

VidMine: Digging Deep into YouTube Insights

Amulya Ambati
Sneha Sasanapuri
Lakshmi Prasanna Poluru

March 23, 2024

Abstract

With the exponential growth of digital content consumption, understanding user behavior and content trends on platforms like YouTube has become paramount for content creators, marketers, and businesses. "VidMine" presents a comprehensive data mining project aimed at unearthing valuable insights from YouTube's vast repository of data. Leveraging the YouTube Data API and advanced data mining techniques, this project delves deep into YouTube's ecosystem to extract, analyze, and interpret key metrics, trends, and patterns. The project's primary objectives include data extraction, analysis, insight generation, visualization, and recommendations. Through a systematic data collection, "VidMine" uncovers actionable insights pertaining to user engagement, content performance, audience demographics, and emerging trends. Natural Language Processing (NLP) techniques are employed to analyze textual data such as comments and descriptions, providing valuable insights into audience sentiment and engagement.

Contents

1	Introduction	3
2	Structure of the Report	3
3	Understanding the YouTube Landscape	3
4	Data	3
4.1	Introduction to Data Source	3
4.1.1	Data Retrieval from YouTube Data API	3
4.2	Attributes for analysis	4
5	Methods	4
5.1	Data Cleaning	4
5.1.1	Dealing with Memory issues	4
5.1.2	Parsing Date Column	5
6	Exploratory Data Analysis(EDA)	5
6.1	Data Visualisations	5
6.1.1	Top channels based on views	5
6.1.2	Insights into the Distribution of Content Creation Time Among Channels	6
6.1.3	Analyzing the relationship between Video Duration and View Count	7
6.1.4	Exploring videos durations Across Years: A Box Plot Analysis	8
6.1.5	Distribution of Video Generation on a Monthly Basis	8
6.1.6	Exploring the Relationship Between Likes and Comments: Insights into Viewer Engagement Patterns	9
6.1.7	Word Cloud of videos descriptions	10
6.1.8	Analyzing Viewer Engagement Based on Captions Presence in Videos	11
6.1.9	Exploring Video Duration Distribution Over Time with Captions Analysis	12
6.1.10	Determining association between Duration and Number of Likes	12
6.1.11	Exploring Univariate Distribution of Video Durations	13

1 Introduction

In the digital age, where online video consumption is at its peak, platforms like YouTube have emerged as ubiquitous sources of entertainment, education, and information. With billions of users and an ever-expanding library of content, YouTube offers a treasure trove of data waiting to be explored. Understanding the dynamics of this vast ecosystem is essential for content creators, marketers, and businesses striving to thrive in the digital landscape.

Through advanced data mining techniques, we embark on a quest to dig deep into the wealth of insights hidden within YouTube’s vast repository of content, engagement metrics, and user interactions. In this report, we present the culmination of our efforts—a comprehensive exploration of YouTube’s data terrain. From data extraction to analysis, visualization, and actionable insights, “VidMine” offers a holistic view of YouTube’s ecosystem and its implications for content creators, marketers, and stakeholders.

2 Structure of the Report

This report is structured to provide a comprehensive overview of the “VidMine” project, encompassing its methodology, findings, insights, and recommendations. We begin by outlining the methodology employed in data extraction, analysis, and visualization. Subsequently, we delve into the insights derived from the analysis, shedding light on user behavior, content trends, and audience engagement metrics. Finally, we conclude with actionable recommendations for leveraging these insights to enhance content strategy, optimize audience engagement, and drive growth on YouTube.

3 Understanding the YouTube Landscape

This paper researches the influence of venue dynamics on cricket match outcomes and prompts the investigation of whether playing at home holds a substantial impact on a team’s likelihood of securing victory, with a consideration of potential exceptions. It also investigates to see if there is an increased likelihood of winning the match for teams that win the toss. Additionally, we aim to discern patterns surrounding teams that demonstrate exceptional prowess or face challenges when competing in away matches.

4 Data

4.1 Introduction to Data Source

Data Retrieval from YouTube Data API

In order to analyze the impact and engagement surrounding the topic of keywords entered in search bar on YouTube, we utilized the YouTube Data API to access various resources such as activities, videos, search, and comments. Through these API endpoints, we were able to gather a comprehensive dataset containing valuable information pertaining to videos related to the specified keyword. For our project, we have taken keywords as “neural networks” and fetched data of top 50 channels.

4.1.1 Data Retrieval from YouTube Data API

Upon querying the YouTube Data API with the search term “neural networks,” we retrieved a multitude of videos matching the search criteria. From these search results, we extracted pertinent information to construct our dataset. The YouTube Data API provides resources to extract data under different module. It is organised effectively to easily identify the different types of resources that can be retrieved using the API.

A “search” resource result contains information about a YouTube video, channel, or playlist that matches the search parameters specified in an API request. An “activity” resource contains information about an action that a particular channel, or user, has taken on YouTube. The actions reported in activity feeds include rating a video, sharing a video, marking a video as a favorite, uploading a video, and so forth. A “channel” resource contains information about a YouTube channel. A video resource represents a YouTube video and is uniquely identified by an id. A “commentThread” resource contains

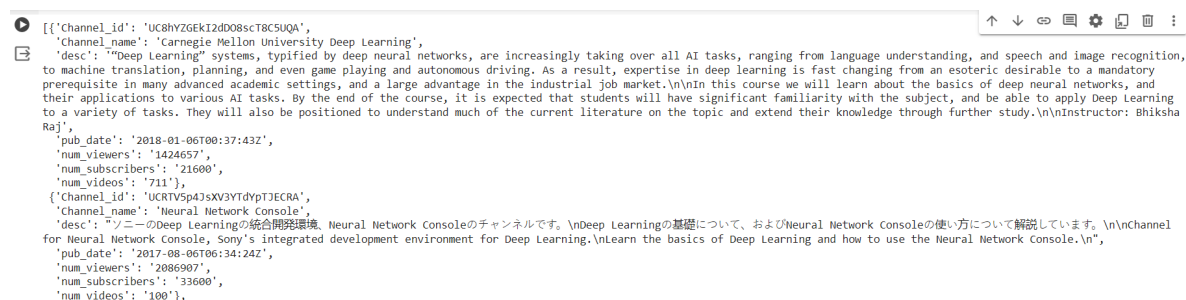
information about a YouTube comment thread, which comprises a top-level comment and replies, if any exist, to that comment. This resource can represent comments about either a video or a channel based on the id supplied.

