

Data Analysis in Trending YouTube Videos for Category People and Blog

Formatted: Indent: Left: 1.27 cm

TUNKU KHAIRI BIN TUNKU HANIZD¹ & MOHAMMAD BIN HOSSIN²

¹Faculty of Computer Science and Information Technology, Universiti Malaysia Sarawak, 94300 Kota Samarahan, Sarawak, Malaysia

Corresponding authors: tunkukhairi@gmail.com; hmohamma@unimas.my

ABSTRACT

The study of users' behaviors on YouTube has been an interesting topic in the research world since YouTube started to rapidly grow into one of the largest video-sharing platforms on the internet. In this project, we aim to study the factors that can give positive impacts on videos in the category "People and Blog" to attract the viewers' interest to interact with the videos using data analysis.

Keywords: YouTube, Trending Videos, Opinion mining, Exploratory Data Analysis, Sentiment Analysis, data science, content creators

Copyright: This is an open-access article distributed under the terms of the CC-BY-NC-SA (Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License) which permits unrestricted use, distribution, and reproduction in any medium, for non-commercial purposes, provided the original work of the author(s) is properly cited.

INTRODUCTION

YouTube as a medium of expression provides public statistics for each uploaded video which are upload date, number of views, number of likes, and dislikes. YouTube also provides a comment section for sentimental engagement of the videos. YouTube determines the high engagement videos called "trending videos" by indicating the level of popularity of the videos based on the high number of views, likes, and positive comments for each video. Although YouTube is a platform where content creators could freely upload their videos based on YouTube's guidelines, it is difficult to have high interactions in a single video. YouTube categorizes the videos by which the topic is covered in the video. One of the categories is "People and Blogs". This category covers videos that are related to people's lifestyles, news about people, promotions, reviews, blogs, and topics that correlate with people.

YouTube has massive datasets that are categorized as Big Data. According to an Alexa report, YouTube has become one of the most preferred digital video platforms that intercept more than 30 million visitors in a day. The largest portion of viewers comes from the USA with 16.4% of traffic followed by India (9.2%) and Japan (4.8%) respectively. Because of high traffic and more videos being added every day, unveiling the viewership pattern is considered a complex process for a normal content creator with no Data Science background. Thus, it is time-consuming and potentially can be misleading when interpreting the data. Assist content creators for better insights will also potentially open the opportunity to greatly increase the engagement quality of their videos thus improving the marketing strategies of their videos. This research intends to uncover the hidden pattern by answering various questions about the trending YouTube videos using co-related analysis of Exploratory Data Analysis and Sentiment Analysis.

MATERIALS & METHODS

Specifically, Language used is Python and Jupyter Notebook extension is utilized for code documentation. The library used for this project will be explained according to each method step. The project overview is shown in Figure 1.

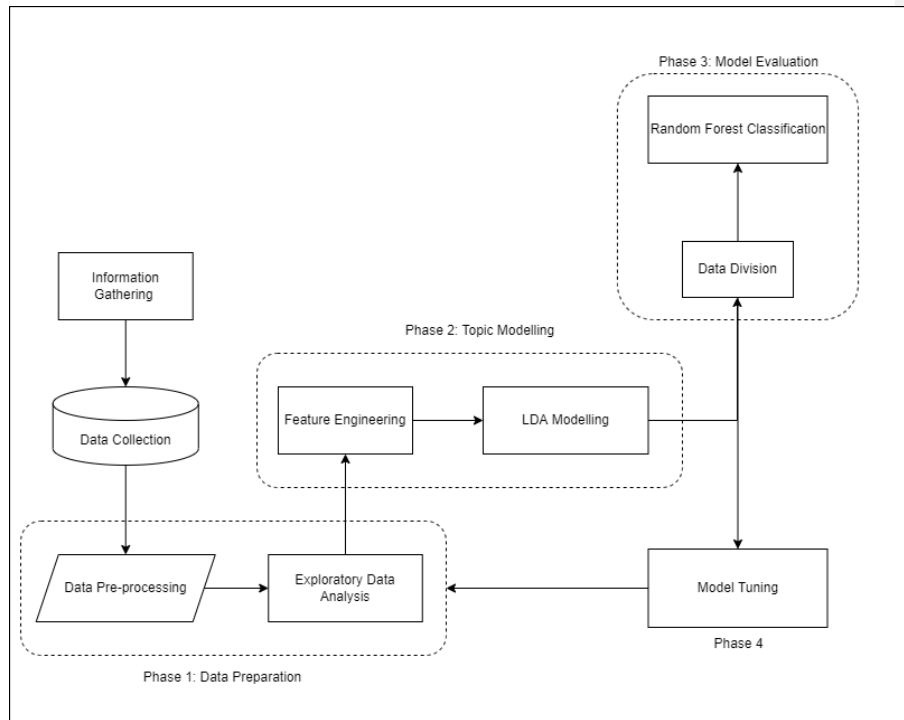


Figure 1. The detailed Project Pipeline.

