

DSA210

Term Project

Analysis of Youtube Data

Alp Çetintaş - 32550



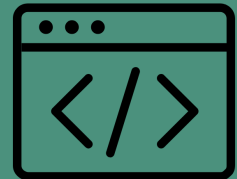
Index & Content

1. Motivation - Why?
2. Data Source - Preprocessing
3. Data Visualization and Exploratory Data Analysis (EDA)
4. Hypothesis Testing
5. Machine Learning Model
6. Findings
7. Limitations and Future Work

Motivation - Why?

- Most used social media platform
- Reflects my interests w/ videos
- Uncover my habits

- What type of content do I watch
- Percentage of likes over videos viewed
- Is there any preferred/favored content creators
- Daily/Weekly watching habits
- My subscriptions
- Search history



Data Source - Preprocessing

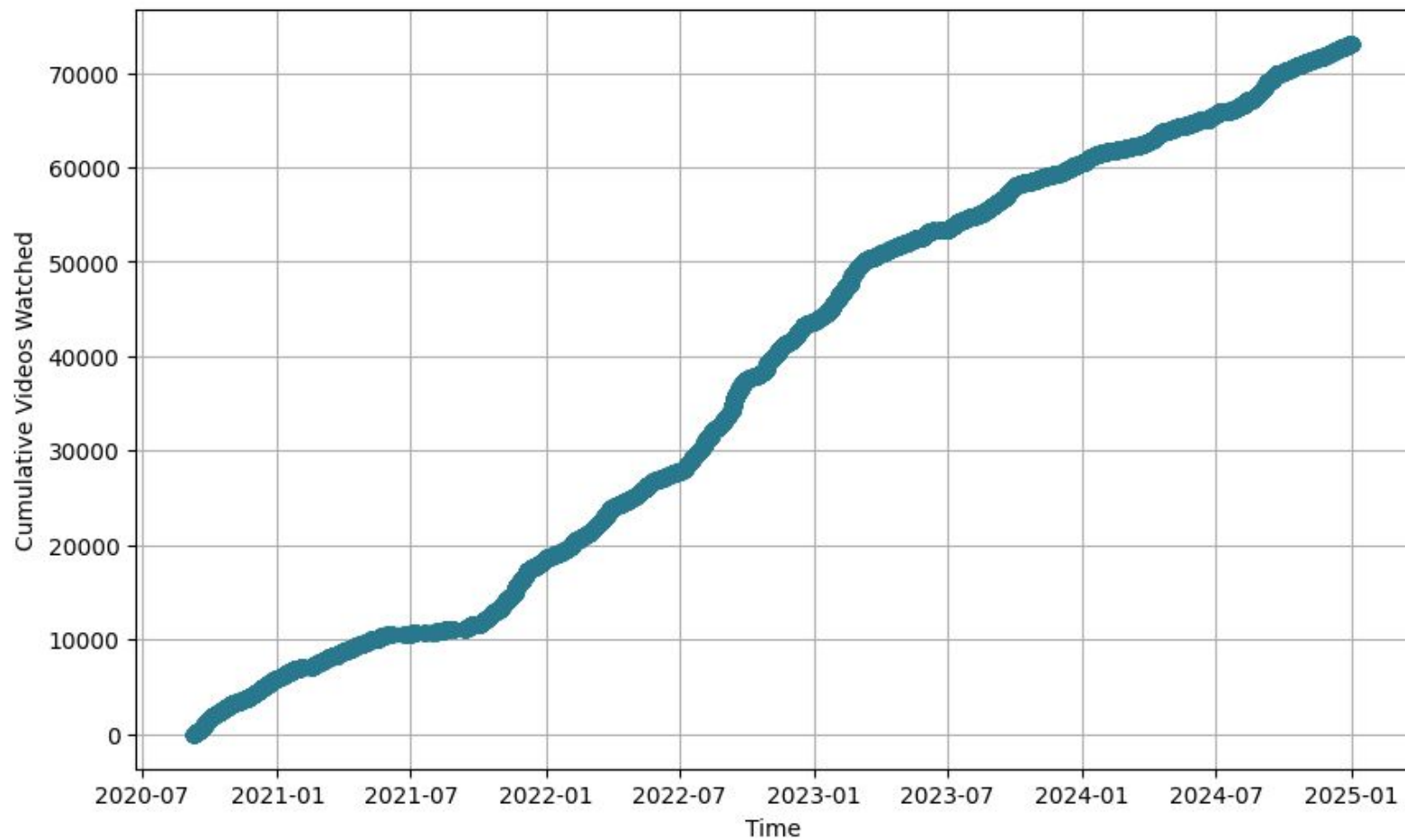
- Google Takeout
 - Watch History
 - Search History
 - Subscriptions
 - Comments
 - Playlists