4.2 Attributes for analysis

The key features extracted include:

- Number of Views: The total number of views garnered by each video, indicating its popularity and reach among viewers.
- Number of Comments: The count of comments posted on each video, providing insights into the level of engagement and interaction generated by the content.
- Date of Published: The timestamp indicating the date and time when each video was published on the platform, enabling temporal analysis and trend identification.
- Number of Likes: The quantity of likes received by each video, reflecting the audience's positive reception and appreciation towards the content.
- Captions Availability: A boolean indicator denoting whether each video has captions or not, facilitating accessibility and catering to diverse audience needs.
- Description: The textual description accompanying each video, offering context and additional information about the content.
- Title: The title of each video, serving as a concise representation of its subject matter and attracting viewer attention.
- Comments: Extracted comments from viewers, providing valuable feedback, opinions, and discussions surrounding the topic of neural networks.
- Time of comment: The extraction of comment timestamps provides valuable insights into the temporal dynamics of user engagement and interaction with YouTube content related to the topic of interest.

The Youtube API returns response in JSON format. Below is the snippet of data before processing.



```
[[{"Channel_id": "UC8hYzGEkI2dD08scT8C5UQA",
  "Channel_name": "Carnegie Mellon University Deep Learning",
  "desc": "\"Deep Learning\" systems, typified by deep neural networks, are increasingly taking over all AI tasks, ranging from language understanding, and speech and image recognition, to machine translation, planning, and even game playing and autonomous driving. As a result, expertise in deep learning is fast changing from an esoteric desirable to a mandatory prerequisite in many advanced academic settings, and a large advantage in the industrial job market.\n\nIn this course we will learn about the basics of deep neural networks, and their applications to various AI tasks. By the end of the course, it is expected that students will have significant familiarity with the subject, and be able to apply Deep Learning to a variety of tasks. They will also be positioned to understand much of the current literature on the topic and extend their knowledge through further study.\n\nInstructor: Bhiksha Raj",
  "pub_date": "2018-01-06T00:37:43Z",
  "num_views": 142457,
  "num_subscribers": 216000,
  "num_videos": 711},
{"Channel_id": "UCRTVSp4JsXV3YtdypTJECRA",
  "Channel_name": "Neural Network Console",
  "desc": "\"ソニーのDeep Learningの統合開発環境、Neural Network Consoleのチャンネルです。\\nDeep Learningの基礎について、およびNeural Network Consoleの使い方について解説しています。\\n\\nChannel for Neural Network Console, Sony's integrated development environment for Deep Learning.\\nLearn the basics of Deep Learning and how to use the Neural Network Console.\\n",
  "pub_date": "2017-08-06T06:34:24Z",
  "num_views": 2086907,
  "num_subscribers": 33600,
  "num_videos": 100},
...]]
```

Figure 1: Snapshot of data before processing

5 Methods

5.1 Data Cleaning

5.1.1 Dealing with Memory issues

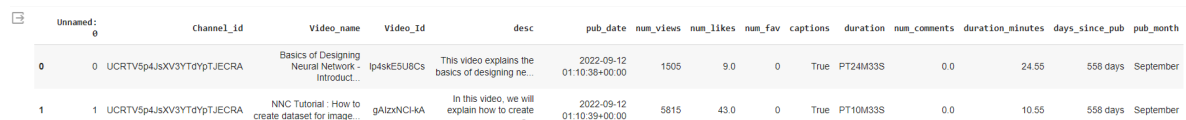
We have encountered memory issues while reading data of comments into working space. This has made us to read comments data into small chunks of 1000 size at a time and then looping through them

in order to convert into dataframe chunk by chunk. As we have pulled data from API, there is very minimal chances of missing values. Null values are present only in 'desc' column which is categorical in nature and doesn't need any imputation method.

5.1.2 Parsing Date Column

The date columns in the dataset were subjected to parsing and standardization to ensure consistency and facilitate meaningful temporal analysis. This process involved converting the date information from its original format into a standardized date format that is easily interpretable and compatible with analytical tools.

After all steps of data cleaning and processing, the dataframe looks like as follows:



	Channel_id	Video_name	Video_Id	desc	pub_date	num_views	num_likes	num_fav	captions	duration	num_comments	duration_minutes	days_since_pub	pub_month
0	UCRTV5p4JxXV3YTdYpTJECRA	Basics of Designing Neural Network - Introd...	Ip4skESU8Cs	This video explains the basics of designing ne...	2022-09-12 01:10:38+00:00	1505	9.0	0	True	PT24M33S	0.0	24.55	558 days	September
1	UCRTV5p4JxXV3YTdYpTJECRA	NVC Tutorial: How to create dataset for image...	gAlzxNCH-KA	In this video, we will explain how to create a...	2022-09-12 01:10:39+00:00	5815	43.0	0	True	PT10M33S	0.0	10.55	558 days	September

Figure 2: Snapshot of data after processing

6 Exploratory Data Analysis(EDA)

Our goal is to unearth meaningful stories embedded in the data and set the stage for more targeted analyses. We aim not only to uncover the inherent characteristics of the dataset but also to lay the groundwork for informed decision-making in subsequent phases of our data exploration.

Using the data insights acquired through the analysis of different features, the creation of various info graphics, and the incorporation of our existing knowledge of how Youtube gives search results and what criteria makes it possible.