Data Collection

The dataset used will be a pre-cleaned dataset from the Kaggle website, “YouTube Trending Video Dataset (updated daily)” uploaded by Rishav Sharma because of the high dimensionality of the dataset with collected data from 11 different countries which are India, the US, Great Britain, Germany, Canada, France, Russia, Brazil, Mexico, South Korea, and Japan respectively. The dataset is consisted of up to 200 listed trending videos per day. The dataset is fetched is separated by regions. The list of data attributes is as below:

Video ID
 Video title
 Date published
 Channel ID
 Channel title
 Category ID
 Trending date
 Video tags
 Count of views

Likes

Data Preparation

As an overall pre-processing method, the dataset is inspected to identify missing values. The missing values are then treated. The data preparation process is then divided by numerical columns pre-processing and textual columns pre-processing. For numerical columns, the dataset is pre-processed by observing outliers and treating them to avoid lower accuracy of model training later steps. For the textual column, the dataset is pre-processed by removing punctuations and special characters, data segmentation, normalization of out-of-vocabulary words, conversion of lowercase, and lemmatization in the column 'title', 'description', and 'tags'. The cleaning dataset will be adjusted according to usability for Exploratory Data Analysis and Topic Modelling in the later steps. Python Libraries involved; Klib, NumPy, Pandas, NLTK, and SpaCy.

Exploratory Data Analysis

Exploratory Data Analysis is then used for more comprehensive views of what given data in the dataset is about. By doing so, the best practice of machine learning and the roadmap of the project is identified. The purpose of doing Exploratory Data Analysis is to provide better insights from statistical evidence in Exploratory Data Analysis and utilize it to determine the results in findings.

The analysis involved are:

- Univariate Analysis
- Bivariate Analysis
- Time-series Analysis

Univariate analysis is the simplest form of analysis which involves summarization and pattern of only one "Uni" variable in the dataset dimension. This analysis is conducted to discover the meaning of each column of data, whether it is categorical or continuous, or independent or dependent on other variables in the dataset. Compared to univariate analysis, Bivariate analysis is a correlation analysis that is conducted to gather insights into the causal relationship between two "Bi" variables. Because the dataset contains timestamps, time series analysis is necessary to discover hidden insights based on time intervals. Time series analysis is a specific approach to analyzing a set of data points accumulated over an extended period. Time series analysis is done by manipulating the record of the data points over a set period that is extracted from video published time. Multivariate analysis is the statistical analysis of correlations between several measurements made on each experimental unit and where the relationship between multivariate measurements and their structure is crucial to understanding the experiment. It is done to find suitable features to feed into machine learning models, or to decide whether the data type should be transformed or not in the feature engineering process.

Feature Extraction & Feature Selection

Feature extraction and feature selection will be conducted to reduce the data dimension and further pull relevant data to use for hypothesis testing. The data types are transformed for better data representation. For topic modeling purposes in later steps, the textual column 'title', 'tags', and 'description' is merged into the 'all_text' column. The purpose of merging is to enhance the context of specific videos for model training later. Reducing the dimension of the data helps to improve the accuracy of the model and reduce the execution time when building and tuning the model. Feature Selection includes the process of removing stop words, and removing null rows in the dataset, creating of dictionary and data annotation.

Topic Modeling

For the topic modeling phase, the chosen analysis for the proposed project is using the unsupervised LDA modeling method. The LDA modeling process includes building a base LDA model, computing perplexity, and coherence score. The perplexity and coherence scores are then observed and evaluated by assigning each document in the dataset with its predicted topic.

Model Evaluation

For the model evaluation process, the dataset with its predicted topic is then divided into a training set and a testing set with a ratio of 80% training set and 20% testing set. The supervised Random Forest Classification method is used to evaluate the predicted topic in terms of the Bag-Of-Word model and TF-IDF model. The evaluation is observed from the value of precision, recall F-1 score, and support obtained after the cross-validation process from training and testing sets.

Model Tuning

To increase the performance of the LDA model, there are 2 steps of model tuning:

- Building LDA Mallet Model.
- Finding the optimal number of topics, k .

Hyperparameter is manipulated to determine the optimal number of topics, k . This is done by building many LDA models and then plotting a computed coherence score graph for a range of k . Generally, the k value is determined by the model that has the highest coherence score before flattening out.

The project will then repeat the preprocessing method followed by the topic modeling phase as part of the tuning model process until an optimal number of topic k is achieved. The optimal k value is then used to predict the true dominant topic for each of the documents.

RESULTS & DISCUSSION

Exploratory Data Analysis

Univariate analysis is conducted based on how videos in the dataset are distributed respectively to view counts, like counts and dislikes count, and tags.