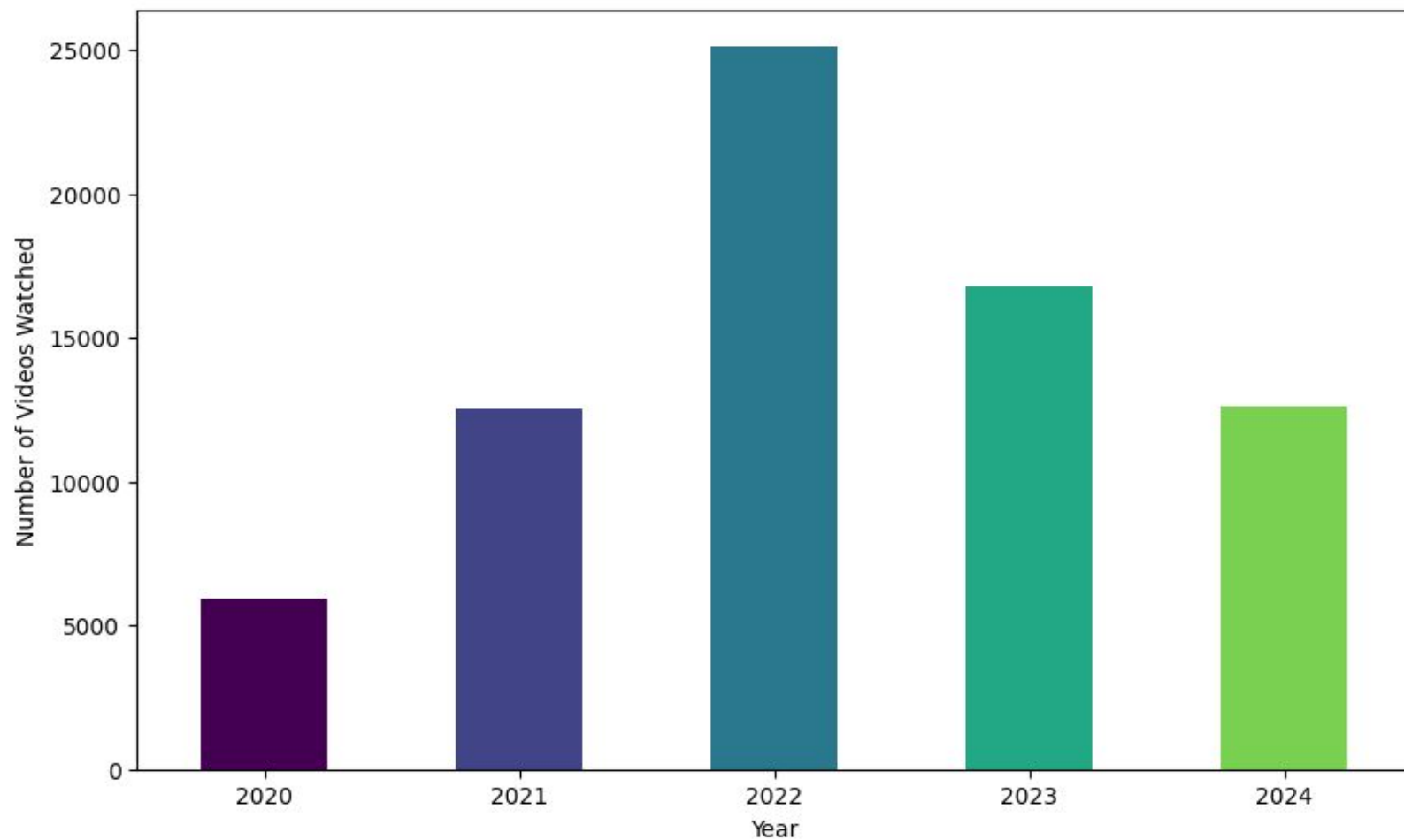
- Preprocessing
 - Merging two accounts
 - Unavailable Title Names
 - Converting time column - date_time
 - Sorting in time order

