

# **YouTube Popularity Scraper and Engagement Prediction**

By: Gabriel Silva

Github: [https://github.com/MrYo10/YouTube\\_Popularity\\_Scraper/tree/main](https://github.com/MrYo10/YouTube_Popularity_Scraper/tree/main)

**Description of the Project:**

The project aimed to predict Youtube video stats which included popularity and engagement using two data sources.

- Scraped Data Which collected straight from the YouTube webpages such as the titles, descriptions, tags, category, etc.
- API data which collected data using the Youtube Data API v3 which are the official statistics such as the view count, likes, comment count, and the channels metrics.

The workflow consisted of the data collection, cleaning & preprocessing, feature engineering, and the machine learning and training and the evaluation for both data sets. The Random Forest Regressor and the optional XGBoost models were trained to predict either the engagement rates which are the (Likes + comments) / views from the api data or the View count from the scrapped data.

## How to use:

- Install Dependencies: `pip install -r requirements.txt`
- Collect data:
  - `python src/collect_api.py --resume`
  - `python src/collect_scraped.py --resume`
- Preprocess: `python src/preprocess.py`
- Feature extraction: `python src/features.py`
- Train Models:
  - `python src/train_api.py`      # API-based
  - `python src/train_scraped.py`      # Scrapped-based
- Evaluate and Visualize: `python src/eval_visualize.py`

The data will appear under:

Data/ interim/ , data/ processed – for the clean and feature data

Models/ - trained models + metrics

Reports/ - plots and comparison figures

## **Training and inferencing:**

Each Training script splits the data in the standard 80% train and the 20% test and it fits to a Random Forest Regressor:

Model 1: Is trained on scrapped features to predict the  $\log_{1p}(\text{viewCount})$

Model 2: is trained on API features to predict `engagement_rate`

The prediction can be done by loading the saved .joblib mode and calling. `predict(new_features)`

## **Data Collection:**

Scraped collector: `collect_scraped` using BeautifulSoup + requests to extract title, description, tags, uploader, category, publish date, and view counts from video URLs

API collector: `collect_api.py` calls the Youtube Data API v3 using videos IDs to fetch `viewCount`, `like count`, `comment count`, `channel subscribers`, and `tags`

Tools used: Python 3.12, pandas, Numpy, BeautifulSoup4, requests, scikit-learn, xgboost, matplotlib.

Collected attributes:

Source	Attributes
Scraped	video_id, url, title, description, uploader, publish_date, category, tags, views (when available)
API	video_id, url, title, description, uploader, publish_date, category, tags, views (when available)

Number of samples:

Scraped: 10,199 rows

API = 2,719 rows

## Data Preprocessing:

### Data cleaning:

Handled bad CSV lines and missing columns

Converted numeric fields like view counts and like count to floats

Parsed ISO 8601 public dates

Removed duplicates by video\_id

Replaced invalid / zero view counts with NaN

Clipped the videos with extreme views to 99.9<sup>th</sup> percentile to limit the outliers.

### Feature Engineering:

Created new columns for:

Duration\_min

time\_since\_upload\_days

tags\_count, desc\_length, had\_links\_in\_desc

Upload\_year, upload\_month, upload\_dow, upload\_hour

Target Variables:

Engagement\_rate = (likes + comments)/views

Target\_views\_log =  $\log_{1p}(\text{viewcount})$

## Model Development and Evaluation

Train / test Partition

Both the data sets are split 80/20 randomly

Model 1- Scraped Data

Algorithm: Random Forest Regressor

Input features: title n-grams, duration\_min, upload time features, text statistics.

Target: target\_views\_log

Train size = 3500 samples with valid views

Performance

train  $R^2 = 0.00$

Test  $R^2 = 0.00$

MAE = 2.5 log scale  $\rightarrow$  = x12 error factor in raw views

Conclusion: The scraped text alone is weakly predictive; channel-level features are missing.

Model 2: API data

Algorithm: Random Forest Regressor (+ optional XGBoost)

Input features: duration\_min, time\_since\_upload\_days, channel stats, title, description features.

Target: target\_engagement

Train size: about 2700 rows

Performance:

Train  $R^2 = 0.91$  MAE = 0.004 RMSE = 0.007

test  $R^2 = 0.41$  MAE = 0.0108 RMSE = 0.0159

This explains the 41% of variance in engagement, which is a strong result for YouTube Behavior data.