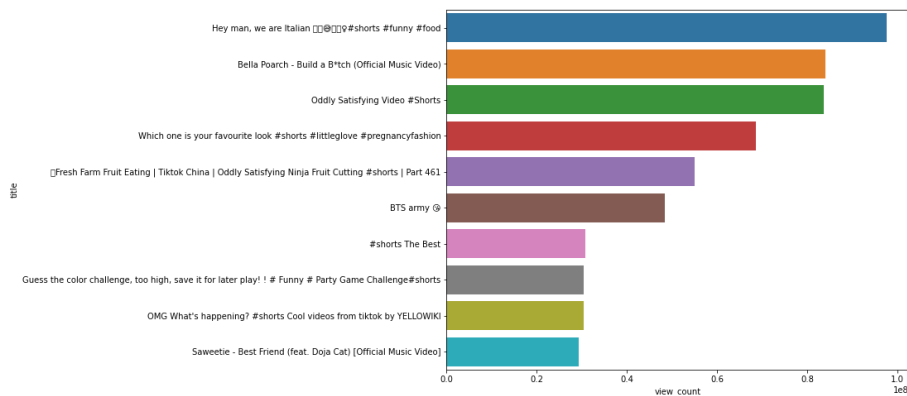


Figure 2. Top 10 Most Viewed Videos in Category *People and Blog*.

Based on Figure 2 above, proves that most viewed videos from the category “People and blog” have wide topics posted. Other than that, the majority of the top viewed videos consist of hashtag #shorts in a count of 7 out of 10 most viewed videos. Roughly, #shorts are used to indicate that the duration of the video is short.

Table 1. Most Frequent Hashtags used in Category *People and Blog*.

Hashtags	Count
#shorts	846
#viral	58
#funny	47
#POV	35
#food	25
#4	32
#meme	19
#trending	16
#minecraft	16
#FYP	16
#Dpeezy2099	16
#1	16
#2	13
#ad	12
#dowehaveaproblem	11

YouTube has a new feature that enables content creators to upload their videos in a form of short videos. Based on Table 1 above, the analysis has proved that usage of the new feature is effective to attract viewers. After removing the word ‘short’ along with other stop words, the WorkCloud for tags used in the videos is generated as shown in Figure 3.



Figure 3. WordCloud for Most Tags used in Category *People and Blog*.

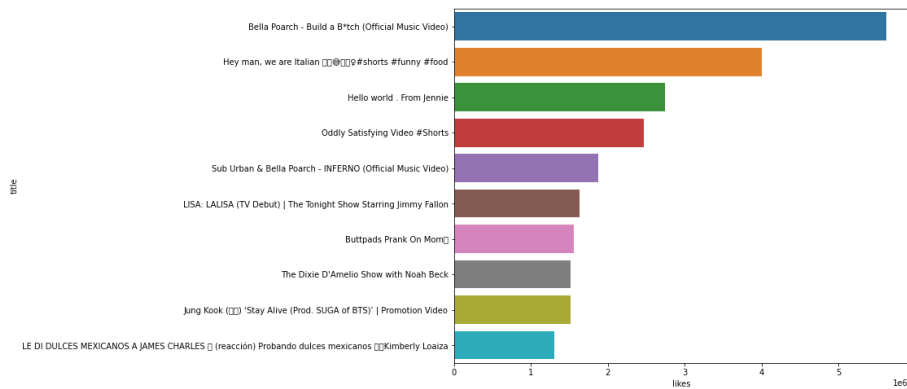


Figure 4. Top 10 Most Liked Videos in Category *People and Blog*.

Based on Figure 4 above, there are slight changes in the order of the Top 10 videos. This indicates that not every trending video has positive feedback despite higher views.

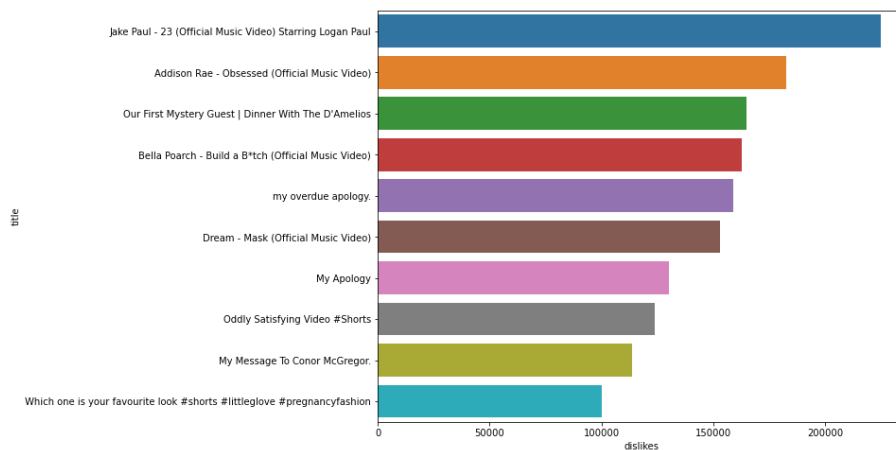


Figure 5. Top 10 Most Disliked Videos in Category *People and Blog*.

Based on Figure 5 above, generally, videos with high views tend to have negative feedback too. For Bivariate Analysis, bicorrelation between views, likes, dislikes, and channels is observed.

