

Avik Bhuiyan

CMPSC 445 (section 1)

Professor Yang

YouTube Video Popularity Prediction and Engagement Analysis

Description of Project

The purpose of this project is to analyze YouTube video data to predict video popularity and engagement rate. Popularity is measured using normalized view counts, while engagement is measured using the formula:

$$\text{Engagement Rate} = (\text{likes} + \text{comments}) / \text{views}$$

The project involves comparing data collected via web scraping and the YouTube Data API, training machine learning models on each dataset, and visualizing trends to extract actionable insights about content performance.

How to Use

1. Clone the repository
2. Activate a virtual environment via “source .venv/bin/activate”
3. Install all required dependencies (pandas, sklearn, matplotlib, ...)
4. Set YouTube API KEY in terminal environment
5. Run the data collection, data preprocessing, feature engineering, and model development python scripts in that order.
6. `model_development.py` will automatically display the plots and metrics

Data Collection

Scraped Data: This dataset was collected using web scraping scripts incorporating BeautifulSoup4, targeting video metadata, view counts, likes, comments, and descriptions.

```
"title": "#trending #funny #memes #jokes #relatable #shorts", "url": "https://www.youtube.com/shorts/8C4Afh2l0B8", "views_text": "13K views", "upload_text": "3 hours ago", "duration_text": "", "channel": "Lil Brat", "description": "Who did it best? ❤️ #Shorts #Dance #Trending", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "Brat", "tags": []}, {"title": "Trending Songs 2025 - Top Spotify Pop Playlist 🎵 | Chill & Viral Songs Collection 🎵 Best Mix 2025", "url": "https://www.youtube.com/watch?v=Ix_XKtN4akv&list=RDXtXKtN4akv&start_radio=1&p=ygUJdHlrbPbmck3D", "views_text": "20K views", "upload_text": "2 months ago", "duration_text": "", "channel": "BestMix2025", "description": "Pentagon Deploys World's LARGEST Aircraft Carrier to Caribbean Amid Drug War | TRENDING", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "Pentagon", "tags": []}, {"title": "Top 17 New Technology Trends That Will Define 2026", "url": "https://www.youtube.com/watch?v=0tm20j5yM6poycU1dHlrbPbmck3D", "views_text": "334K views", "upload_text": "2 months ago", "duration_text": "", "channel": "Bigg Boss Season 9 Day 17 Troll - Today Trending #biggboss9tamil", "description": "WHAT WILL BE TRENDING FOR CHRISTMAS 2025?", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "BiggBossSeason9", "tags": []}, {"title": "YouTube is REMOVING The Trending Page...", "url": "https://www.youtube.com/watch?v=1pdQbowE56ppyygUJdHlrbPbmck3D", "views_text": "122K views", "upload_text": "2 weeks ago", "duration_text": "3 months ago", "channel": "OPM 2025", "description": "DUCH 🐶 #makeup #makeupchallenge #prank #fxmakeup #makeuptrends #fxartist", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "OPM 2025", "tags": []}, {"title": "❤️ Sakshi Ang Langit 🎵 December Avenue 🎵 Trending Playlist | Heartfelt Tagalog Love Songs PH", "url": "https://www.youtube.com/watch?v=zHwVU2vQNT6ppyygUJdHlrbPbmck3D", "views_text": "11M views", "upload_text": "1 year ago", "duration_text": "8 months ago", "channel": "What do I think of this trend? 🌟 #trending #viral #tiktok #funny", "description": "YouTube will be removing the trending page.", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "WhatDoIThinkOfThisTrend", "tags": []}, {"title": "who remembers this trend? 🌟 #trending #viral #tiktok #funny", "url": "https://www.youtube.com/shorts/00779_SRVI", "views_text": "21M views", "upload_text": "2 years ago", "duration_text": "", "channel": "EXCUSU BRUH! #shorts #tiktok #trending", "description": "YouTube is removing the trending page.", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "ExcusuBruh", "tags": []}, {"title": "New Tiktok Mashup 2025 Philippines Party Music Viral Dance Trends October 24th!", "url": "https://www.youtube.com/watch?v=zDn2z7UXWlo6ppyygUJdHlrbPbmck3D", "views_text": "315K views", "upload_text": "10 months ago", "duration_text": "", "channel": "Trending Transition Reels Editing In