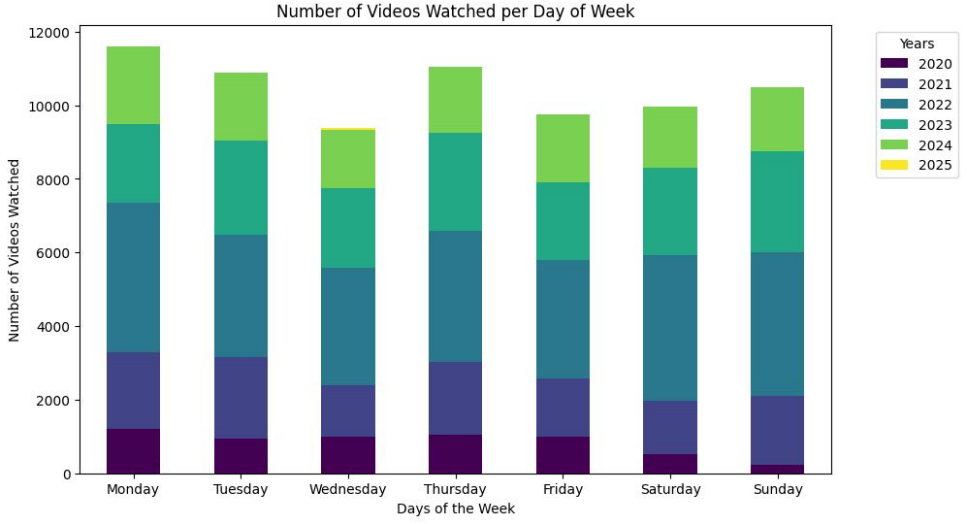
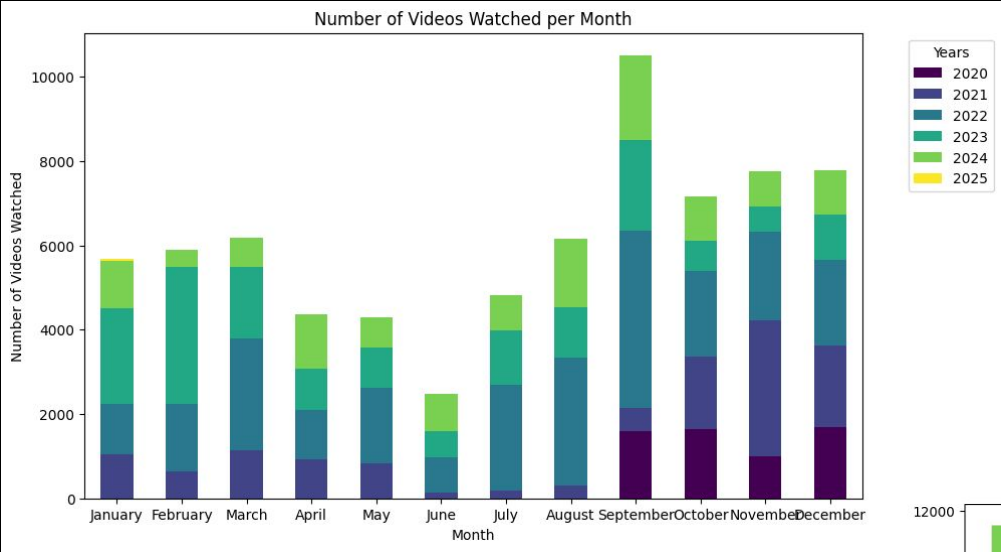
```
{  
  "header": "YouTube",  
  "title": "WW1 and Self-Inflicted Wounds adlı videoyu izlediniz",  
  "titleUrl": "https://www.youtube.com/watch?v\u003dcf9eDFlawL0",  
  "subtitles": [{  
    "name": "Johnny Johnson",  
    "url":  
      "https://www.youtube.com/channel/UCg7Q08KKOdjSrSnXcUL00Jw"  
  }]  
}
```