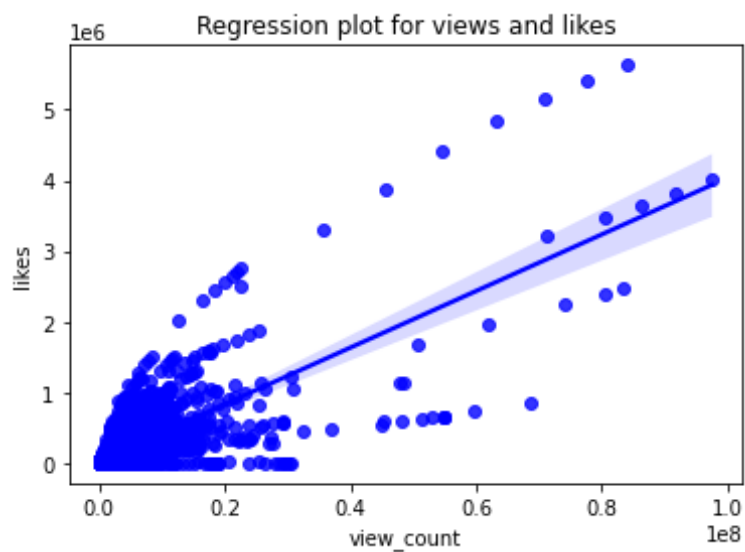


Figure 6. Regression Plot for Views and Likes for Category *People and Blog*.

Figure 6 above shows that the relationship between views and likes is linearly related. From the plot, a large percentile of videos is scattered close to the simple linear regression line while a minor percentile of videos is considered outliers. However, the residual value of the outliers is distributed linearly too. This resulting the videos having more than one pattern. The simple assumption can be made that likes do not always dependent on the view count. Some possibilities are taken from the scattered plot, where view counts of 2 videos are the same, but one is highly skewed from the line of best fit likes and another one has lowlily skewed from the line of best fit likes.

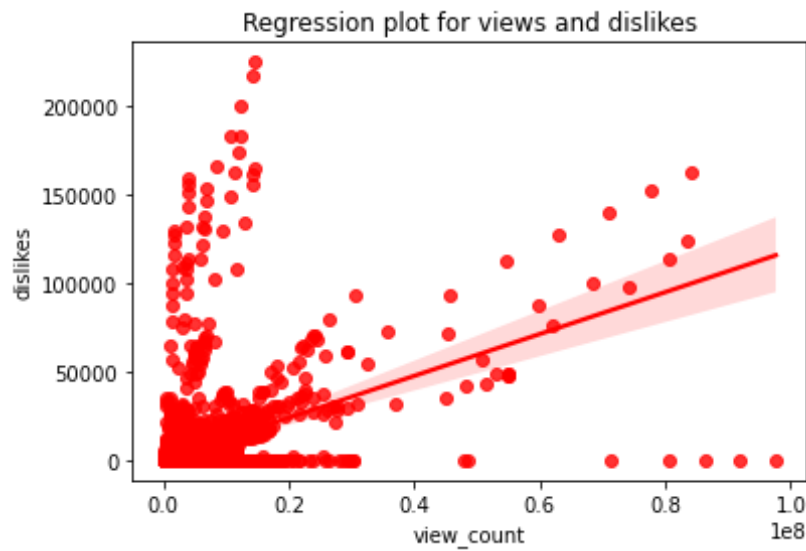


Figure 7. Regression Plot for Views and Dislikes for Category *People and Blog*.

Figure 7 shows that a simple linear regression line does not fit views and likes. While a large percentage of videos is scattered close to the line, some outliers potentially show different behavioral patterns. The usual case is where the greater view counts, the greater amount of dislikes a video has. However, the behavioral pattern mentioned earlier is that there are videos that despite having low views, the videos might have a great number of dislikes. This can be interpreted as videos being extremely disliked by the viewer, but they watch them anyway. Another behavioral pattern is there are videos that despite having high views, the videos might have a lesser number of dislikes (compared to the best of fit line). This can be interpreted as the videos being extremely liked by the viewers. Notes that the dataset is originally extracted from only high-interacted videos where the proposition videos of low views in the statistic are assumed higher than average videos on YouTube. From assumptions made earlier, the deduction is the video does not particularly need to be likable to viewers to reach high interactions.

