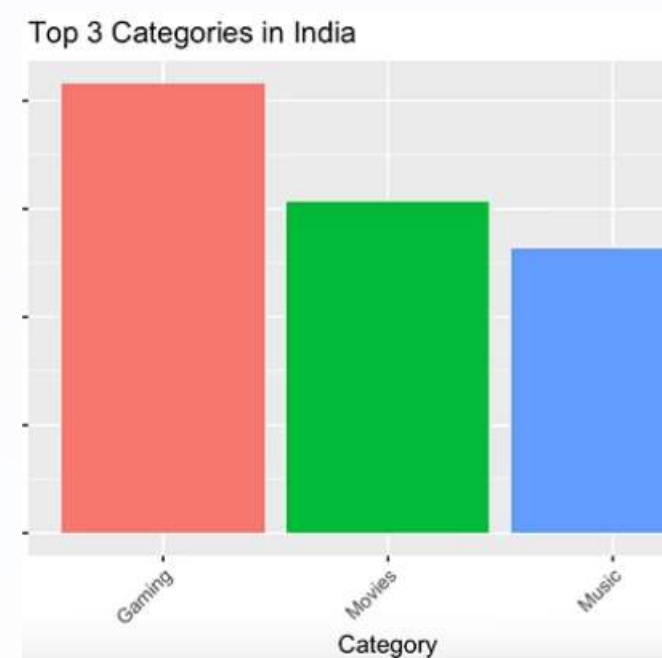
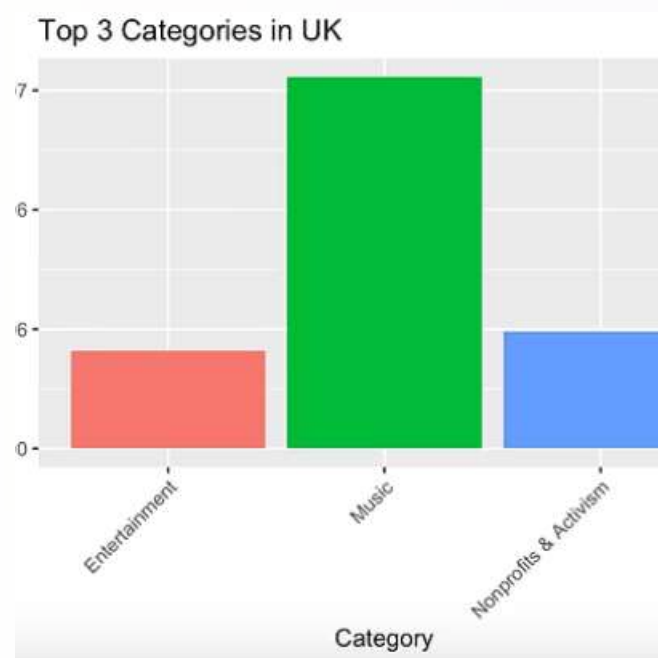
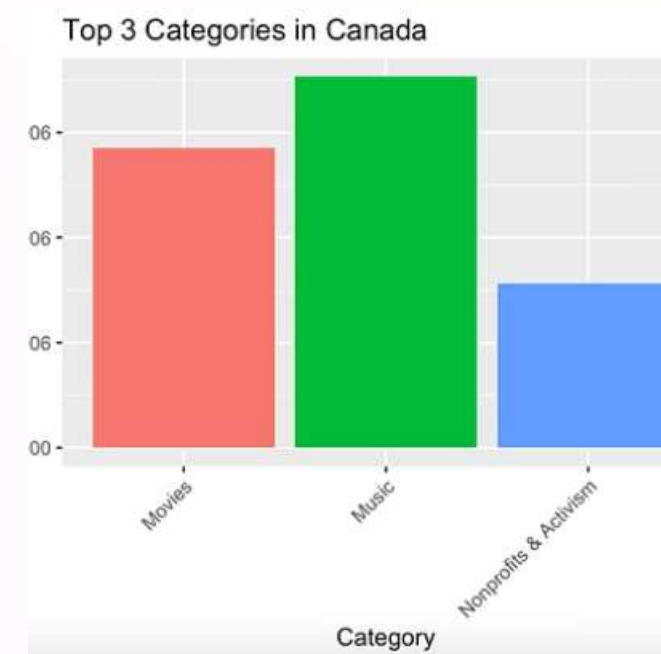
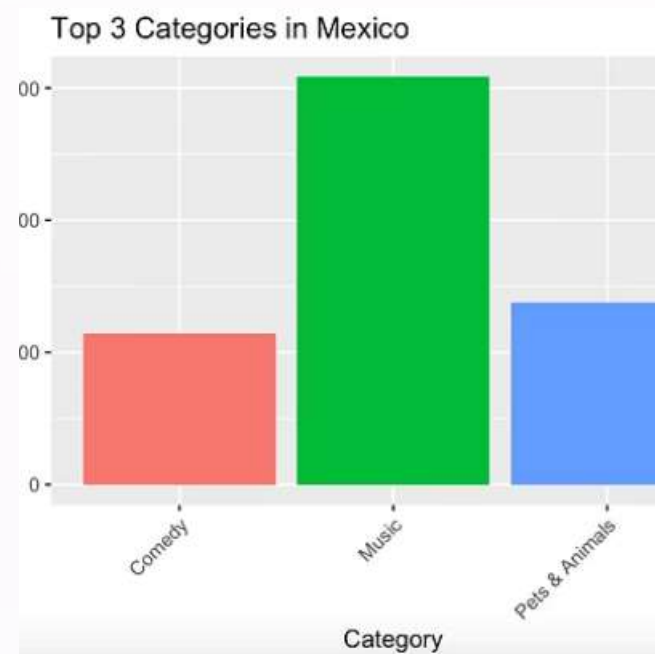
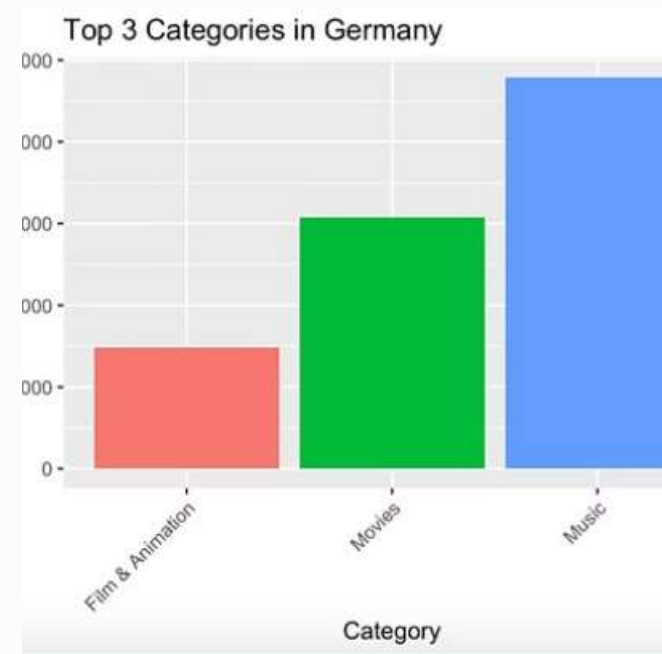
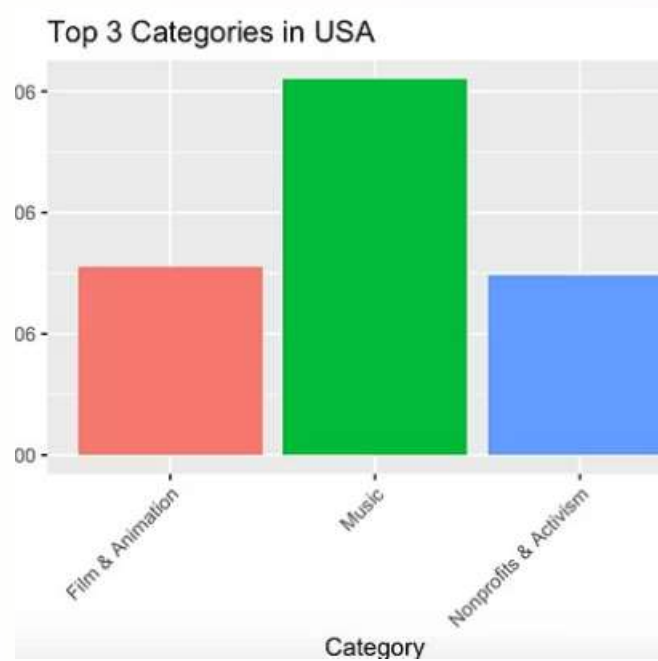
# Sentiment Analysis