## Feature Importance

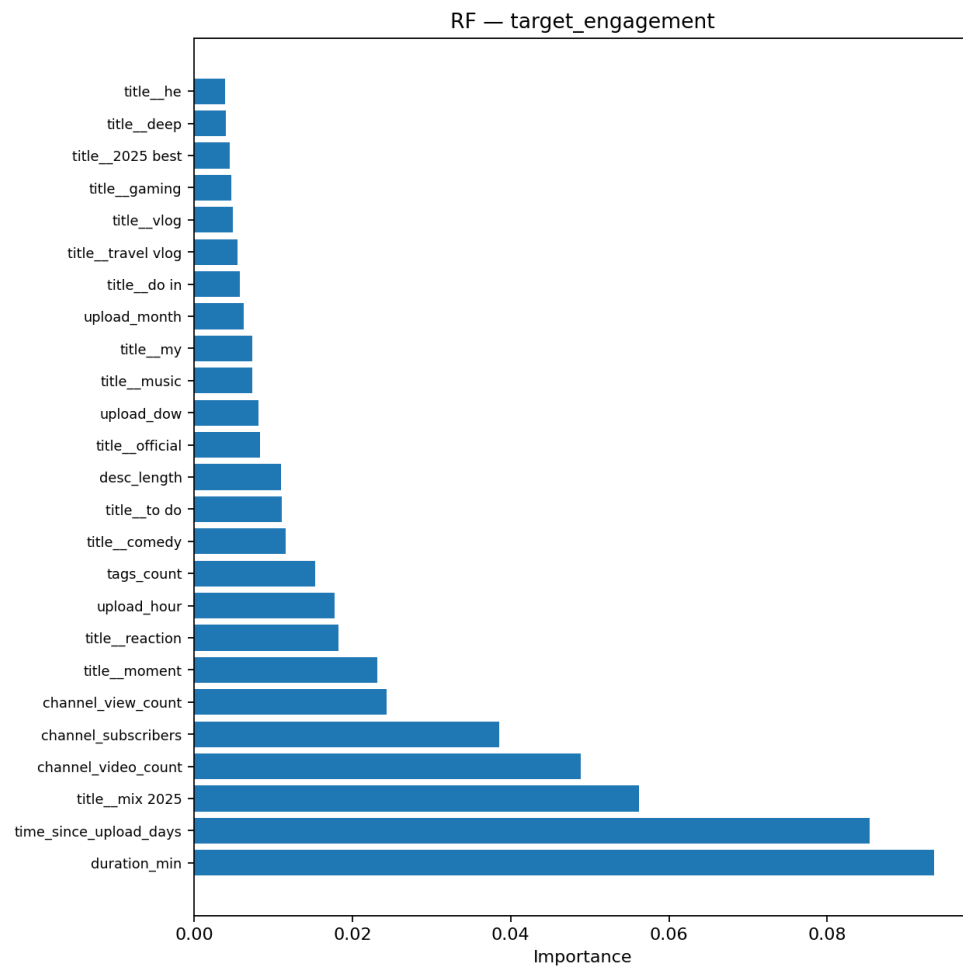
Technique: Mean Decrease in impurity (Random Forest feature importances)

Top predictors for API model:

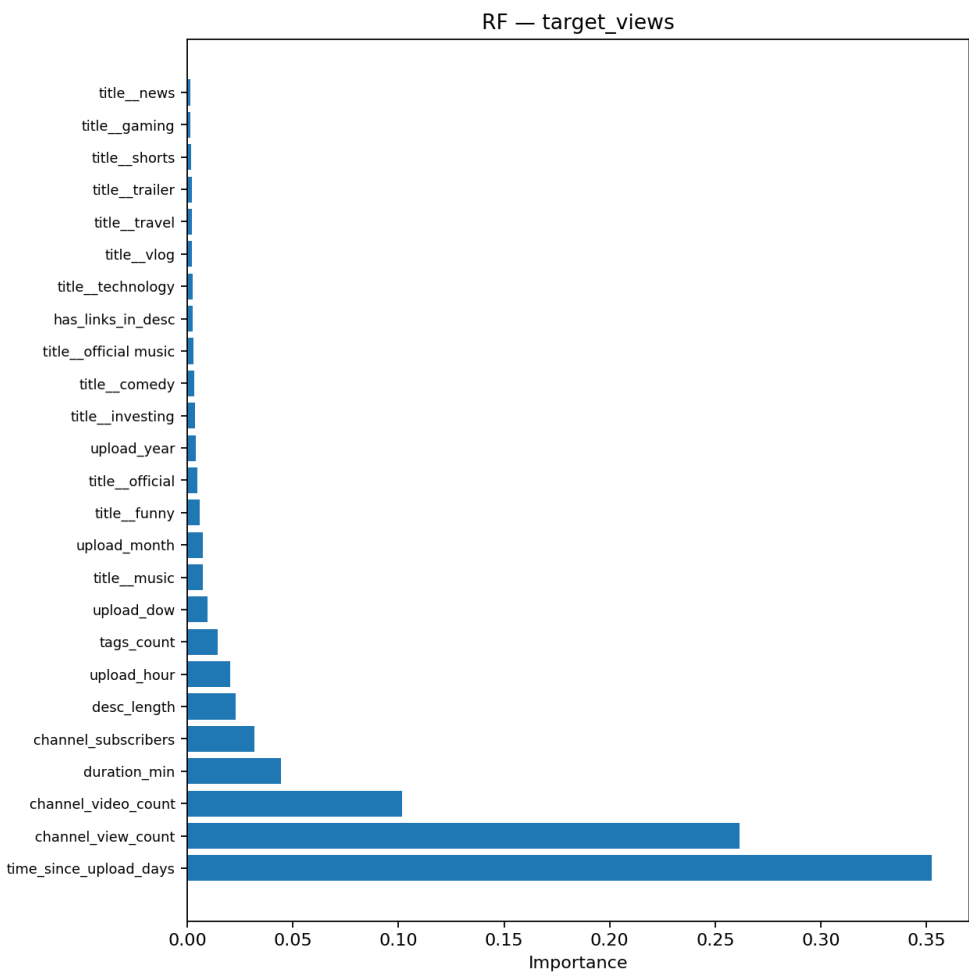
Channel\_subscribers, viewCount\_clipped, time\_since\_last\_upload\_days, desc\_length, tags\_count, and the specific title tokens like official, remix and or live.