Total amount of Videos Watched

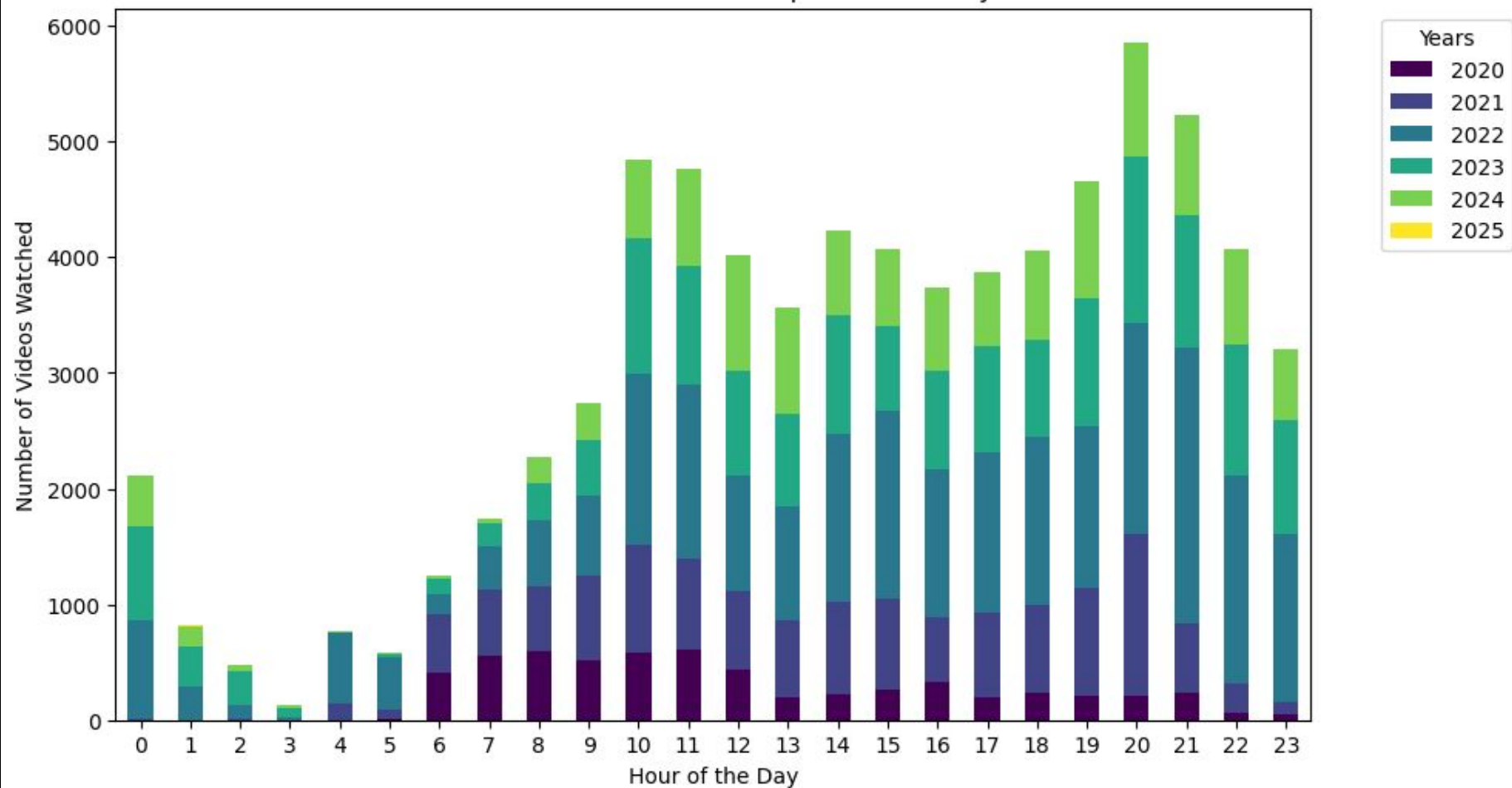


Number of Videos Watched Per Year

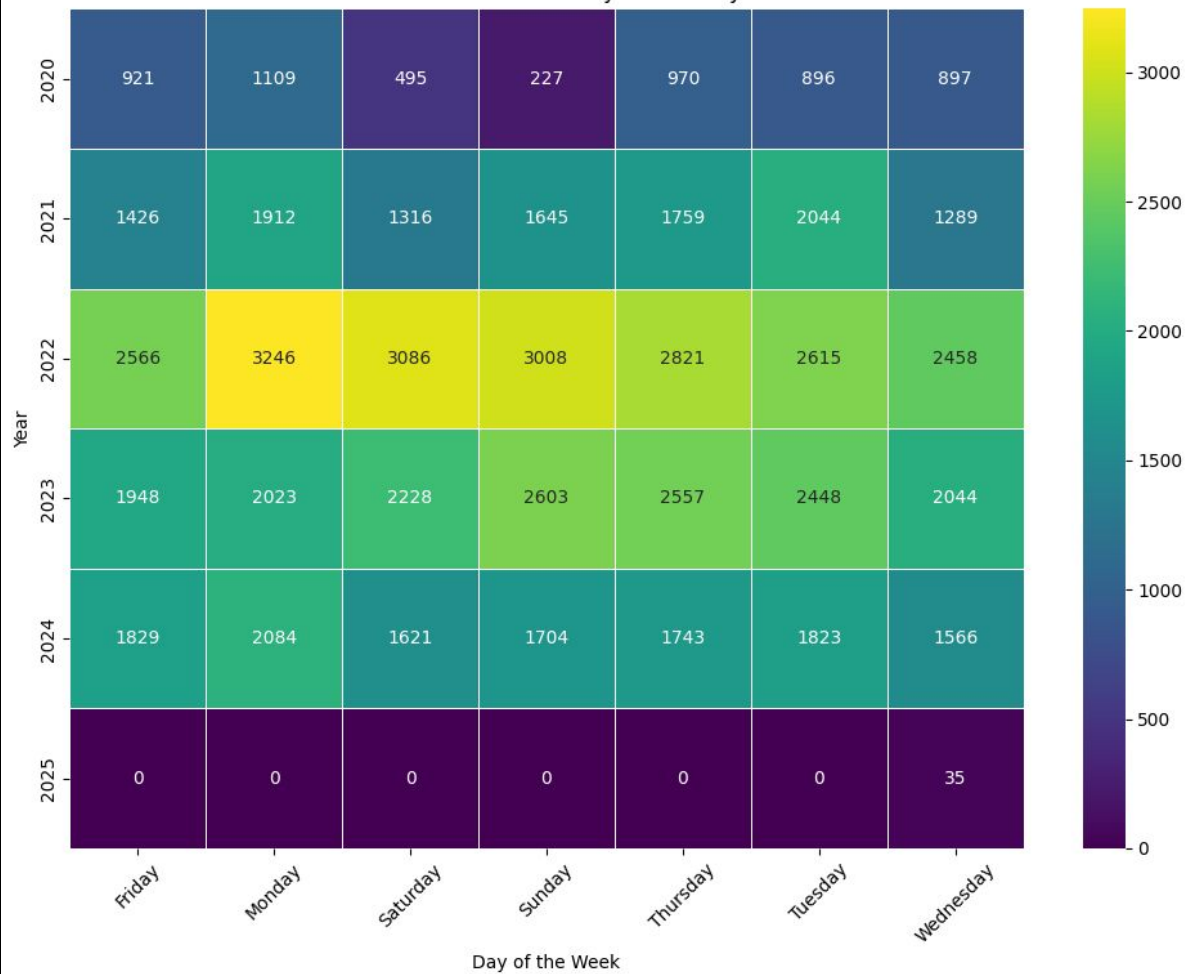