WHAT ARE THE TOP 3 CATEGORIES BASED ON THE NUMBER OF VIEWS, FOR EACH OF THE COUNTRIES?





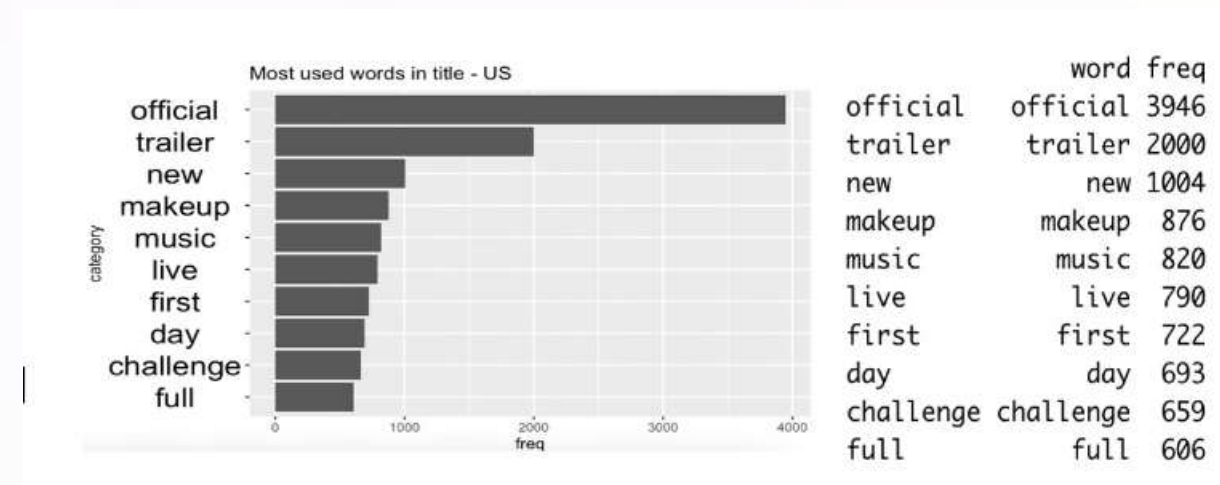
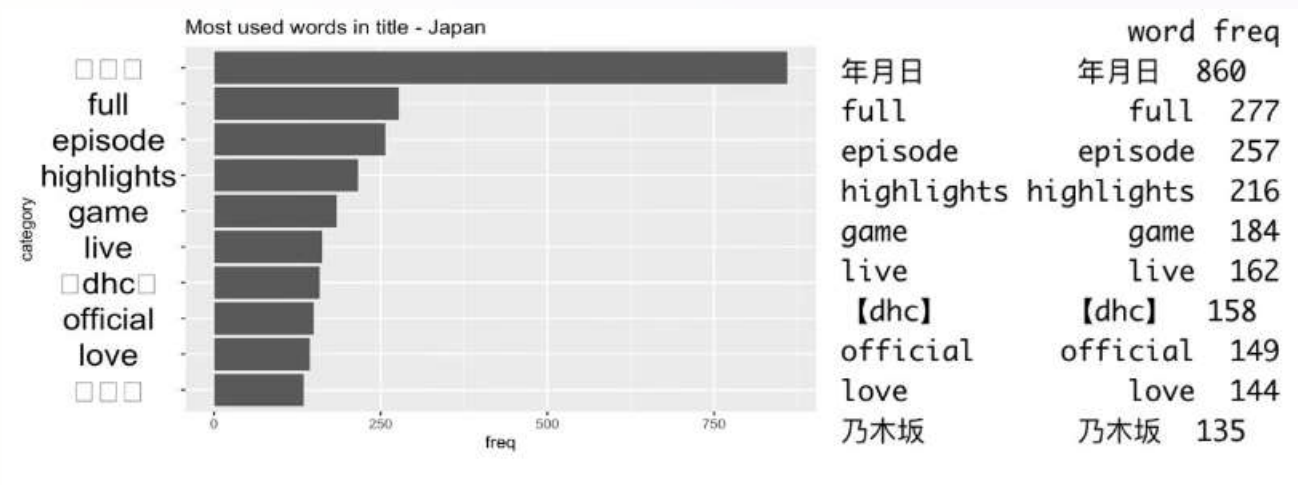
# Countries with highest least views

```
> head(IN_most_category,3)
# A tibble: 3 × 3
  category_name average_views least_views
  <chr>          <dbl>         <dbl>
1 Gaming        4162462.         72964
2 Movies        3065001.        165601
3 Music         2631116.         10971
```