6.1 Data Visualisations

6.1.1 Top channels based on views

In order to identify the most viewed channels within the realm of content related to our topic of interest [keyword (neural networks)], we conducted an analysis of the total number of views garnered by individual channels. This analysis provides valuable insights into the popularity and reach of content creators within our target domain. **3Blue1Brown** channel emerged as the front runner with a remarkable total view count, reflecting its broad appeal and engaging content offerings. Following closely, **Sebastain Lague** channel and **StatQuest with Josh Starmer** channel demonstrated strong viewer engagement and resonance, securing prominent positions in the ranking. **Deep Learning AI** and **Neural Networks & Deep Learning** channels rounded out the top five with consistent performance and audience affinity, highlighting their relevance and influence within the niche content landscape.

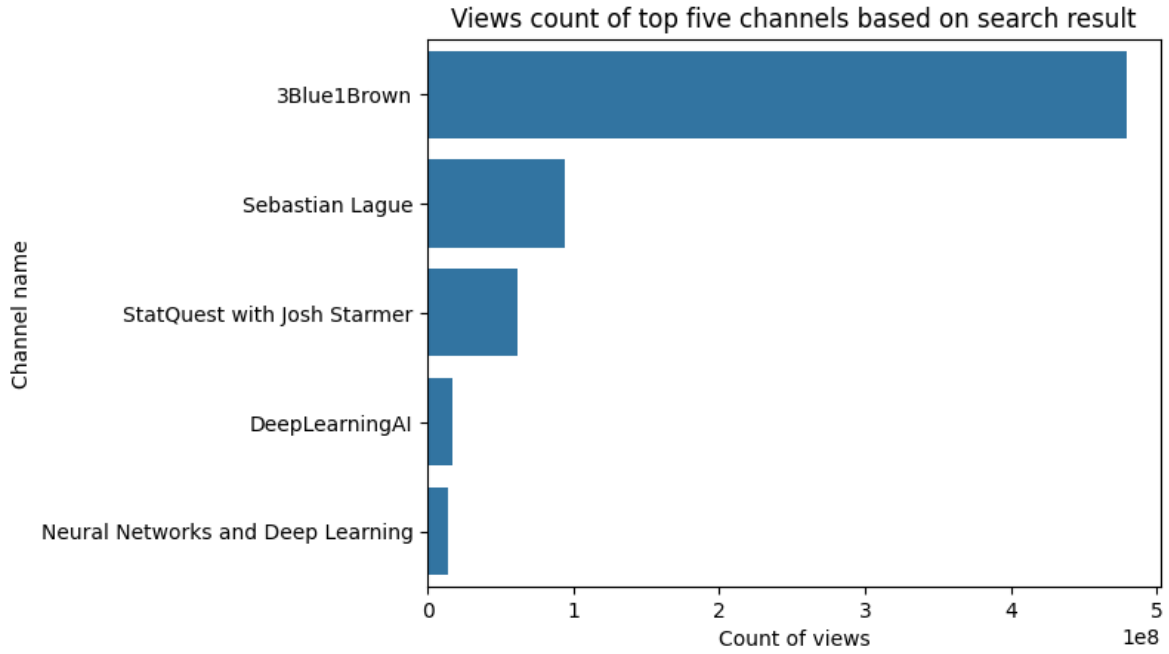


Figure 3: Top Five Channels Based on Number of Views: Insights

6.1.2 Insights into the Distribution of Content Creation Time Among Channels

In this pie chart, we have visualized the top 10 channels by distribution of their content creation time. This visualization provides a clear overview of the proportion of content creation efforts contributed by these channels. This analysis enhances our understanding of the content creation dynamics within our target domain and informs strategic decision-making regarding content partnerships, audience targeting, and engagement initiatives. This pie chart reveals **Alexander Amini** channel as the leading contributor, with 19% of the total content creation minutes, followed closely by **Killian Weinberger** channel with 13%. Remaining channels collectively contribute varying shares, ranging from 5% to 12%.

Percentage distribution of top 10 channels based on minutes of content creation

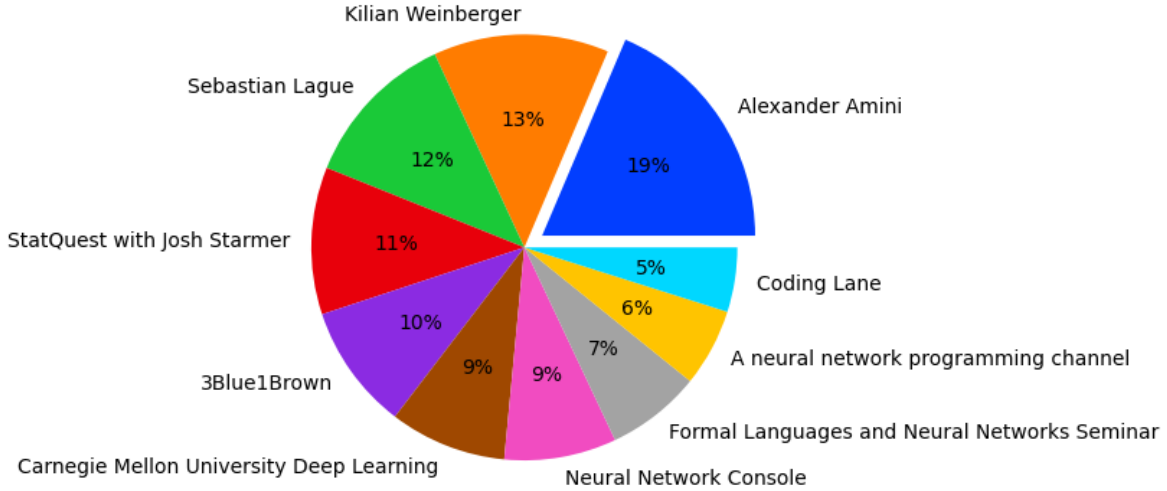


Figure 4: Breakdown of Top 10 Channels by Share of Content Creation Time

6.1.3 Analyzing the relationship between Video Duration and View Count

We utilized Seaborn's regplot function to generate a scatter plot examining the relationship between the duration of videos and the number of views they accrued. The data was filtered to include only videos with a view count of 100,000 or fewer views, ensuring a focused analysis on relatively less-viewed content. From this, we can infer that there is no preference for shorter videos among users which is proved by an almost straight line in the plot. It indicates very weak or no relationship between duration of video length and number of views.