Number of Videos Watched per Hour of Day



Total Videos Watched Per Day of Week by Year



Hypothesis Testing

Hypothesis: The top 5 channels account for more than 10% of the watch time

H0: The top 5 channels make up 10% or less of the videos watched

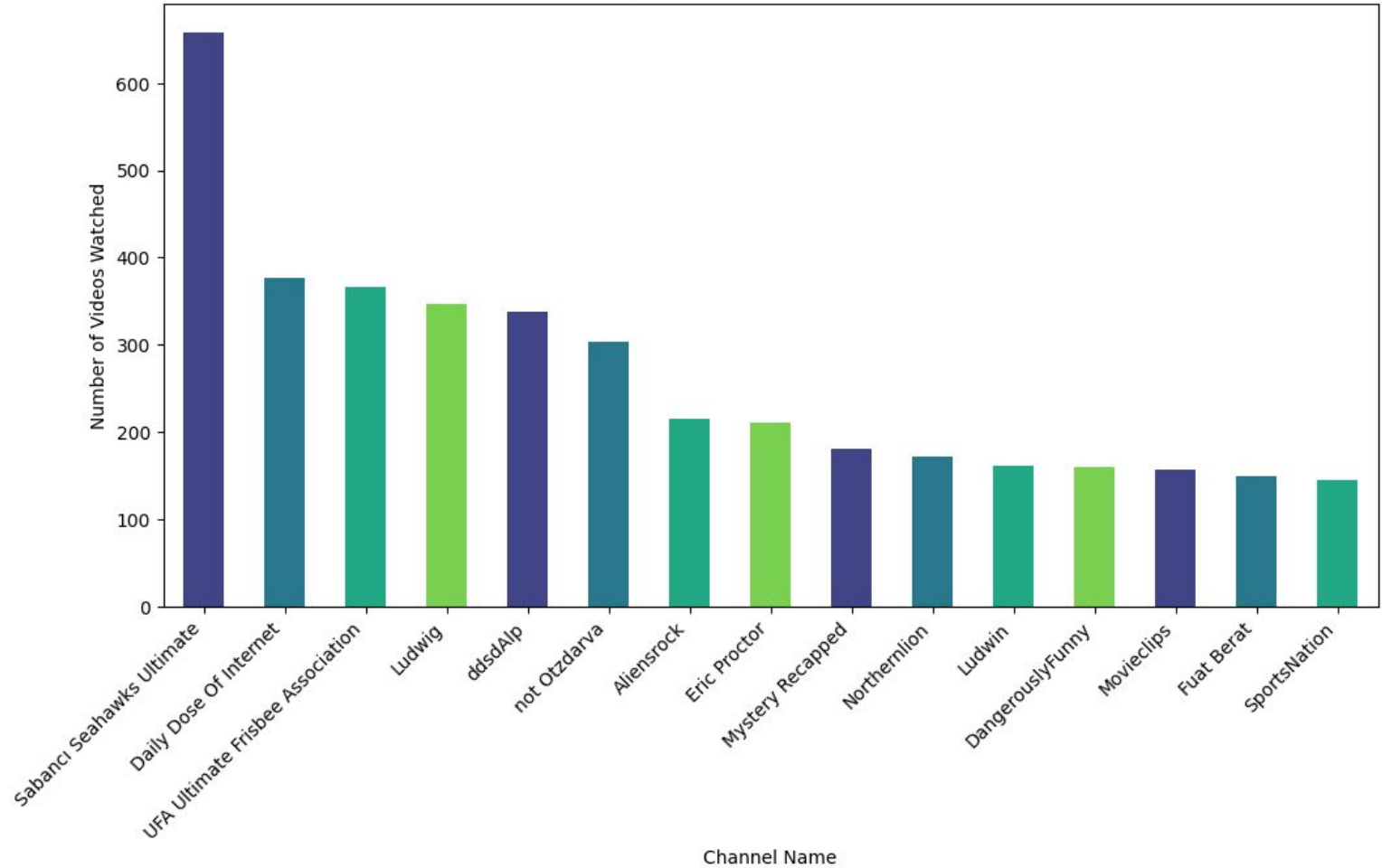
HA: The top 5 channels make up more than 10% of videos watched