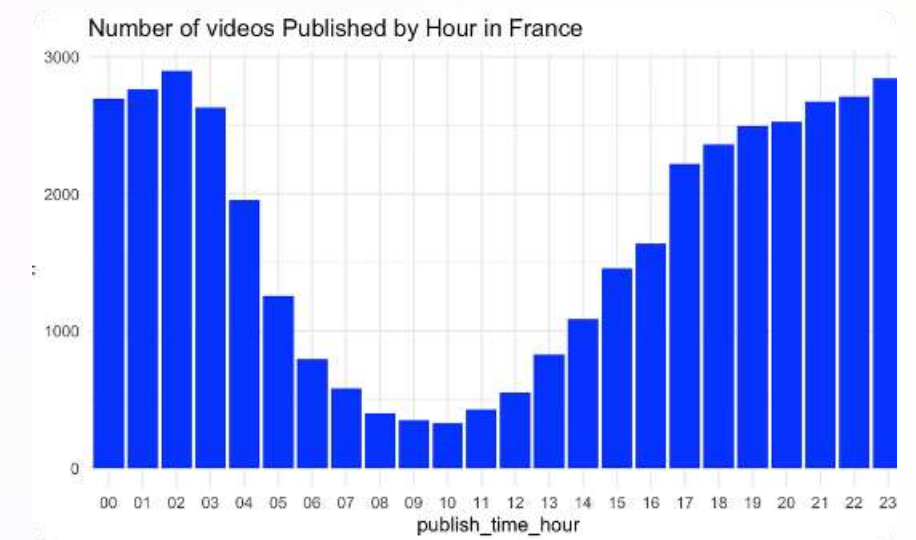
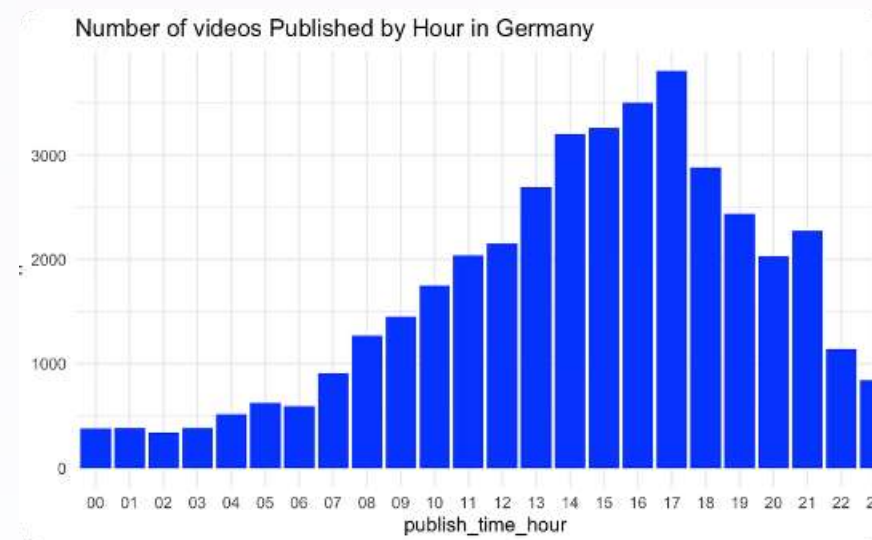
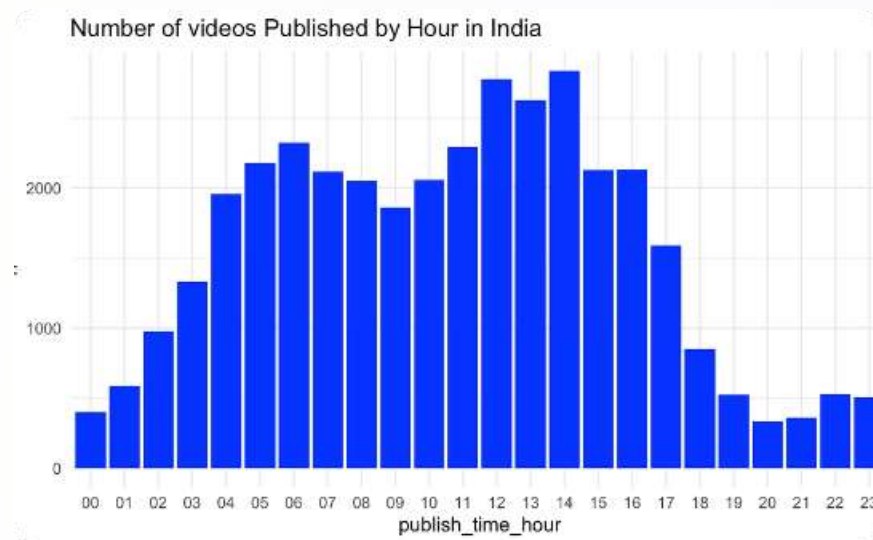
```
> head(CA_most_category,3)
# A tibble: 3 × 3
  category_name average_views least_views
  <chr>          <dbl>         <dbl>
1 Music        3532525.          3201
2 Movies       2853415         225528
3 Nonprofits & Activism 1562184.          1898
```