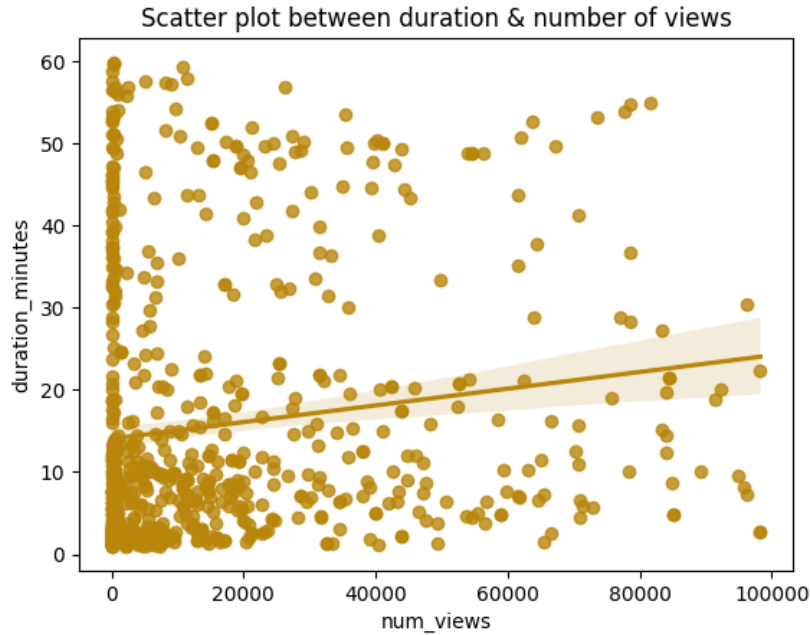


Figure 5: Linear regression between video length & views

6.1.4 Exploring videos durations Across Years: A Box Plot Analysis

In our analysis of the year-wise distribution of video duration, we utilized a box plot visualization to explore how the duration of videos has evolved over time. From the visualization, we observe variations in video duration across different years, with some years exhibiting a wider range of durations than others. It is evident that there is more range in the year **2018** spanning to larger length videos. This analysis provides valuable insights into trends and patterns in video duration over time, informing content creators and stakeholders about audience preferences and content consumption habits across different years.

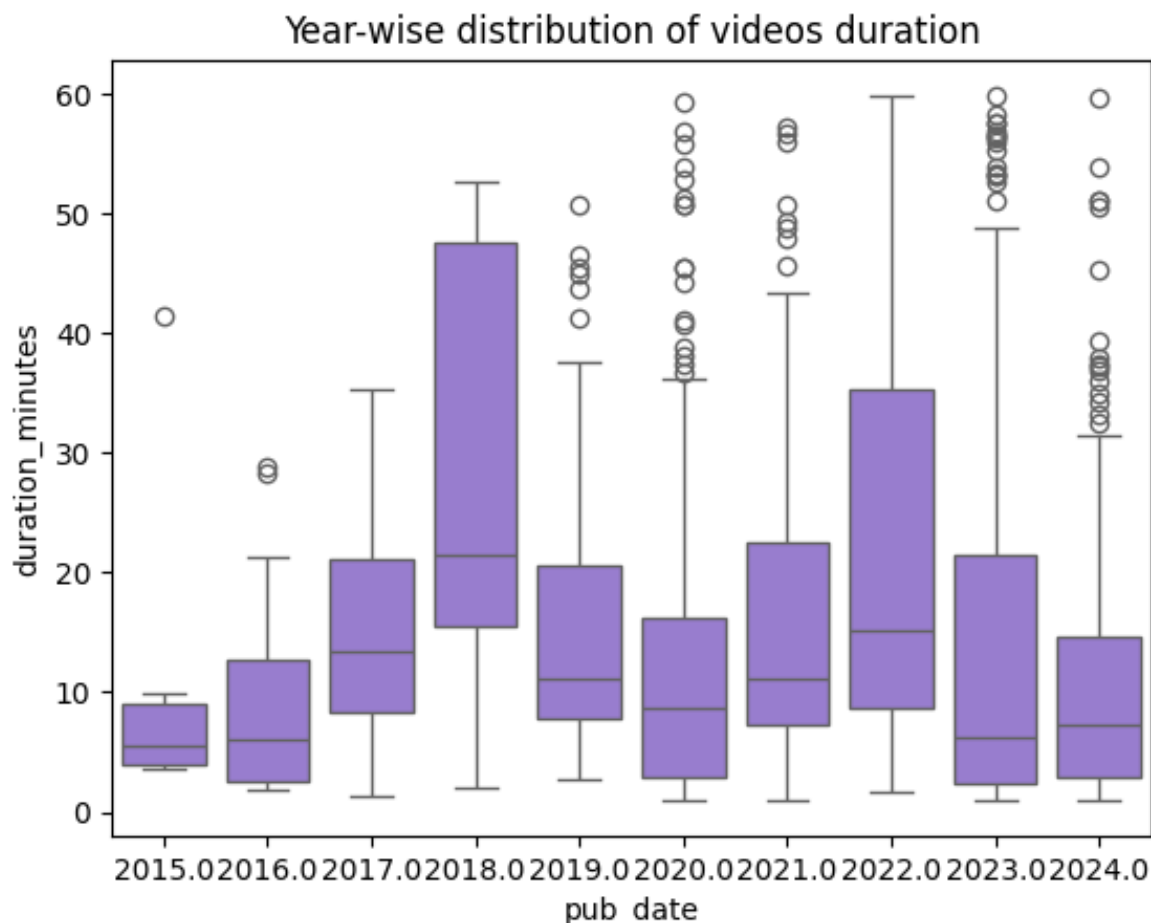


Figure 6: Year-wise Distribution of Video Duration

6.1.5 Distribution of Video Generation on a Monthly Basis

The histogram visualization depicts the month-wise distribution of the number of videos generated within our dataset. We can clearly see that there are more number of videos generated in early months of year compared to mid year months and again picking up pace in later months of the year.