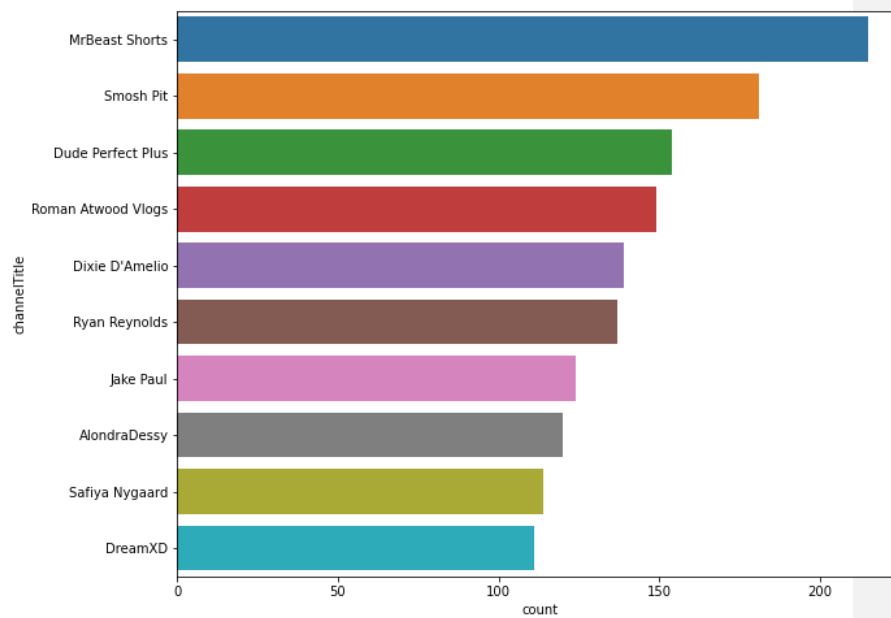


Figure 8. Top 10 Channels with Most Videos on Trending.

Figure 8 above shows that channels that have higher subscribers tend to have more of their uploaded videos trending. For Multivariate Analysis the Pearson Correlation Chart is observed.

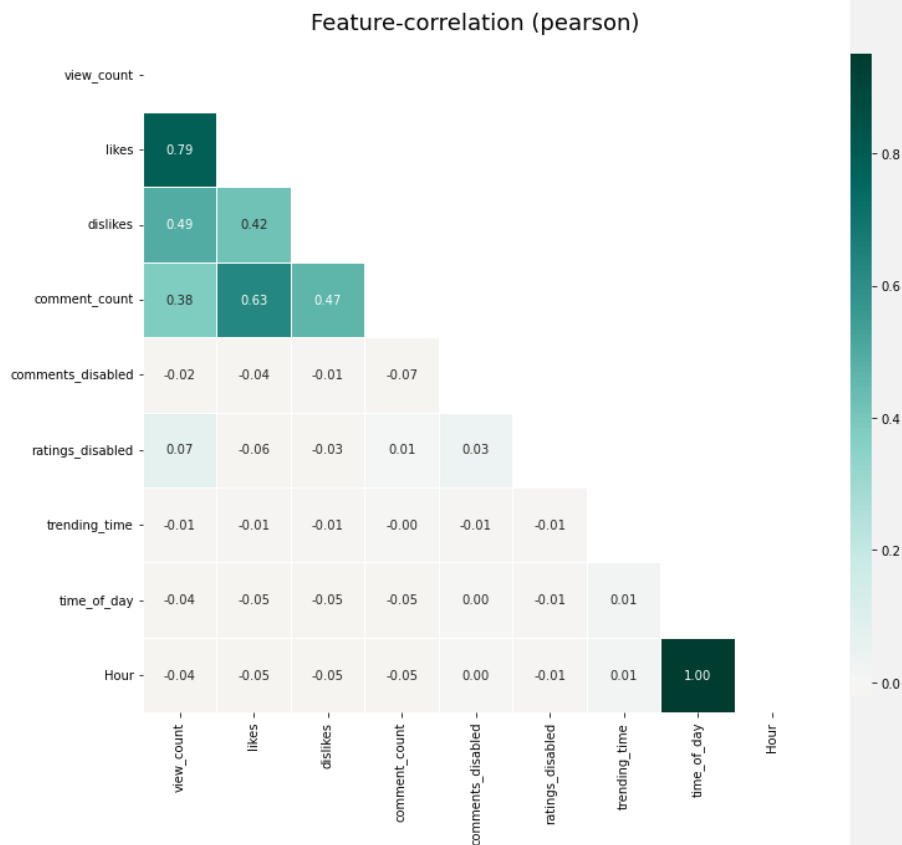


Figure 9. Pearson Feature Correlation Chart for Category *People and Blog*.

Figure 9 shows the Pearson correlation chart for numerical features in the dataset. From the chart, some assumptions can be made. The first assumption is, likes are highly positively correlated to view counts with a 0.79 score meanwhile dislikes are on average positively correlated to view counts. Comments counts are highly positively correlated to likes, however, it has an average influence on view counts and dislikes.

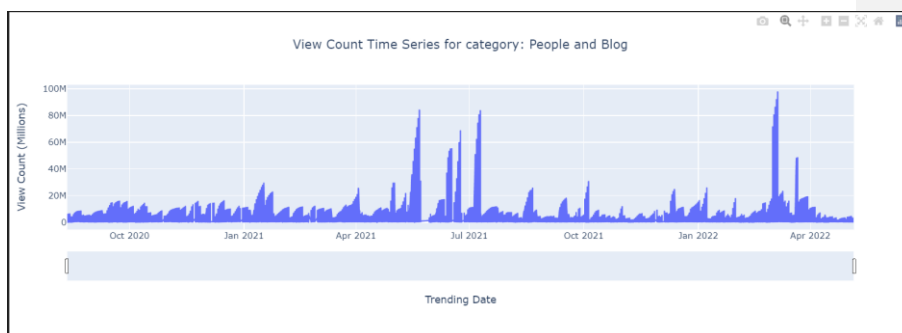


Figure 10. Overall Time Series of View Count for Category *People and Blog*.