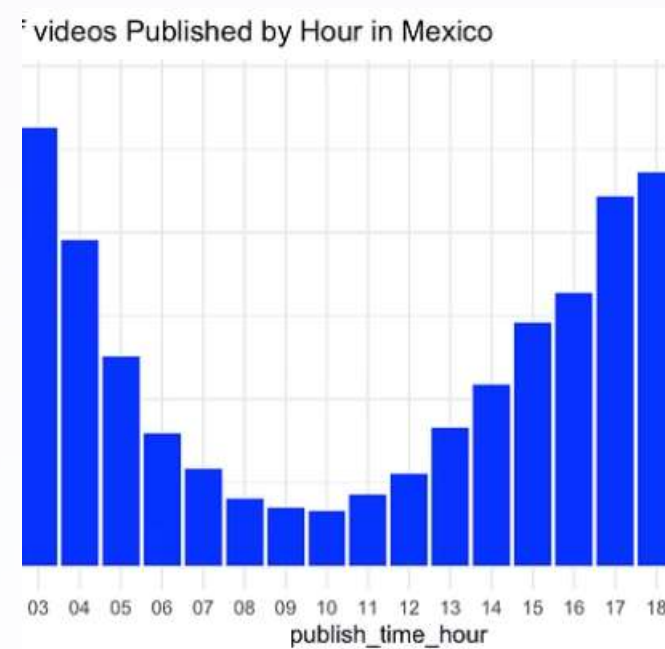
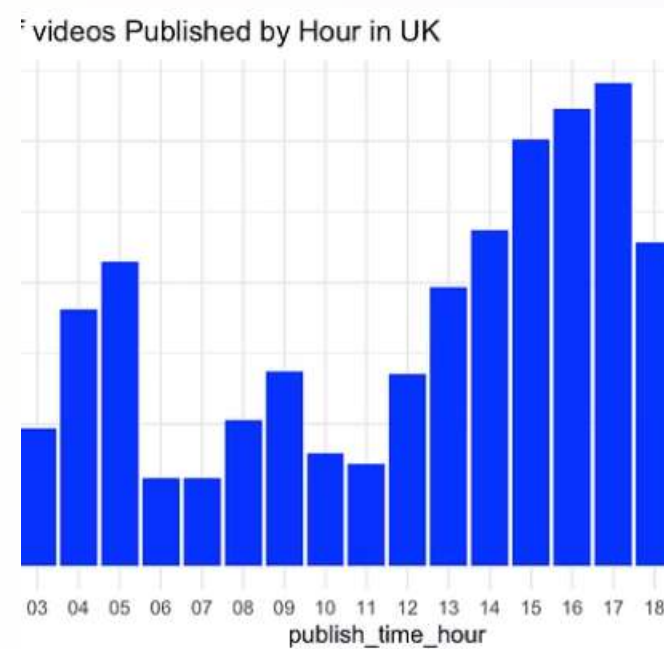
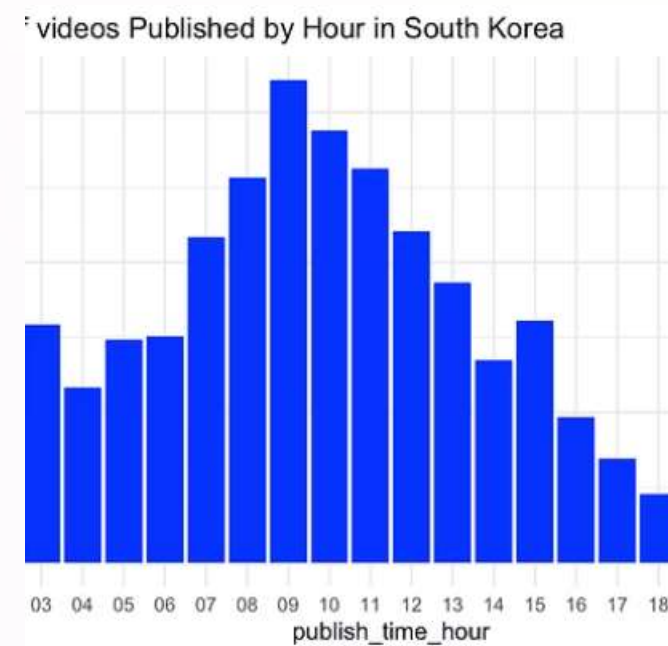
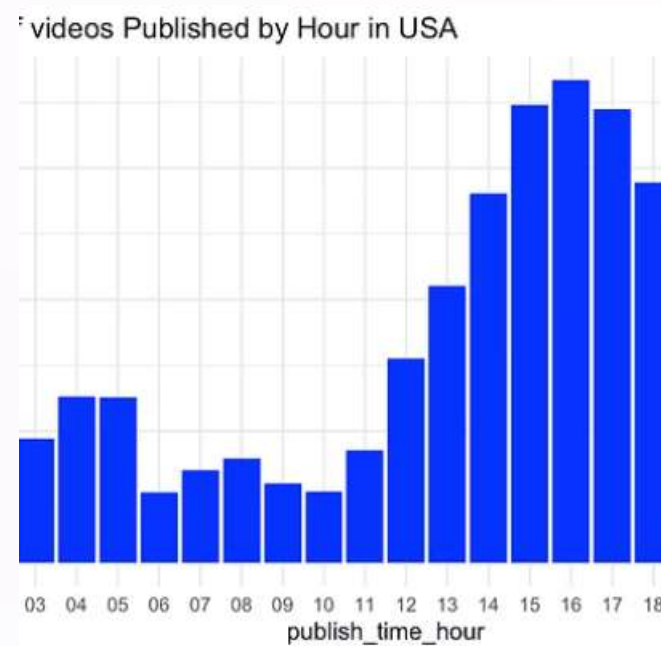
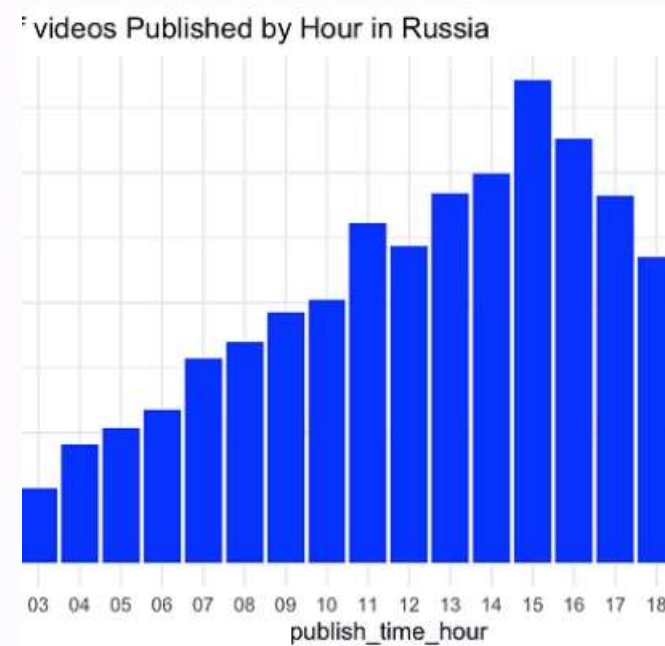
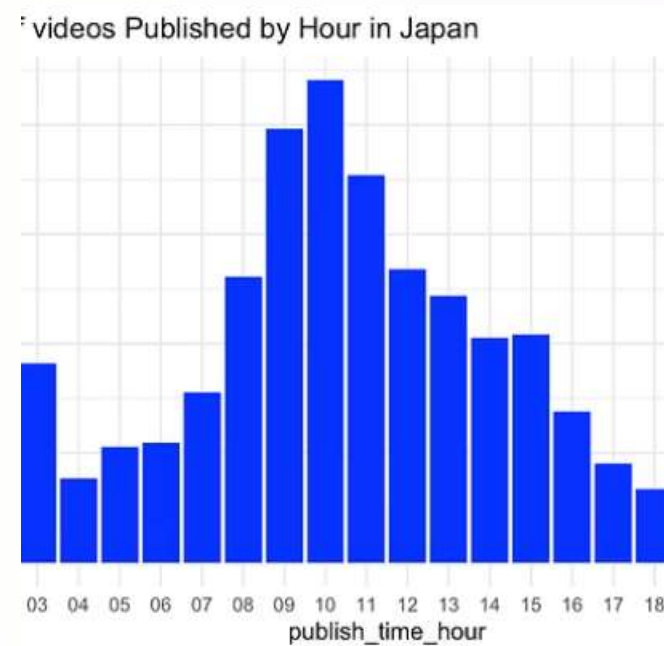
```
> head(GB_most_category,3)
# A tibble: 3 × 3
  category_name average_views least_views
  <chr>          <dbl>         <dbl>
1 Music       12444443.           2152
2 Nonprofits & Activism 3919981.          19270
3 Entertainment 3264608.           2650
```

ARE THERE SPECIFIC KEYWORDS IN THE TITLE THAT ARE MORE LIKELY TO RESULT IN A VIDEO TRENDING?