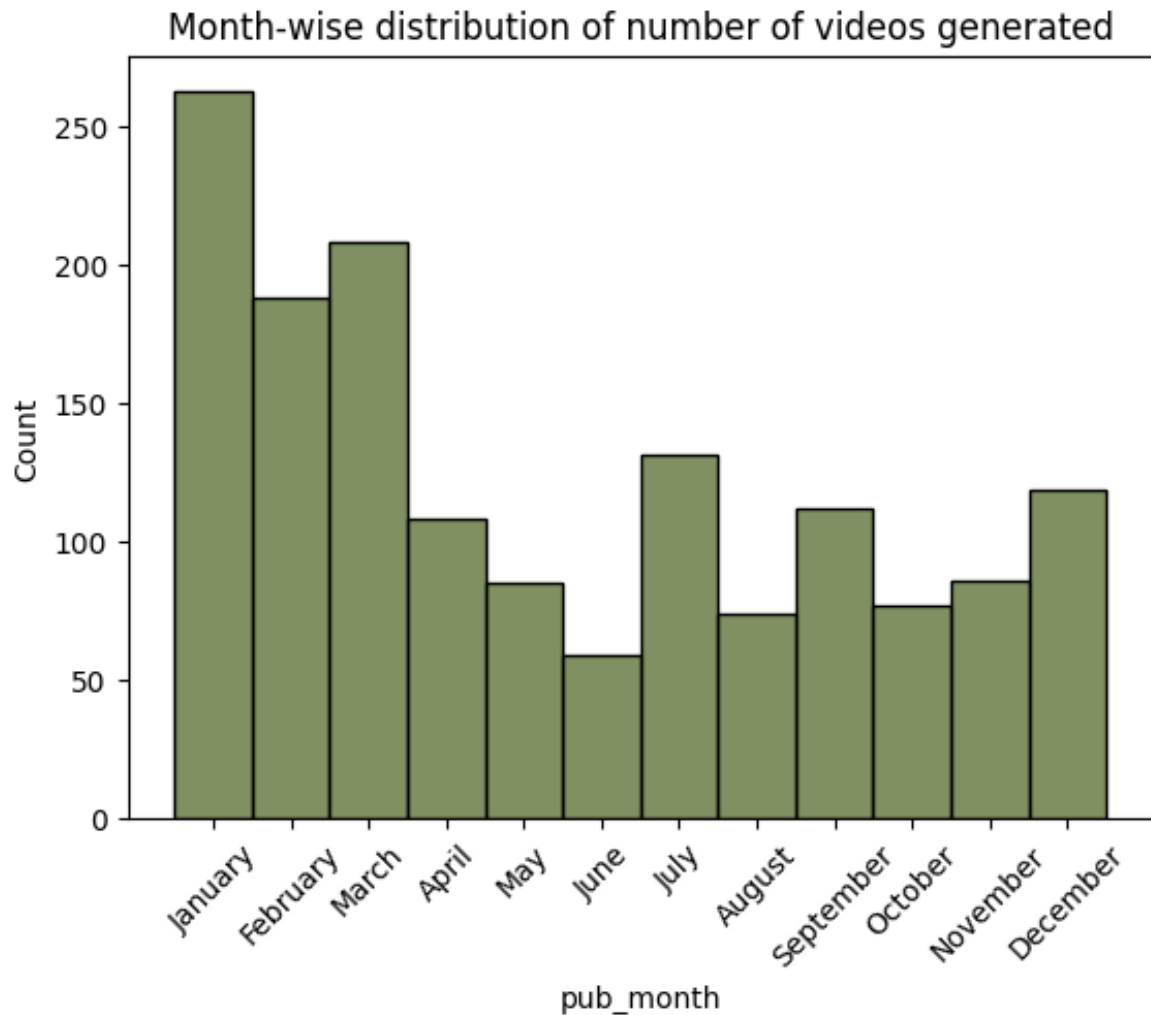


Figure 7: Month-wise distribution of number of videos generated

6.1.6 Exploring the Relationship Between Likes and Comments: Insights into Viewer Engagement Patterns

In our analysis, we utilized a scatter plot visualization as a tool to investigate the correlation between two key engagement metrics: the number of likes and the number of comments garnered by each video within our dataset. While it might seem intuitive that individuals would comment on a video they have watched or liked, the visualization allowed us to delve deeper into the underlying patterns of human behavior. Specifically, we sought to understand the propensity for individuals to engage with content by leaving comments, particularly in instances where they have expressed a positive or negative sentiment through liking or disliking the video. This observation underscores a fundamental aspect of human nature – the tendency to express opinions or reactions when experiencing an emotional response. Consequently, our analysis suggests a symbiotic relationship between the number of comments and likes, where heightened engagement in one metric often corresponds to increased activity in the other.