Feature importance plots saved as reports/fi\_api.png

## Random Forest Feature Importance:

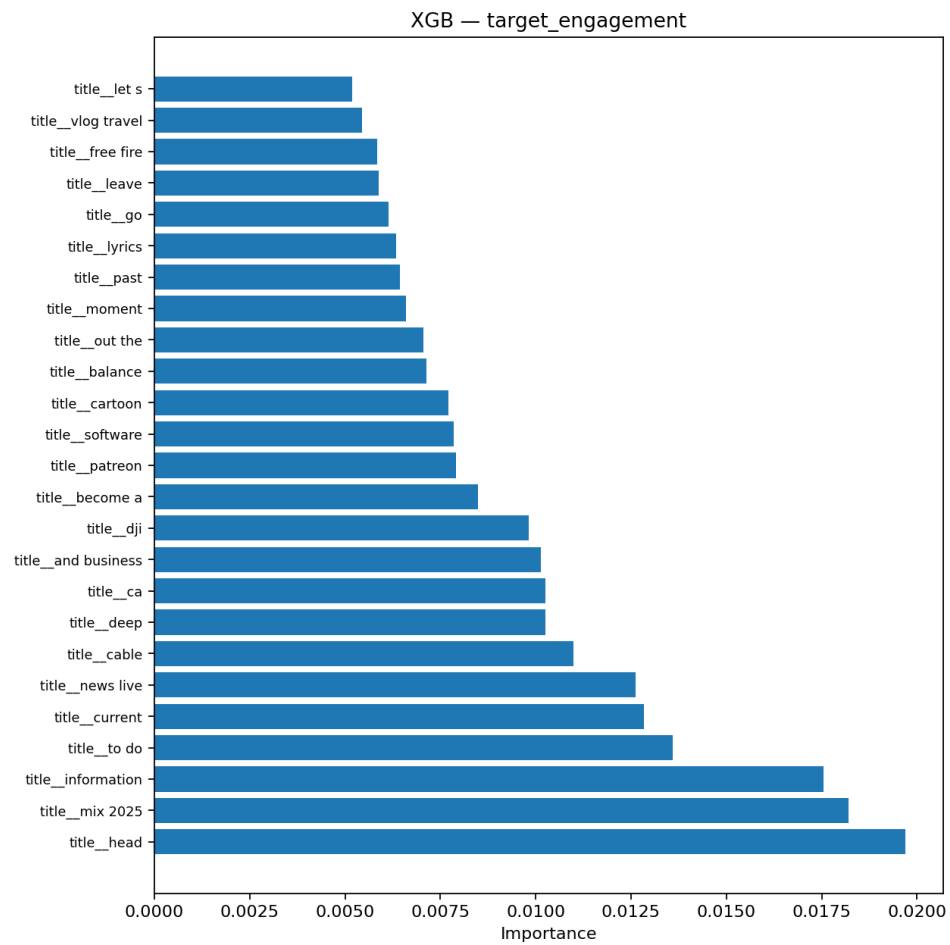


**Random Forest feature importance for View