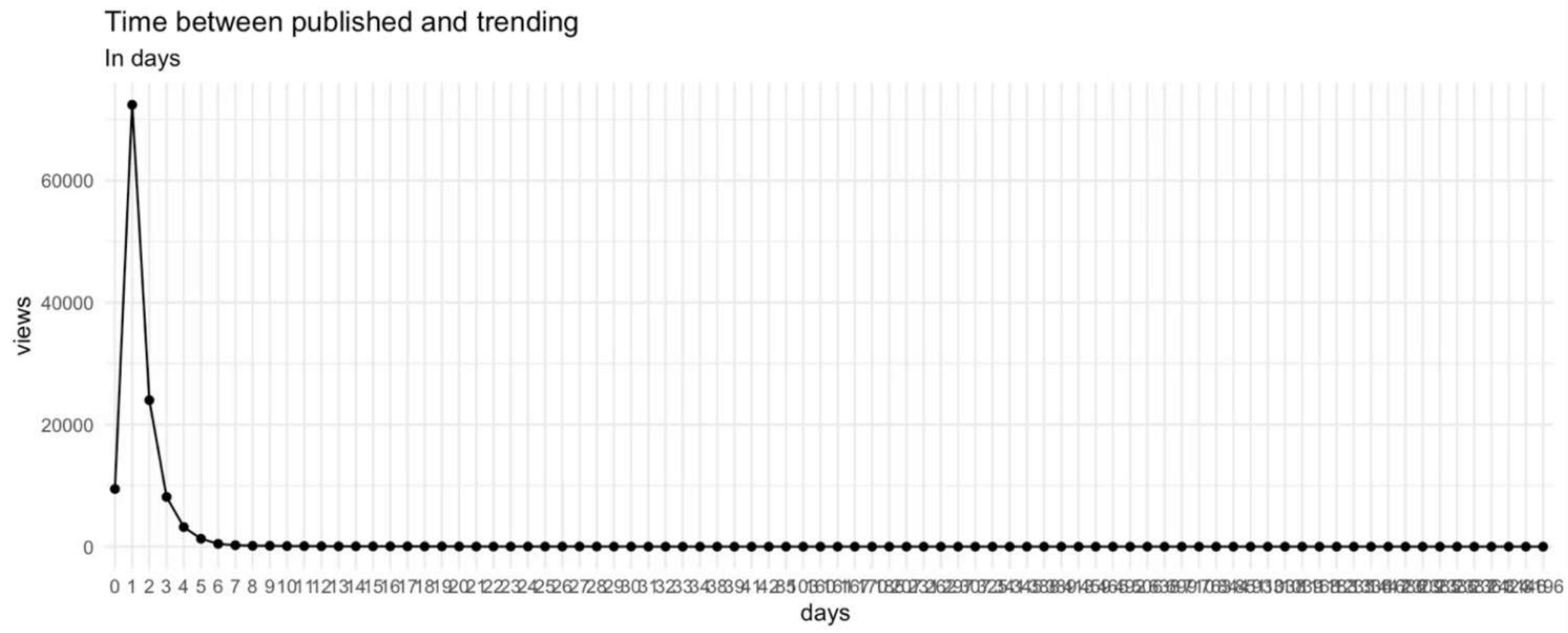


# HOUR OF THE DAY WHEN TRENDING VIDEOS ARE UPLOADED



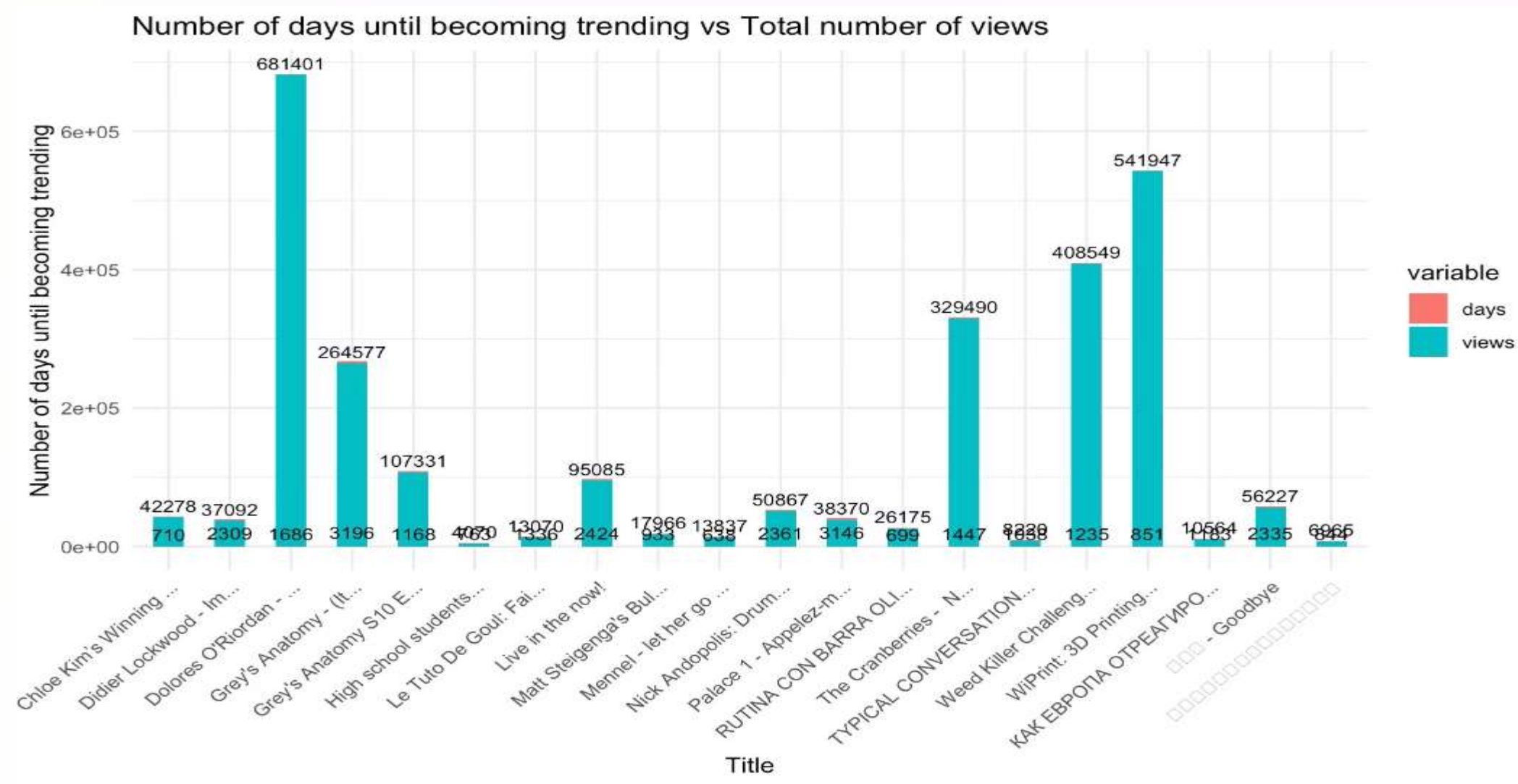


# TIME BETWEEN PUBLISHING AND TRENDING?





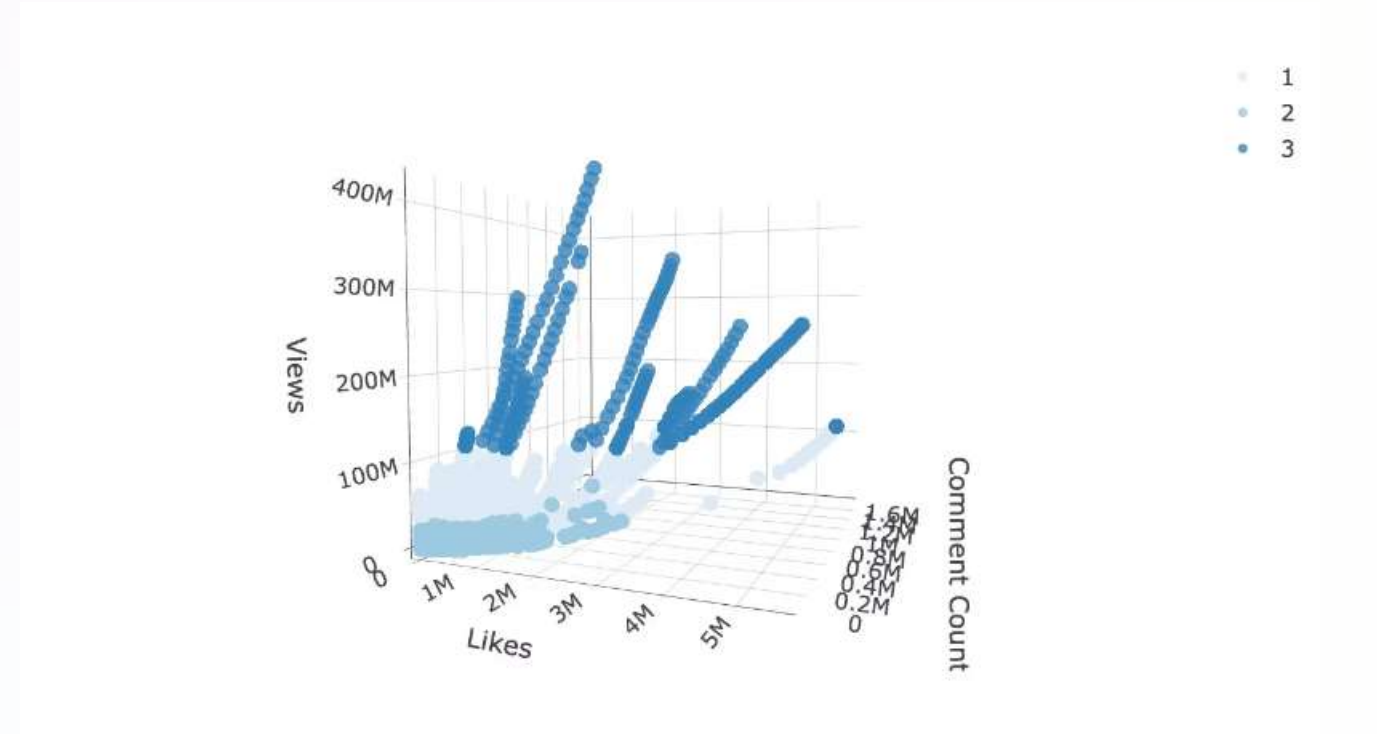
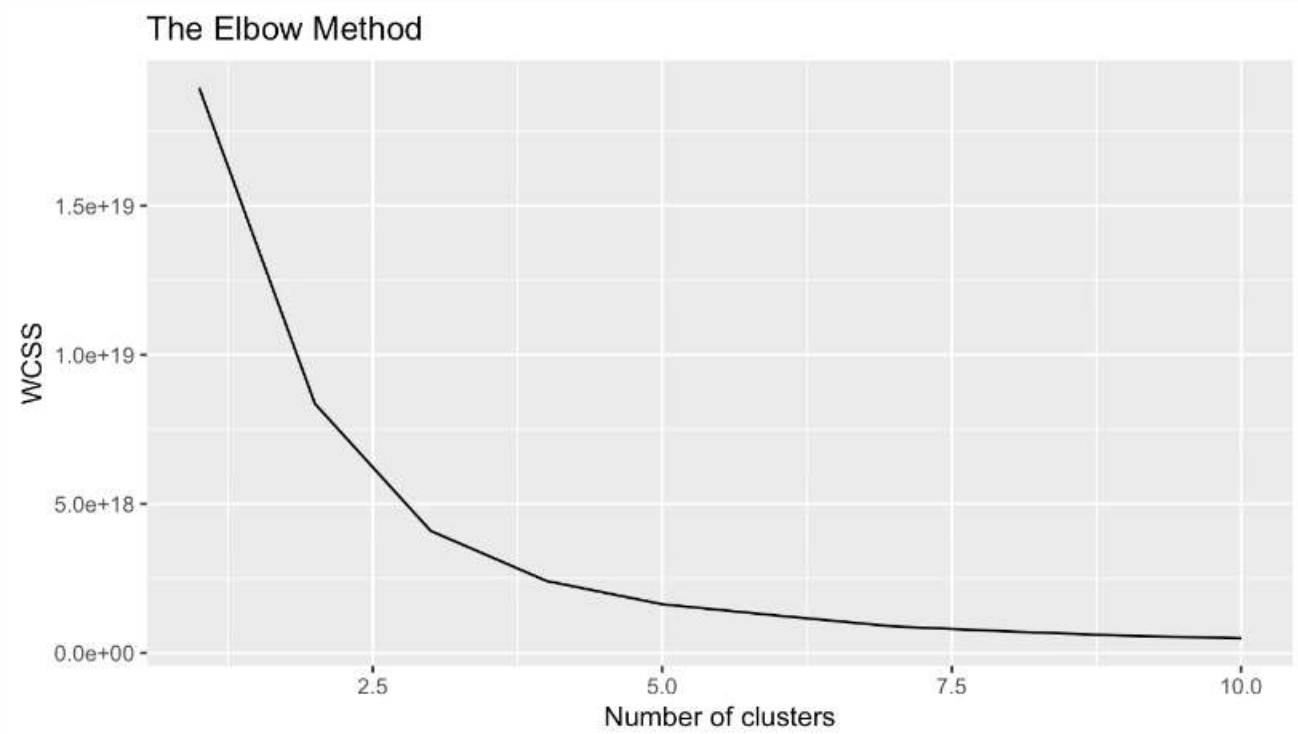
# VIDEOS THAT WAITED THE LONGEST BEFORE THEY BECAME TRENDING





# Model Training and Evaluation

## k-means clustering



# Cluster Analysis

## Cluster 1 (Size: 2946):

- The mean values for `data_total.likes`, `data_total.dislikes`, `data_total.views`, and `data_total.comment_count` are notably high.