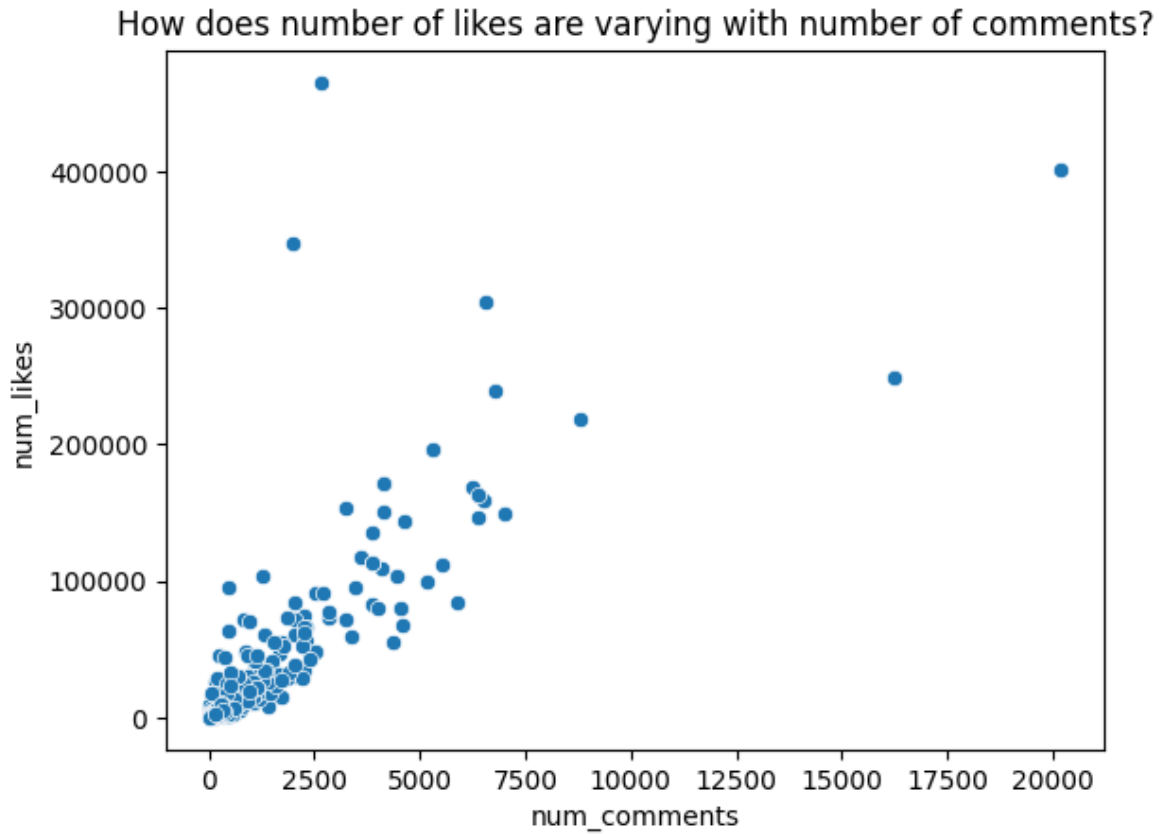


Figure 8: Scatterplot between comments and likes count

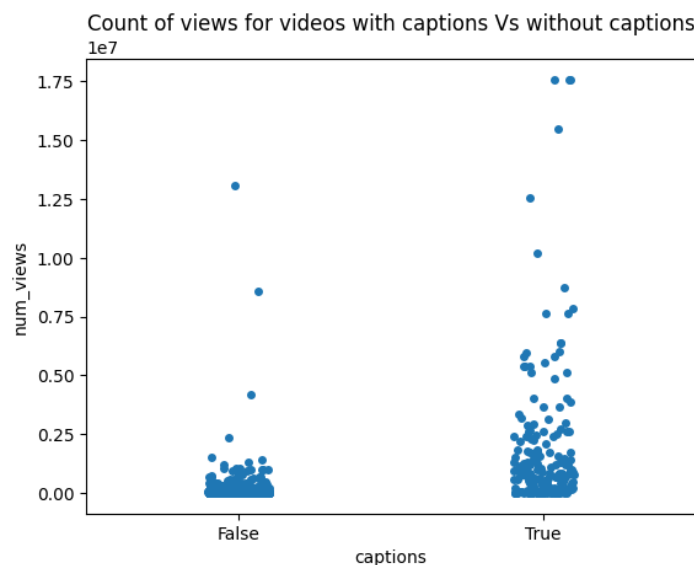
6.1.7 Word Cloud of videos descriptions

Upon analyzing the word cloud generated from the descriptions of videos related to neural networks, several key terms emerged prominently, providing valuable insights into the prevalent themes and topics within this domain. Notably, words such as "machine learning," "deep learning," "beginner," "topic," "AI," and "specialization" recurred frequently, indicating their significance and prevalence in the descriptions of these videos. The prominence of these terms suggests that the videos likely cover introductory concepts and discussions surrounding machine learning and deep learning, catering to individuals with varying levels of expertise, including beginners seeking to explore the topic further. Additionally, the mention of "AI" underscores the intersectionality of neural networks with artificial intelligence, reflecting the broader context in which these videos are situated. Overall, the word cloud provides valuable insights into the dominant themes and focal points of the videos, guiding viewers and content creators alike in navigating the expansive landscape of neural network-related content on the platform.



6.1.8 Analyzing Viewer Engagement Based on Captions Presence in Videos

The strip plot visualization showcases the distribution of views for videos categorized based on the presence or absence of captions. The x-axis represents the captions status, with "Yes" indicating videos with captions and "No" indicating videos without captions, while the y-axis depicts the corresponding count of views. The plot provides a comparative analysis, allowing us to discern any potential differences in viewership between videos with and without captions. The title, "Count of views for videos with captions Vs without captions," succinctly summarizes the main focus of the visualization, facilitating clear interpretation. This visualization serves as a valuable tool for exploring the impact of captions on viewer engagement and accessibility within our dataset of videos, offering insights that can inform content creators and stakeholders in optimizing their video content strategies. We can see that videos having captions have more number of views.