Prediction:**



## XGBoost Feature importance for Engagement Rate

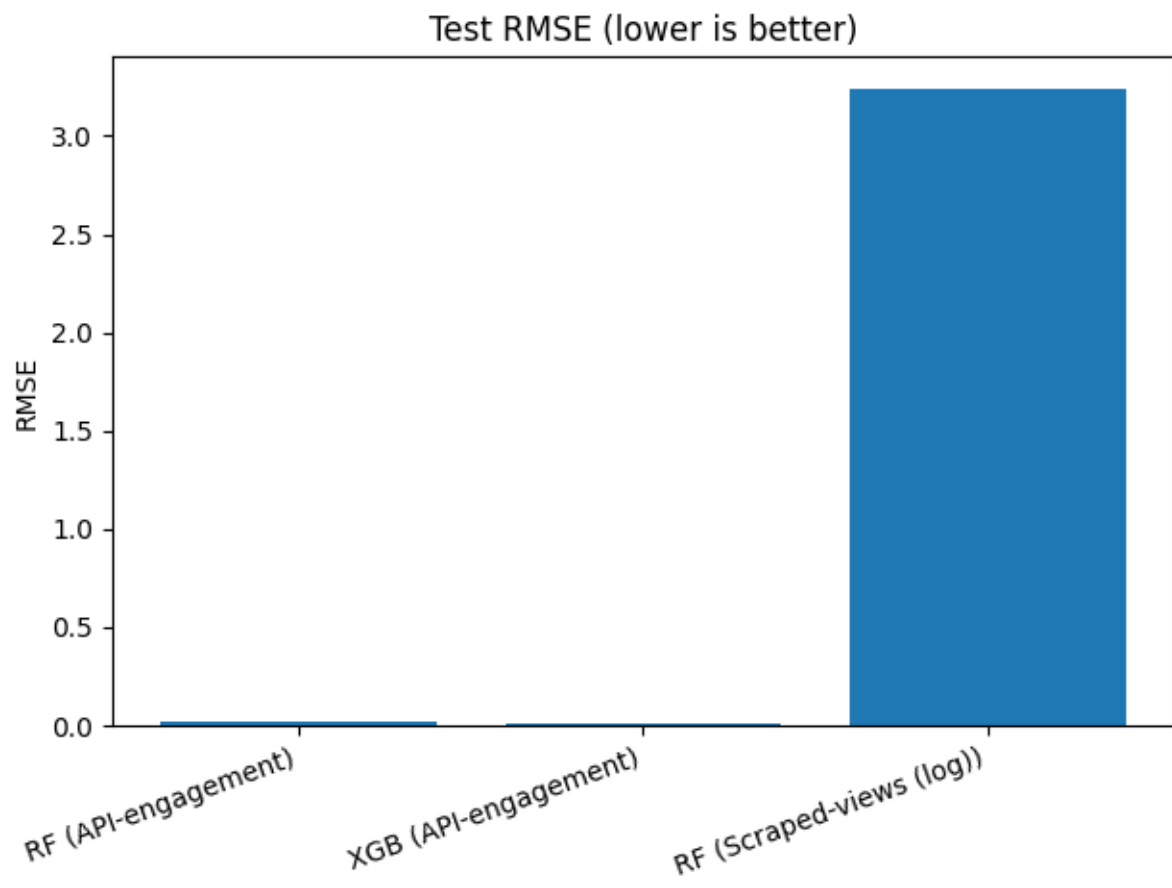




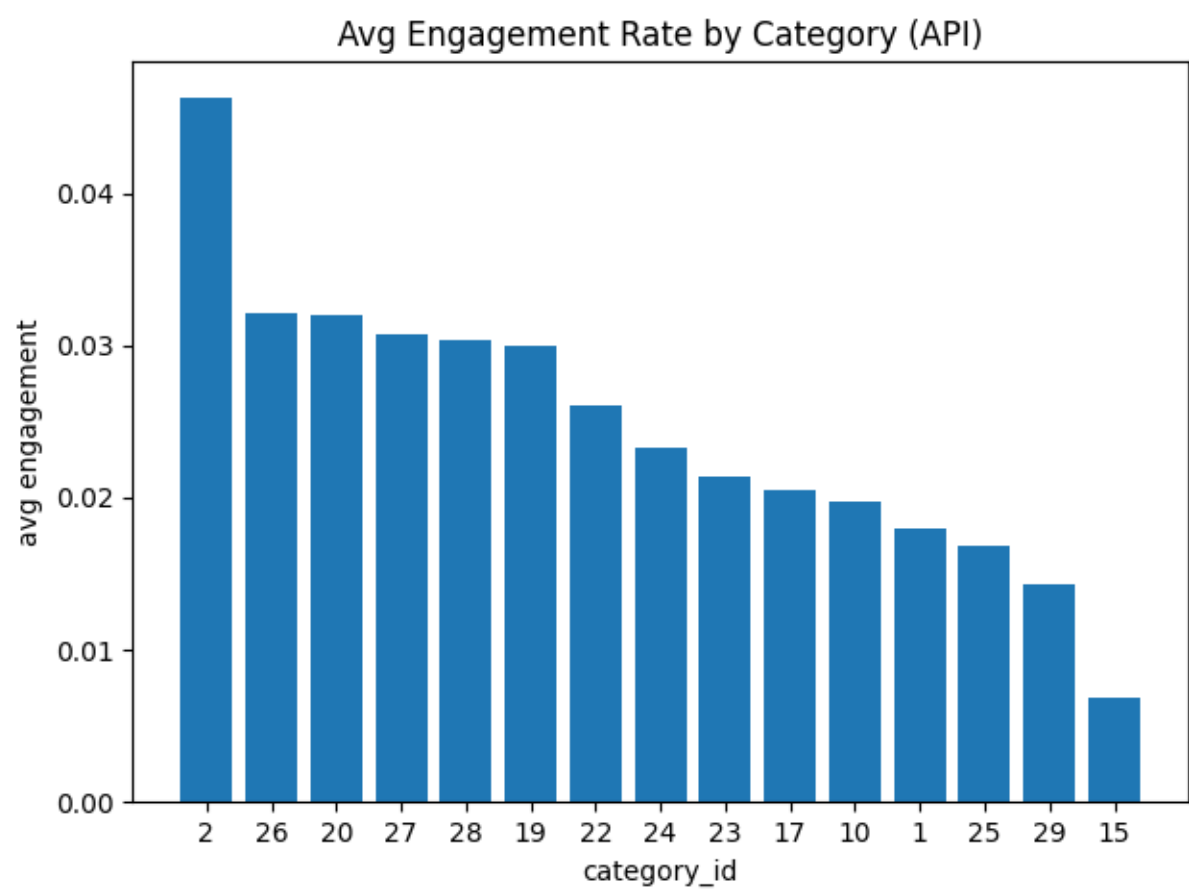
## Visualization:

Test RMSE comparison:

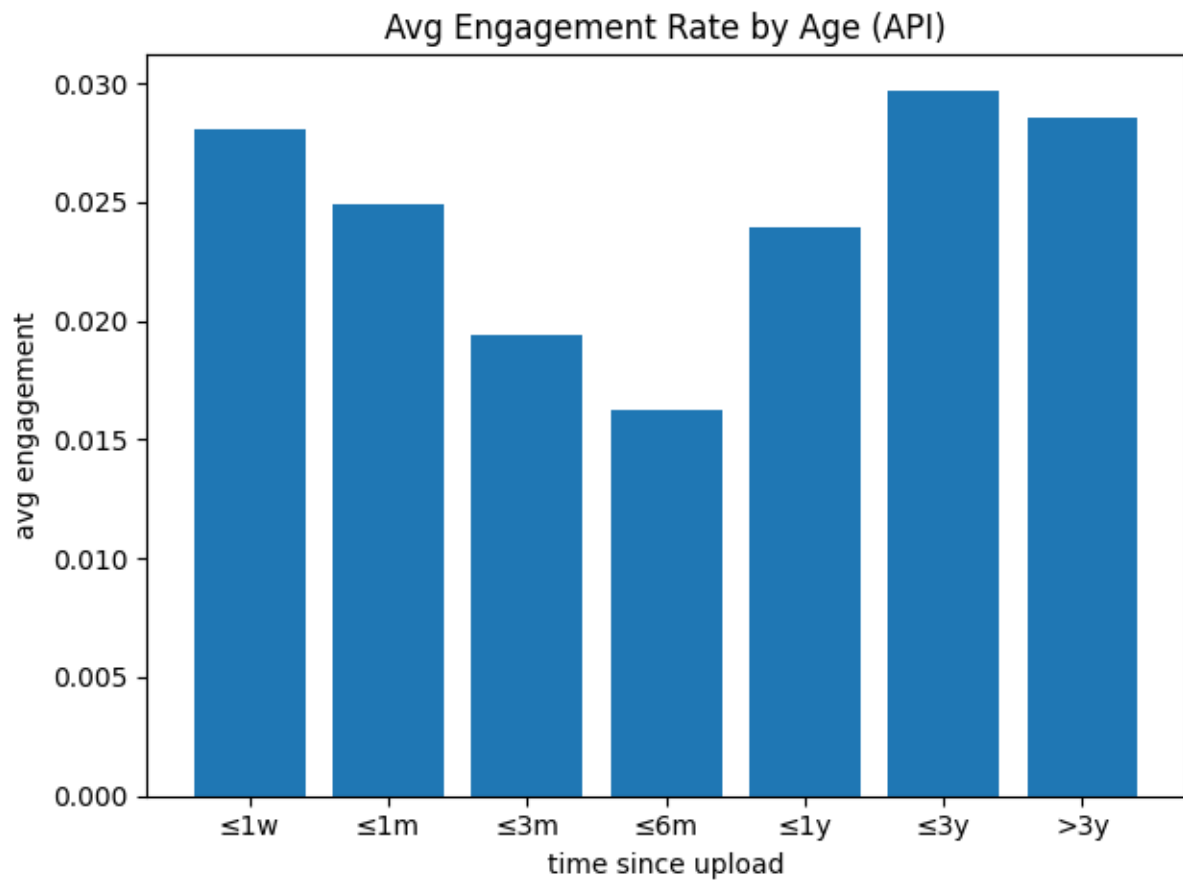
Comparison Test of RMSE between API and Scraped data models:



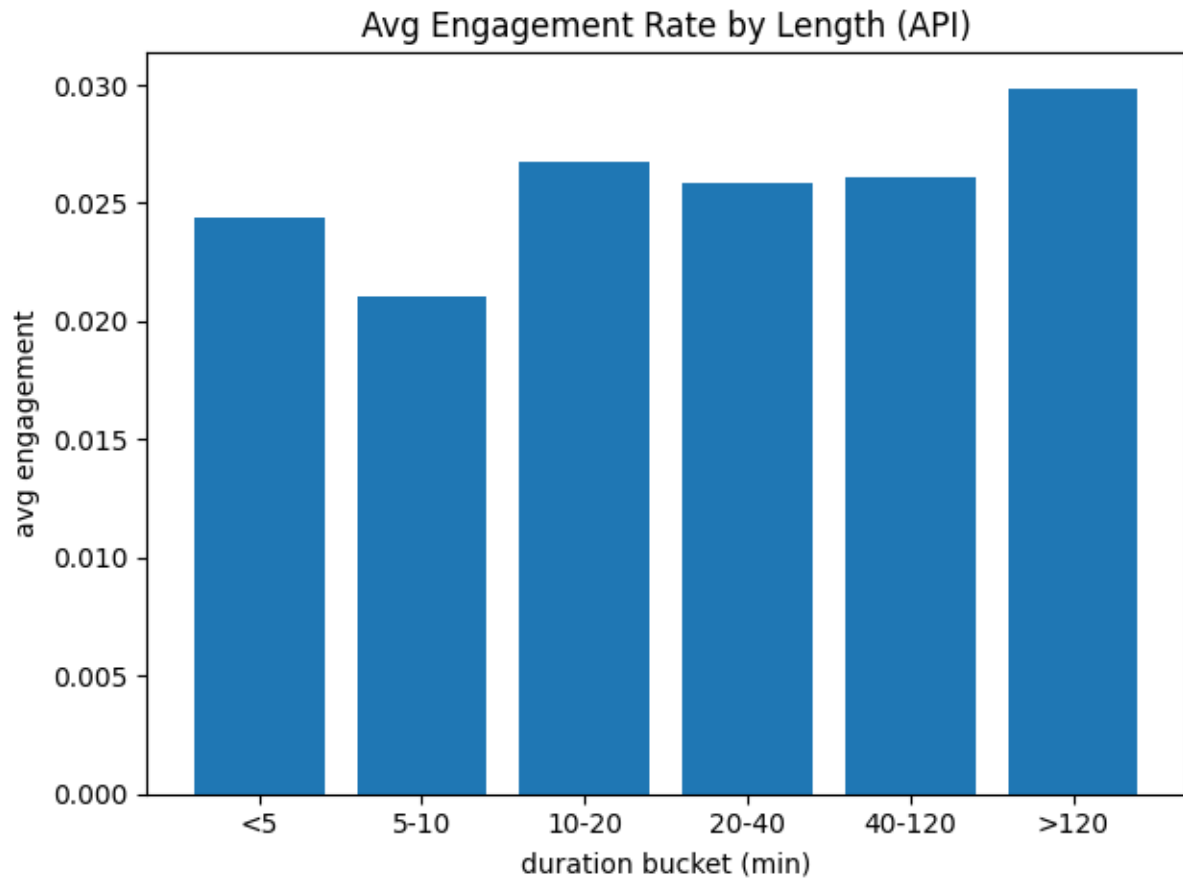
By Category:



By age:



**By Length:**



## Discussions and Conclusions:

Project findings, including insights gained from data analysis:

The results of the project show the differences between the predictive power of the two data sources. The API model outperformed the scraped model in predicting engagement rates by being able to achieve the  $R^2$  of about 0.41 on test data compared to the scraped models near zero correlation when predicting view counts. The outcome tells me that the structured and reliable numerical features of YouTube's API, such as the subscriber count total channel views and how recent the video was uploaded, carries harder when being able to predict values than just textual information. In contrast, the scraped data set, while larger in size, lacked accurate numerical engagement indicators. It mainly relied on text-based metadata such as the titles descriptions and tags which are highly variable across every channel and genre. Because of that, the scraped model struggled to be able to identify consistent patterns in finding viewer behavior.

Several Patterns emerged during the feature of importance and trend analysis. The first videos from music and gaming categorized tended to have the highest engagement, which showed their strong fan communities and frequent sharing behaviors. Secondly, engagement peaks for videos with 5 to 10 minutes in length, which makes it long enough to be valuable to the viewer but short enough to maintain their attention. Third, the recent uploads consistently shower higher engagement, which showed that the freshness and YouTube 's algo's promotion of new content had a high impact.

Challenges encountered during model development:

I encountered a couple of challenges with this project, being technical and data related. The scraped dataset has irregular formatting, broken lines, and missing columns which require more robust preprocessing and error handling. HTML structure changes on YouTube also made the scraping less consistent, often making malformed records. Another issue I faced was the difference in view counts as some videos had millions of views while other videos had only a couple hundred. This made the training model unstable which out log transformation or clipping. The API dataset was a lot cleaner, but it was limited to googles' quota restrictions, which prevented collection of large-scale data in a small amount of time. Finally, since engagement and popularity depend on unobservable factors like thumbnail or algorithm boots, i feel as if even the best model could only partially capture underlying problems.

Ethical and legal considerations of data collection:

All the data collected was limited to publicly available information on YouTube, no private or personal data were accessed or stored. The Project followed YouTube's API terms of Service and used API keys stored in environment variables to protect credentials. Web scraping was used responsibly and only for educational purposes, and I ensured that it caused minimal server load and complied with ethical research standards.

#### Recommendations and Future Work:

I think that several improvements could benefit the performance and scalability of this system. First, the models could incorporate TF-IDF vectorization or transformer-based embeddings to better interpret the context and semantics. Second, i think that expanding the API dataset to include regional view statistics, topic categories, and audience demographics could allow for a better understanding of engagement drivers. Third, i think that applying hyperparameter tuning or ensemble learning could improve predictive accuracy and generalization. Overall, the project successfully demonstrates a machine learning pipeline and provides a solid foundation.