- This indicates that Cluster 1 likely corresponds to highly popular and engaging videos, characterized by substantial numbers of likes, views, dislikes, and comments.

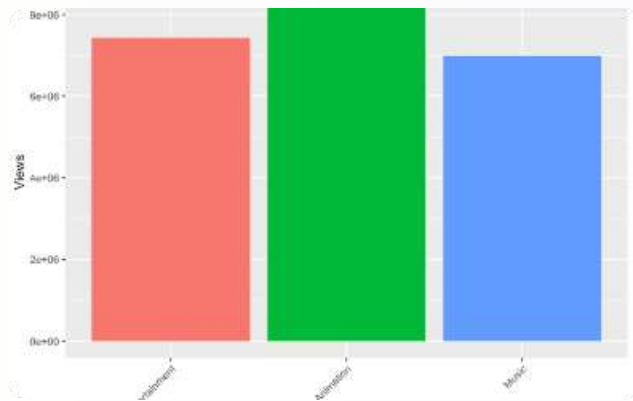
## Cluster 2 (Size: 372768):

- This cluster is the largest among the three identified clusters.
- The mean values for `data_total.likes`, `data_total.dislikes`, `data_total.views`, and `data_total.comment_count` are significantly lower than those of Cluster 1 but higher than those of Cluster 3.
- This suggests that the cluster likely represents moderately popular and engaging videos, characterized by a moderate number of likes, views, dislikes, and comments.