6.1.9 Exploring Video Duration Distribution Over Time with Captions Analysis

This catplot visualization employs violin plots to explore the distribution of video durations across different years, with a distinction made between videos with and without captions. These plots display the distribution of video durations within each year, with the width indicating the density of videos at different duration values. The hue parameter separates the violin plots based on the presence or absence of captions, allowing for a comparative analysis of video durations between the two categories. The mean of duration of videos with captions is relatively higher than those videos without captions. By utilizing violin plots, this visualization effectively illustrates the variability and distribution of video durations over time, while also highlighting any discernible patterns or differences in duration between videos with and without captions.

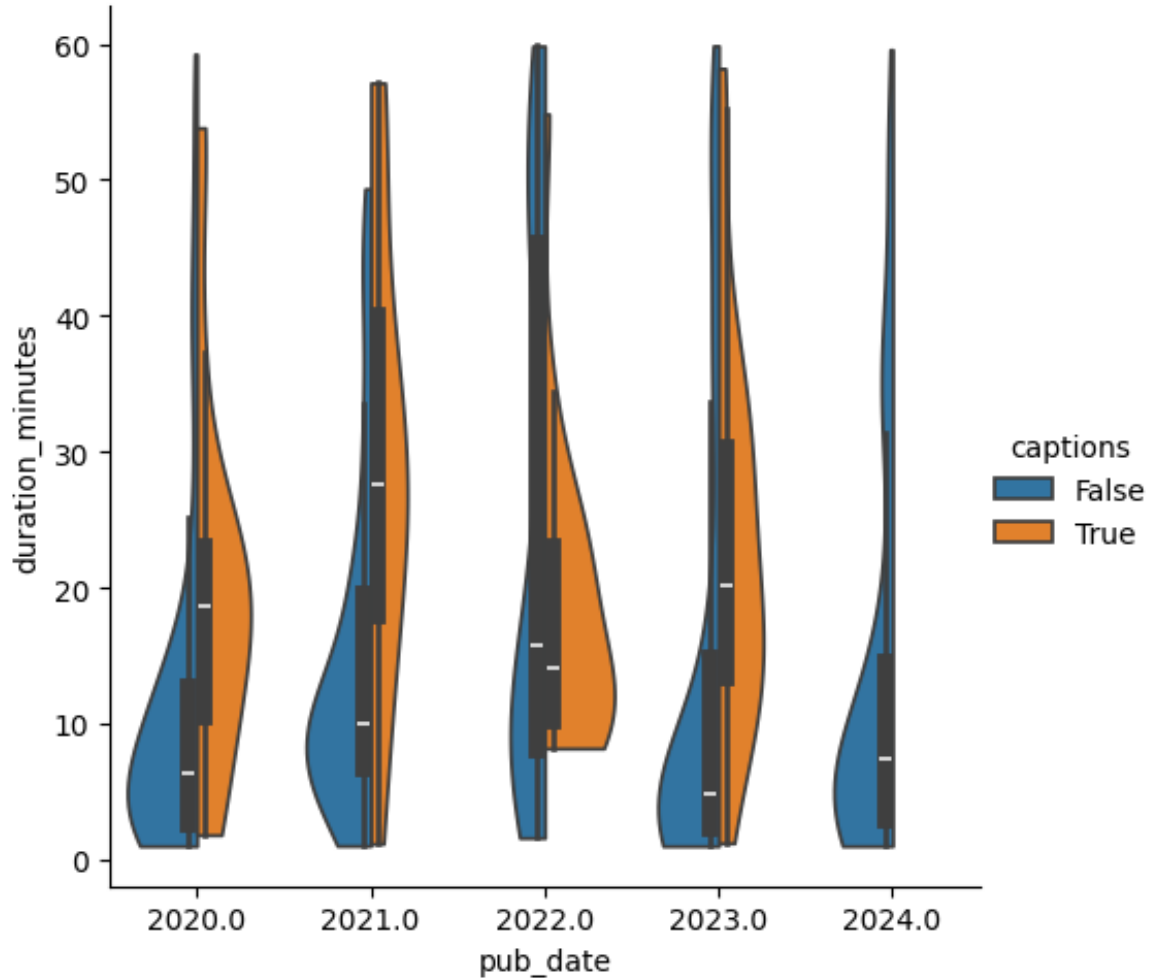


Figure 11: Plot of number of views against caption

6.1.10 Determining association between Duration and Number of Likes

This jointplot reveals that there is no clear conclusion to draw but this plot helps in understanding independence between these two features. Also, it is not very clear to interpret the relationship between these features. This visualization serves as a valuable tool for understanding how video duration influences viewer engagement, as measured by the number of likes, while also considering the presence of captions as a potential factor influencing viewer interaction.