Figure 10 above shows that there are some unusual spikes of view counts over 80 million views at several timeframes.

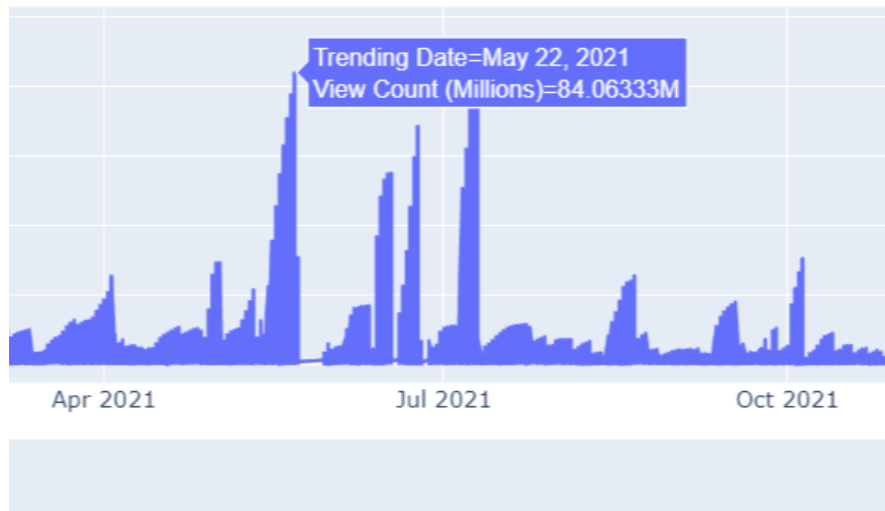


Figure 11. Views Spiked up to 84.0633 million on May 22, 2021.

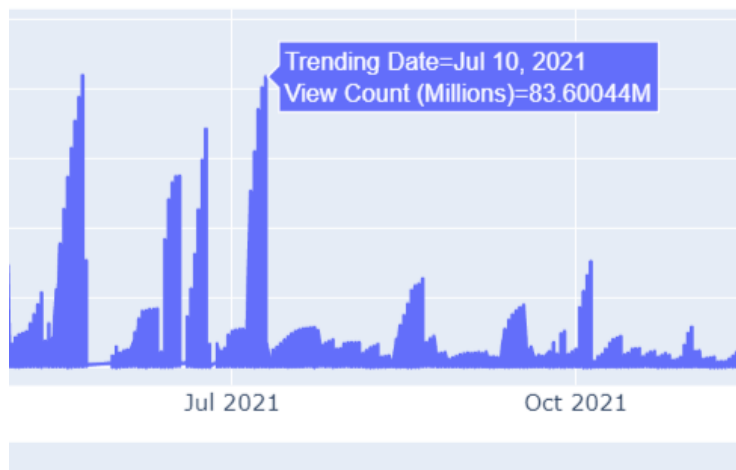


Figure 12. Views Spiked up to 83.6 million on July 10, 2021.

decided to change or quit their current jobs. Word “Cardboard” may indicate that videos about manipulating cardboards are trending at the time.

From 3 generated WordCloud, the finding is real-world events do impact greatly on view counts. The videos that are related to exact events in the real world tend to get more interactions than videos that are not related.

Topic Modelling

For the base model, the number of topics is randomly assigned for $k = 10$ (in Python num_of_topic = 10). The result of the base LDA model is shown in Figure 18 below.

```
# Compute Perplexity
print('\nPerplexity: ', lda_model.log_perplexity(dt_corpus)) # a measure of how good the model is. lower the better.

# Compute Coherence Score
coherence_model_lda = CoherenceModel(model=lda_model, texts=data_lemmatized, dictionary=id2word, coherence='c_v')
coherence_lda = coherence_model_lda.get_coherence()
print('\nCoherence Score: ', coherence_lda)

Perplexity: -8.075182183006767

Coherence Score: 0.35944259744198476
```

Figure 17. The Perplexity and Coherence Score for Base LDA Model.

Based on Figure 17 above, the perplexity and coherence scores are observed. Observation is made to measure how good the LDA model is. For perplexity score, the lower the score, the higher chance that the base LDA model is good. A coherence score is used to evaluate a single topic's score by gauging the degree of semantic similarity between the topic's top-scoring words. The higher the coherence score, the better the chance the LDA model is good.