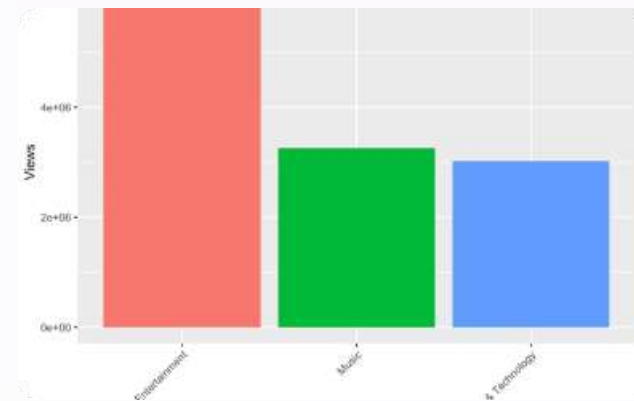
## Cluster 3 (Size: 228):

- This cluster is the smallest among the three identified clusters.
- It exhibits the lowest mean values for `data_total.likes`, `data_total.dislikes`, `data_total.views`, and `data_total.comment_count` compared to the other clusters.
- This suggests that the cluster likely represents videos or content that are less popular and less engaging, characterized by comparatively fewer likes, views, dislikes, and comments.

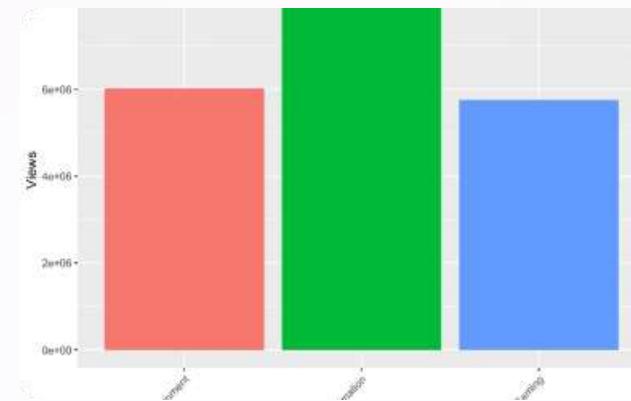
# Analysis on scraped data



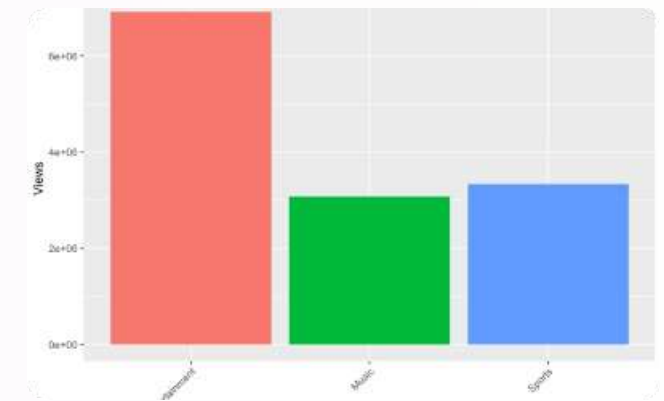
Top 3 Categories in Canada



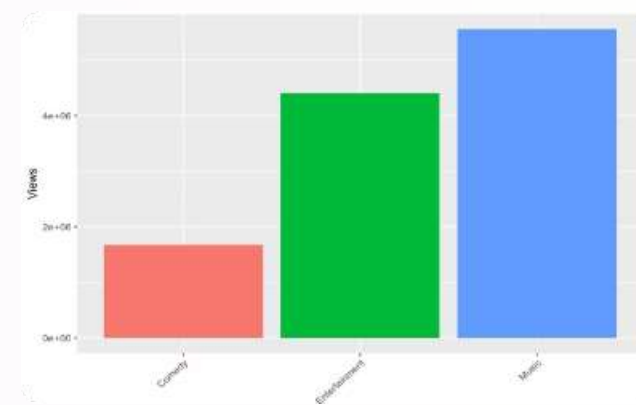
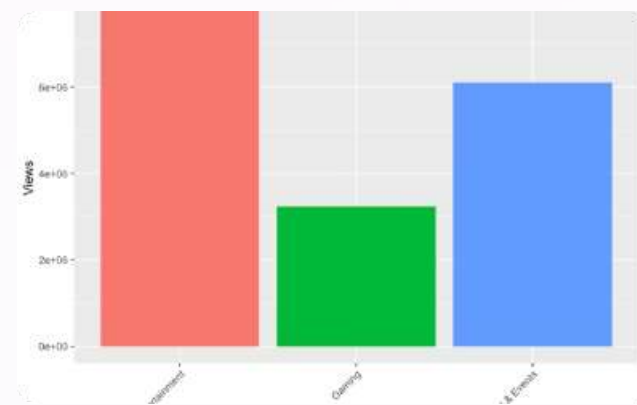
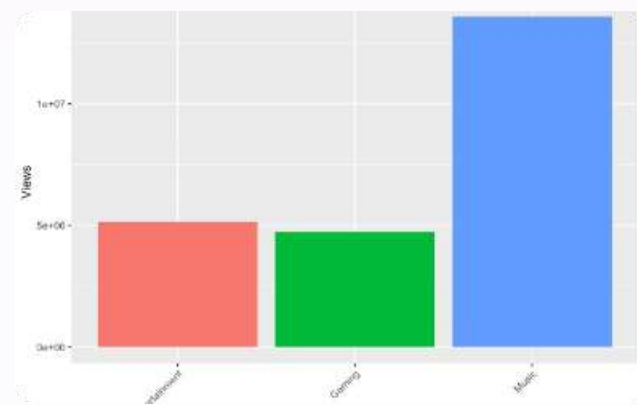
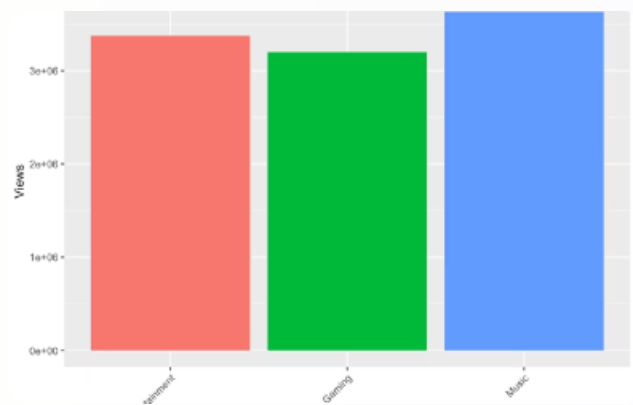
Top 3 Categories in France