$$H_0 : p \leq 0.10\%$$

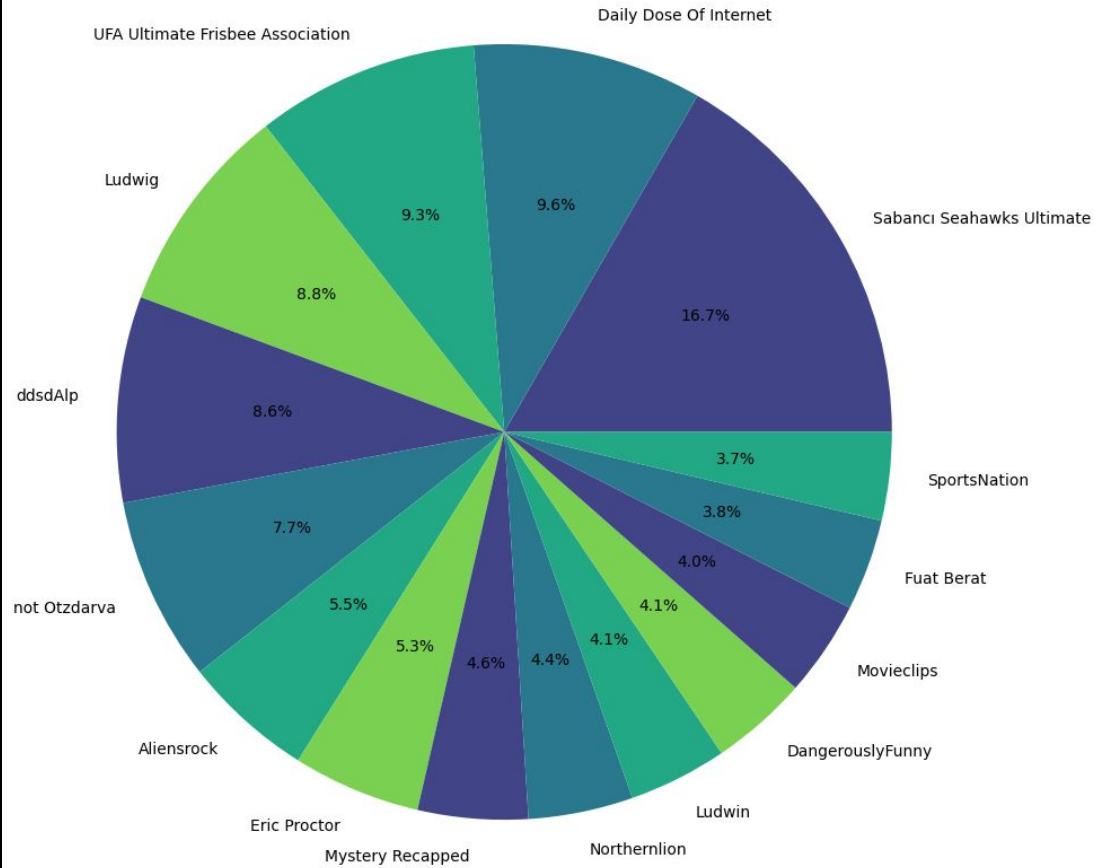
$$H_a : p > 0.10\%$$

(The value was 50% at first but I realized it was too much and lowered it down)