Model Tuning

The base LDA model is compared to the base LDA Mallet model that uses the Gibbs Sampling method. Figure 19 below shows the result of the coherence score for the base LDA Mallet model.

```
# Compute Coherence Score for Mallet
coherence_model_lda = gensim.models.CoherenceModel(model=ldamallet, texts=data_lemmatized, dictionary=id2word, coherence='c_v')
coherence_lda = coherence_model_lda.get_coherence()
print('\nCoherence Score: ', coherence_lda)

Coherence Score: 0.365307724851063
```

Figure 18. The Perplexity and Coherence Score for Base LDA Mallet Model.

Based on Figure 18, the coherence score for the base LDA Mallet model is slightly higher than the base LDA model. Hence, the LDA Mallet model is chosen to compute the optimal model. Then, the next step is tuning the model by finding the optimal number of topics. This is done by computing multiple LDA Mallet models and calculating their coherence score 'c_v' vs the number of topics, k for k is in the range of $2 \leq k \leq 40$. The results are plotted in a graph as shown in Figure 20 below.

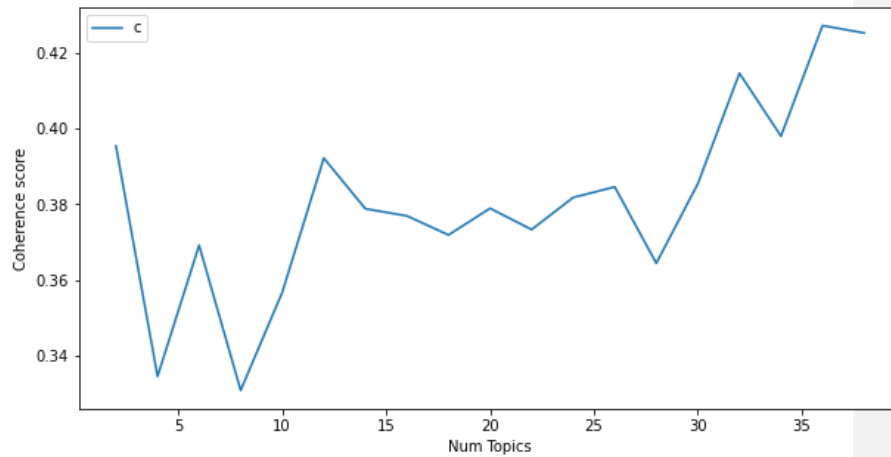


Figure 19. Graph of LDA Mallet Model's Coherence Score vs Number of Topic.

Based on the graph in Figure 19, the optimal number of topics is chosen subjectively at the point where the coherence score is at the peak before flattening out. Hence the chosen optimal number of topics, $k = 12$. By choosing $k = 12$, the optimal LDA Mallet is then built. Next, each dominant topic is computed with the weightage of the keywords as shown in Table 1 below.

Table 2. The Dominant Topic Computed based on Number of Topics, $k = 12$

Keywords	Dominant Topic
baby, vlog, pregnant, code, boy, couple, find, pregnancy, girl, birth	1
TikTok, music, show, official, love, song, happy, hope, fun, tik_tok	2
dream, kid, official, make, support, time, day, long, update, back	3
short, funny, story, part, prank, among_us, viral, challenge, moment, game	4
good, watch, friend, click, big, play, scene, dude_perfect, laugh, comedy	5
family, vlog, day, life, surprise, challenge, royalty, wedding, house, watch	6
eat, make, food, cook, ASMR, recipe, cooking, babish, steak, Minecraft	7
link, find, free, home, store, business, worth, content, speed, episode	8
move, makeup, year, fashion, check_out, give, car, world, beauty, make	9
season, watch, voice, live, podcast, team, episode, series, clip, highlight	10
Getty, make, dog, room, hour, leave, real_life, challenge, turn, man	11
life, home, camera, build, water, house, live, cabin, off_grid, tiny	12

The dominant topics are then used to predict each video in the dataset. To evaluate the performance of the optimal model, the confusion matrix is calculated by feeding the dominant topic and combining text (video title, video description, and video tags) in each video into two models, the Bag-of-Word model and TF-IDF model using supervised Random Forest Classification.

The classification report for Bag-of-model is shown in Figure 20 below.

	precision	recall	f1-score	support
0	0.34	0.88	0.49	40
1	0.66	0.80	0.72	50
2	0.74	0.37	0.49	38
3	0.57	0.77	0.66	53
4	0.92	0.50	0.65	24
5	0.58	0.66	0.61	29
6	1.00	0.65	0.79	37
7	0.96	0.57	0.72	40
8	0.48	0.34	0.40	32
9	0.65	0.39	0.49	28
10	0.85	0.71	0.77	31
11	0.87	0.45	0.59	29
accuracy			0.61	431
macro avg	0.72	0.59	0.61	431
weighted avg	0.70	0.61	0.62	431

Figure 20. Bag-of-Word Model Classification Report.