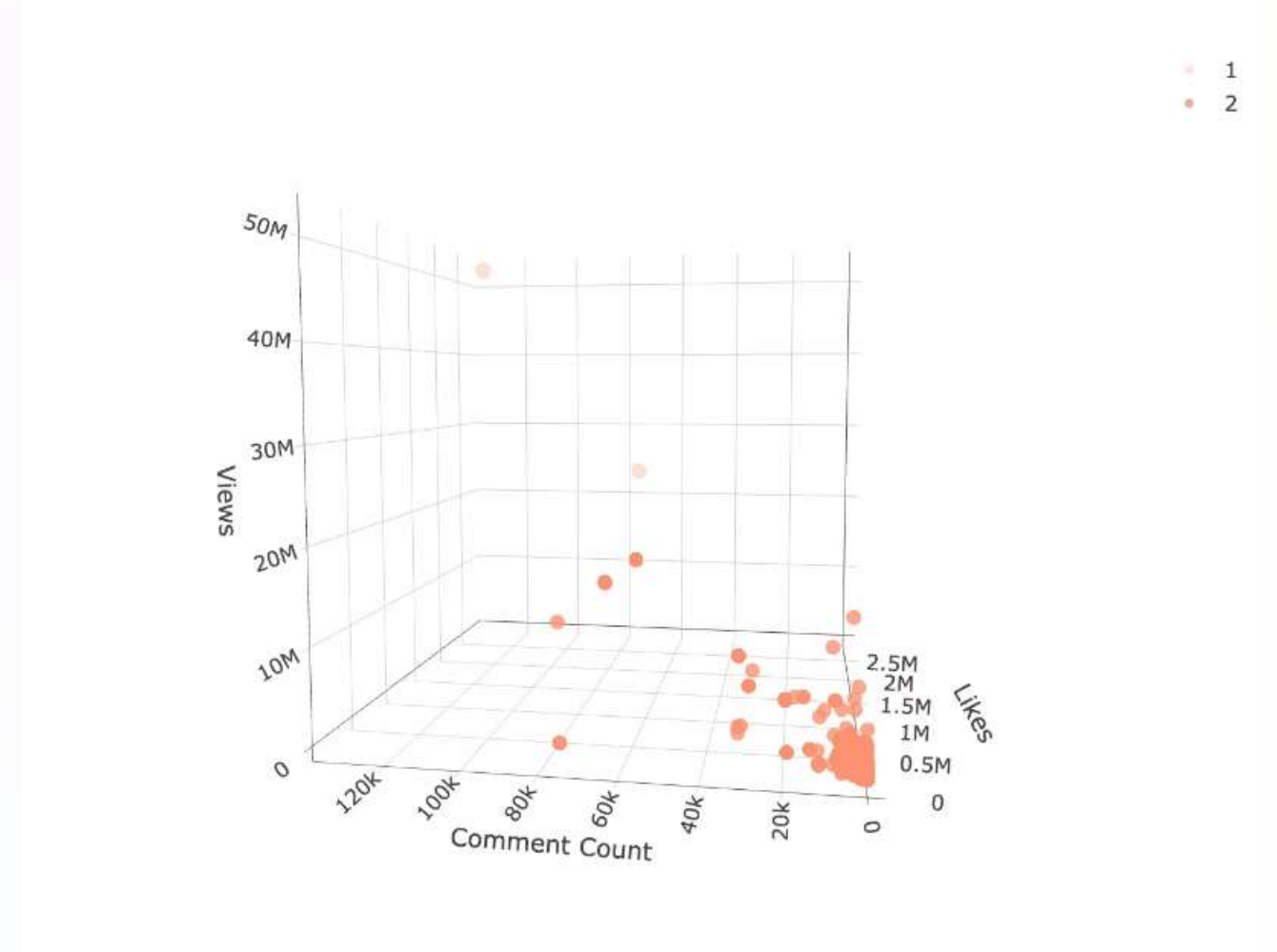
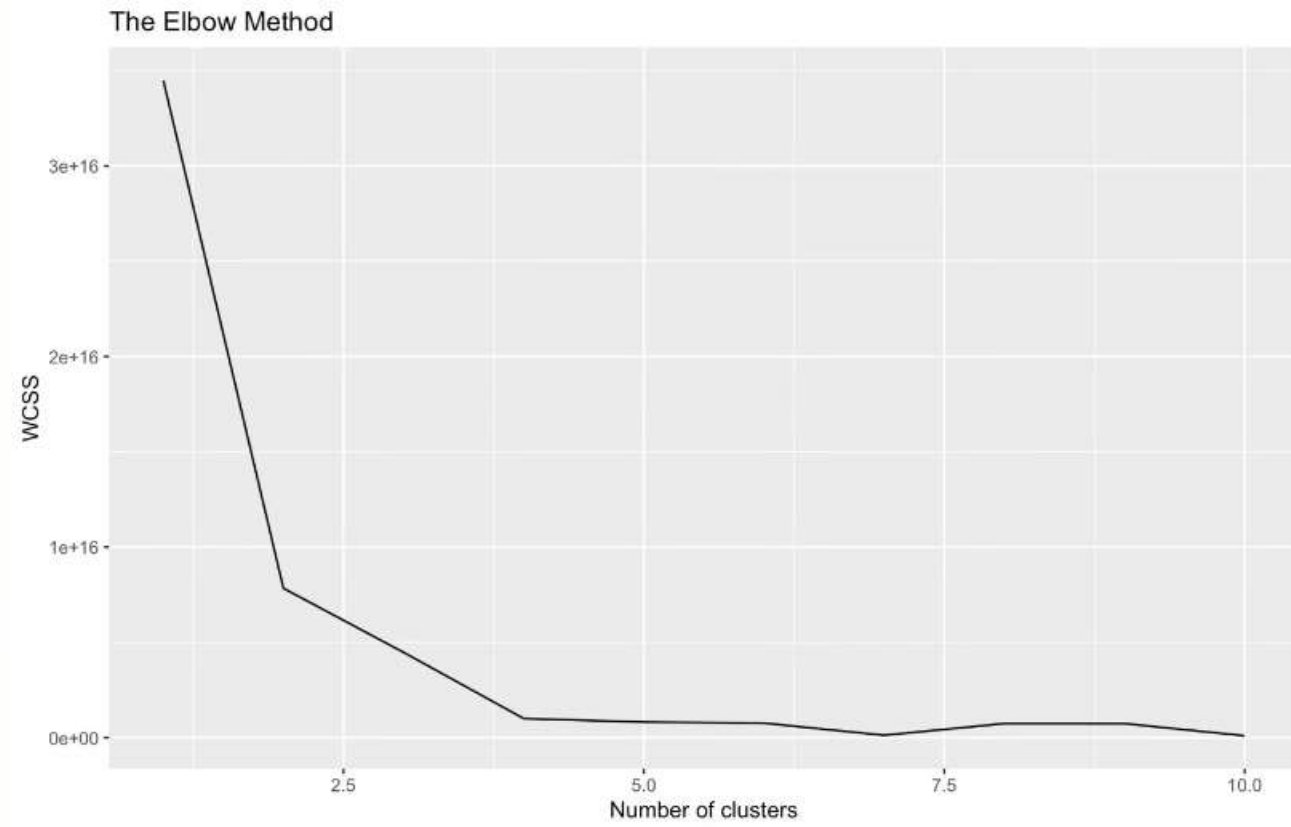
Top 3 Categories in Germany



Top 3 Categories in Mexico



# k-means clustering on scraped data



# Conclusions

Gained insights into the optimal choice of video category, the ideal title length, and the potential influence of description sentiment on viewership.

Compared the findings with current trends to obtain a comprehensive view of the data. Overall, the trends observed in the current data closely align with those identified in past data.

An attempt was made to perform cluster analysis to uncover more nuanced insights; however, the results were not as successful as anticipated.

Future work involves a deeper exploration of the identified clusters and the development of methods to predict video trends.

Thank You!