Top 15 Most Watched Channels



Top 15 Most Watched Channels

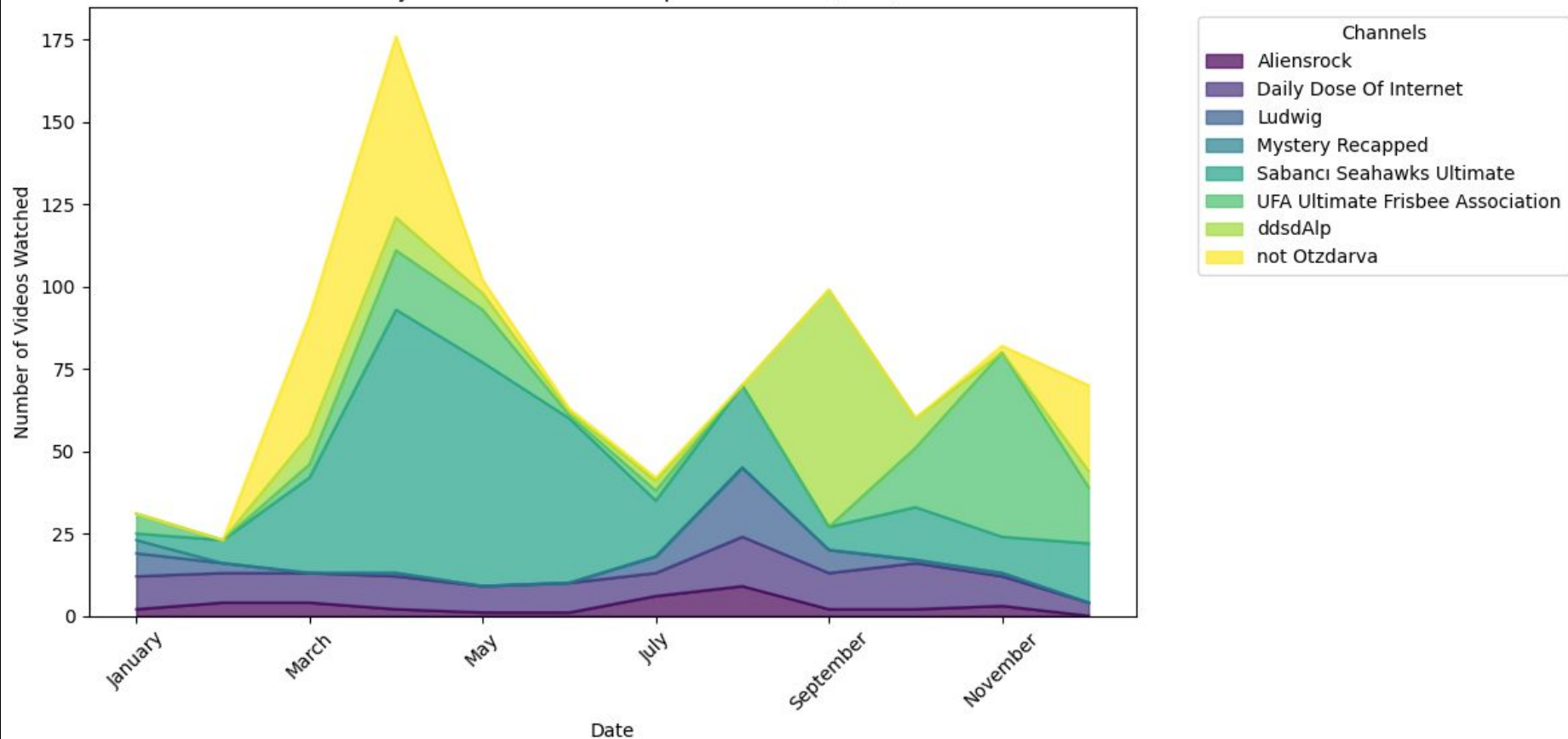


The chart displays the monthly video watch counts for the top 10 channels in 2023. The total watch count fluctuates significantly throughout the year, with major peaks occurring in March, June, and September/October. The channel 'not Otzdarva' (yellow) is a significant contributor to the total watch count, especially in the first half of the year. Other prominent contributors include 'UFA Ultimate Frisbee Association' (light green), 'Northernlion' (teal), and 'Aliensrock' (dark purple).

Month	Aliensrock	Daily Dose Of Internet	Eric Proctor	Ludwig	Mystery Recapped	Northernlion	Sabancı Seahawks Ultimate	UFA Ultimate Frisbee Association	ddsdAlp	not Otzdarva
January	18	10	5	5	10	10	10	10	20	25
February	5	10	5	5	10	10	10	10	20	10
March	10	10	5	5	10	10	10	10	20	70
April	5	10	5	5	10	10	10	10	20	10
May	5	10	5	5	10	10	10	10	20	10
June	5	10	5	5	10	10	10	10	20	10
July	5	10	5	5	10	10	10	10	20	10
August	5	10	5	5	10	10	10	10	20	10
September	20	10	5	5	10	10	10	10	20	10
October	15	10	5	5	10	10	10	10	20	10
November	10	10	5	5	10	10	10	10	20	10

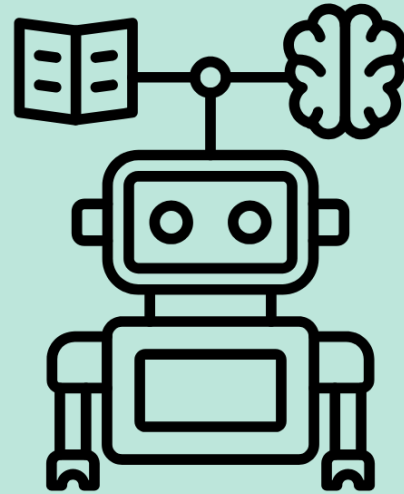


Monthly Videos Watched for Top 10 Channels (2024)



Machine Learning Model

- Creating a ML model to predict how many videos I'll watch in a day
- Features of the ML model
 - day_of_week
 - month
 - is_weekend
 - lag



Findings

- Almost never comment
- Quick change of interests and its correlation with my consumed content
- Not sticking to one/a group of channels
- My own channel is in my top 15
- The year 2022
- Active days, heatmap, active months



Limitations and Future Work

- I couldn't get the packets for the textual analysis to work
- The dataset not having the watchtime of the video
- The confusion and problem this creates
- Future work if watch time was available
 - Avg. watch time of a video
 - Avg. percentage of the videos watched
 - Time spent on categories (gaming, education, coding, cars)
 - Time spent on channels