Based on Figure 20, the performance of the BoW model is 0.61. The purpose of using the BoW model is to measure how well the document classification is based on the frequency occurrence of each word when it is used as a feature in the model classifier. While having a decent precision model, the low recall score made the F-1 score lowered to 0.61. For TF-IDF model evaluation, the classification result is shown in Figure 22 below.

	precision	recall	f1-score	support
0	0.30	0.90	0.45	40
1	0.67	0.84	0.74	50
2	0.61	0.37	0.46	38
3	0.59	0.70	0.64	53
4	0.92	0.46	0.61	24
5	0.65	0.59	0.62	29
6	1.00	0.62	0.77	37
7	0.96	0.60	0.74	40
8	0.64	0.44	0.52	32
9	0.64	0.32	0.43	28
10	0.96	0.74	0.84	31
11	0.87	0.45	0.59	29
accuracy			0.61	431
macro avg	0.73	0.59	0.62	431
weighted avg	0.72	0.61	0.62	431

Figure 21. TF-IDF Model Classification Report.

The purpose of using the TF-IDF model is to evaluate how well the document classification is based on the most relevant word in the document. Based on Figure 21, the performance of the TF-IDF model is also 0.61. While having a decent precision model, the low recall score made the F-1 score lowered to 0.61.

Based on the two evaluation model discussed, the performance of the model is assumed to be low bias but have high variance. The model is overfitted due to the fact when data pre-processing, only “English” words are selected to be fed into the model while the dataset should consist of different other languages’ text. In the testing phase, the recall score has become low due to the incapability of the model to predict the text in other languages text. Hence the model only fitted for English Languages text.

CONCLUSION

Exploratory Data Analysis has shown several findings. Firstly, the usage of the new feature is effective to attract viewers. Secondly, likes are not always dependent on the view count and the deduction is the video does not particularly need to be likable to viewers to reach high interactions. Some possibilities are taken from the analysis where exceptional videos are made without having many likes. Next, the finding is real-world events do impact greatly on view counts. The videos that are related to exact events in the real world tend to get more interactions than videos that are not related. Using unsupervised machine learning, the main topic of trending video is discovered. Content creators might use the main topic produced as a guideline to generate title, description, and tags for their video as it was used by other successful content creators to make their video appearance stands out.

REFERENCES

- Snickars, P., & Vonderau, P. (2009). *The youtube reader*. Kungliga biblioteket.
- Holmbom, M. (2015). The YouTuber: A qualitative study of popular content creators.
- Rowley, J. (2004). Online branding. *Online Information Review*, 28(2), 131-138.
- Biel, J. I., & Gatica-Perez, D. (2010, November). Vlogcast yourself: Nonverbal behavior and attention in social media. In *International Conference on Multimodal Interfaces and the Workshop on Machine Learning for Multimodal Interaction* (pp. 1-4).
- Biel, J. I., Aran, O., & Gatica-Perez, D. (2011, July). You are known by how you vlog: Personality impressions and nonverbal behavior in youtube. In *Proceedings of the International AAAI Conference on Web and Social Media* (Vol. 5, No. 1).
- Luers, W. (2007). *Cinema without show business: A poetics of vlogging* (Vol. 5, No. 1). Ann Arbor, MI: MPublishing, University of Michigan Library.
- Wauters, Robin. (2010, January 5). State of the vlogosphere. Retrieved from <http://techcrunch.com/2010/01/05/mefeedia-state-of-the-vlogosphere-2010/>
- Mitchell, J. C. (1969). The concept and use of social networks. In J. C. Mitchell (Ed.), *Social networks in urban situations*. Manchester, England: University of Manchester Press, 1969.
- Arthurs, J., Drakopoulou, S., & Gandini, A. (2018). Researching youtube. *Convergence*, 24(1), 3-15.
- Dehghani, M., Niaki, M. K., Ramezani, I., & Sali, R. (2016). Evaluating the influence of YouTube advertising for attraction of young customers. *Computers in human behavior*, 59, 165-172.
- Usmani, Z. (2017). *Kaggle for Beginners: with Kernel Code*. Gufhtugu Publications.
- Ghojogh, B., Samad, M. N., Mashhadi, S. A., Kapoor, T., Ali, W., Karray, F., & Crowley, M. (2019). Feature selection and feature extraction in pattern analysis: A literature review. *arXiv preprint arXiv:1905.02845*.
- Everitt, B. S., Landau, S., & Leese, M. (2001). *Cluster Analysis*, ch. 4.

Wu, K. (2016). YouTube marketing: Legality of sponsorship and endorsements in advertising. *JL Bus. & Ethics*, 22, 59.

Firat, D. (2019). YouTube advertising value and its effects on purchase intention. *Journal of Global Business Insights*, 4(2), 141-155.

Bendersky, M., Garcia-Pueyo, L., Harmsen, J., Josifovski, V., & Lepikhin, D. (2014, August). Up next: retrieval methods for large scale related video suggestion. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining* (pp. 1769-1778).

Choudhury, S., & Breslin, J. G. (2010). User sentiment detection: a YouTube use case. Rogers, R. (2013). *Digital methods*. MIT press.