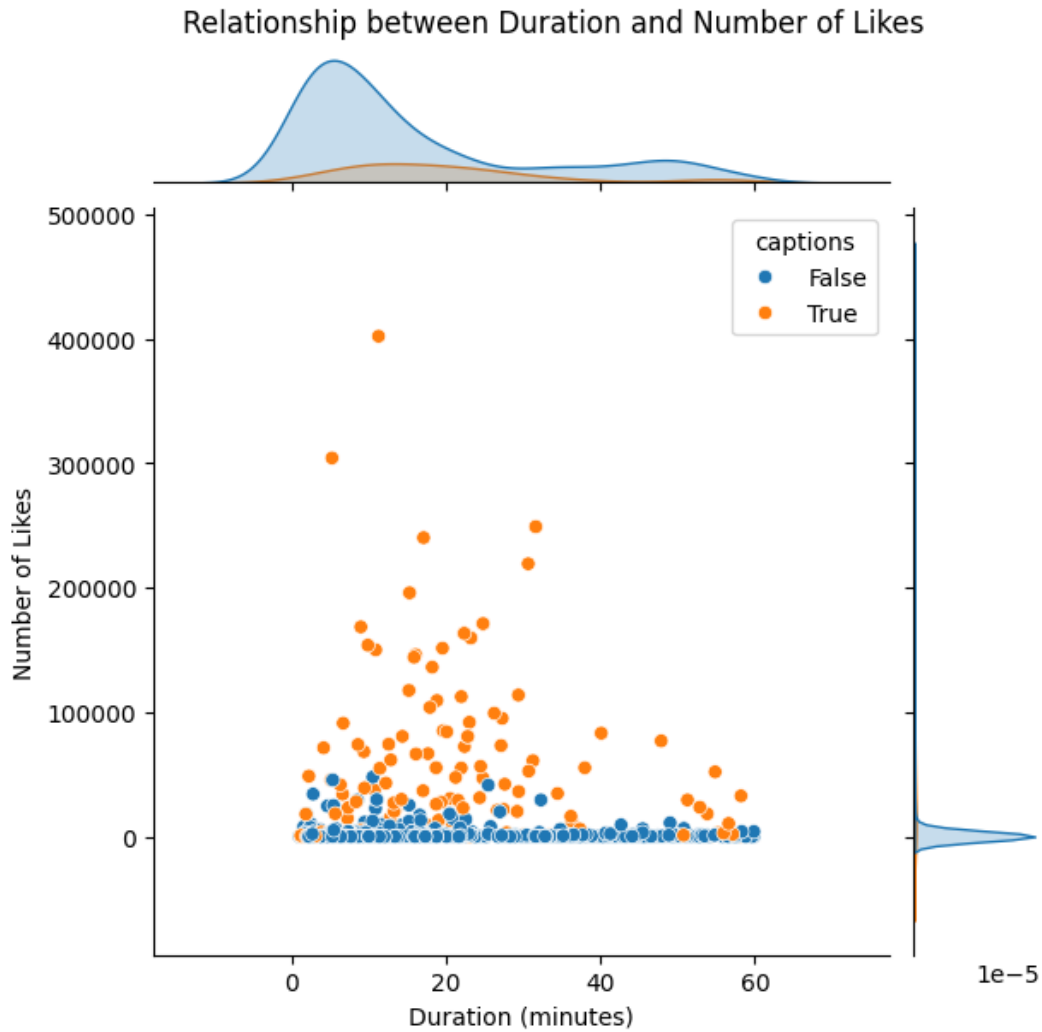


Figure 12: Jointplot of number of likes and duration of videos

6.1.11 Exploring Univariate Distribution of Video Durations

The swarm plot visualization provides a univariate distribution of the duration of videos within our dataset. Each data point on the plot represents an individual video, with the x-axis depicting the duration of the videos in minutes. From the visualization, we observe a wide range of durations, spanning from shorter to longer videos. Additionally, the swarm plot highlights the density of videos at lower duration values, with clusters of data points indicating areas of higher prevalence. Overall, this visualization offers valuable insights into the distribution and variability of video durations within our dataset, informing our understanding of the content landscape and helping to identify potential trends or patterns in video duration.

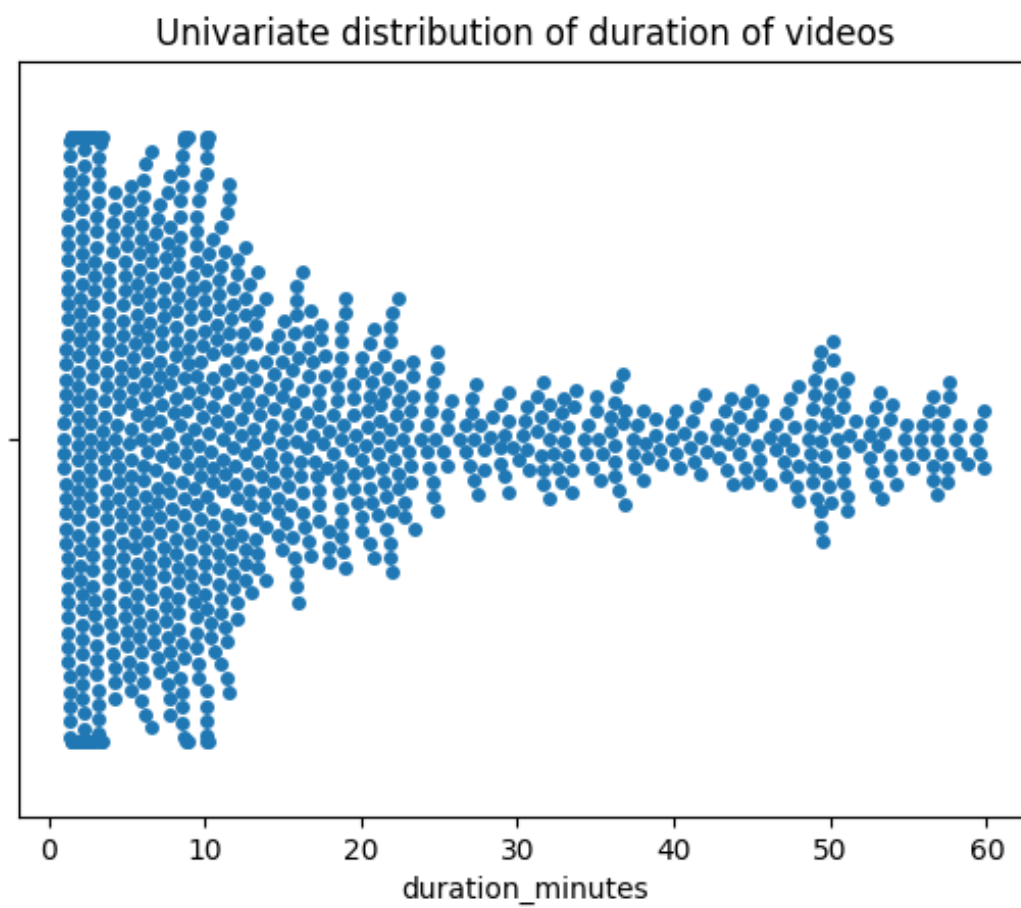


Figure 13: Swarmplot of number of likes and duration of videos