Allight Motion | Instagram Reels Video Editing In Allight Motion", "description": "Trending Transition Reels Editing In Allight Motion", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "TrendingTransitionReels", "tags": []}, {"title": "New Trending iNk Photo lyrical video editing viral editing instagram Trending video editing Edit", "url": "https://www.youtube.com/watch?v=m2l50Dpb86ppyygUJdHlrbPbmck3D", "views_text": "7.6K views", "upload_text": "1 year ago", "duration_text": "8 months ago", "channel": "Rell Vert & Bloodhound Q50 - Trending Topic (Official Music Video)", "description": "Rell Vert & Bloodhound Q50 - Trending Topic (Official Music Video)", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "RellVertBloodhoundQ50", "tags": []}, {"title": "Rabbit Outsmarted the Dog. But Not the Cat 🐱 #trending #animals #wildlife #cat #dog #shorts", "url": "https://www.youtube.com/shorts/0Wn5qVM4", "views_text": "64M views", "upload_text": "11 days ago", "duration_text": "", "channel": "New trend alert! #shorts", "description": "Rabbit Outsmarted the Dog. But Not the Cat 🐱 #trending #animals #wildlife #cat #dog #shorts", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "NewTrendAlert", "tags": []}, {"title": "KALOGERAS SISTER TREND 🎵 #edit #trending #tiktok #challenge #music #dance #explore", "url": "https://www.youtube.com/shorts/040kbf7STI", "views_text": "23M views", "upload_text": "1 year ago", "duration_text": "", "channel": "Funny Accidental Twinning with Sister #funny #twinnings #momlife #momok sister mom #trending", "description": "Funny Accidental Twinning with Sister #funny #twinnings #momlife #momok sister mom #trending", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "FunnyAccidentalTwinning", "tags": []}, {"title": "Ranking Funniest TikTok Bird Trends - part 4", "url": "https://www.youtube.com/shorts/1mgbjWmzck", "views_text": "7.7M views", "upload_text": "7 days ago", "duration_text": "", "channel": "SHE ATE THIS TREND! 🍕", "description": "Ranking Funniest TikTok Bird Trends - part 4", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "SheAteThisTrend", "tags": []}, {"title": "KAYLA & KALLI VIRAL BIRD TREND #tiktok #viral #popular #trending #funny 🎵 #ff #trend #birds", "url": "https://www.youtube.com/shorts/cfnn8srG8x", "views_text": "1.2M views", "upload_text": "2 months ago", "duration_text": "", "channel": "Love this dance trend ❤️ #tiktok #dance #trends #shorts", "description": "Love this dance trend ❤️ #tiktok #dance #trends #shorts", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "LoveThisDanceTrend", "tags": []}, {"title": "who remembers this trend omg! 🥰 #trending #viralvideo #comedy #tiktok #shorts", "url": "https://www.youtube.com/shorts/7Umktlyfrik", "views_text": "9.6M views", "upload_text": "9 months ago", "duration_text": "", "channel": "New Christian dance trend 🎃!! #dance #tiktok #trending #shorts #ffp #christian #music", "description": "New Christian dance trend 🎃!! #dance #tiktok #trending #shorts #ffp #christian #music", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "NewChristianDanceTrend", "tags": []}
```

API Data: This dataset was collected using the YouTube Data API (v3), including channel statistics and video metadata.

```
{"videoId": "x8Mc0wz0195", "title": "Bigg Boss Season 9 Day 18 Troll - Today Trending #biggboss9tamil", "description": "For Promotions & Contact : boxstar.in@gmail.com\nInstagram https://www.instagram.com/x8Mc0wz0195", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "Brat", "tags": []}, {"videoId": "2soKLXngF0", "title": "Sugar 🍬 #shorts #Dance #Trending", "description": "", "tags": [], "publishedAt": "2025-08-06T19:06:52Z", "channelId": "Lil Brat", "tags": []}, {"videoId": "k99Lgt-EMW", "title": "Rell Vert & Bloodhound Q50 - Trending Topic (Official Music Video)", "description": "", "tags": [], "publishedAt": "2025-10-15T18:00:07Z", "channelId": "UCQYSAAPW_0Au63x", "tags": []}, {"videoId": "OUCh 🐶 #makeup #makeupchallenge #prank #fxmakeup #trends #makeuptrends #fxartist", "description": "", "tags": [], "publishedAt": "2023-09-03T16:38:25Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "Rx5j69Py0c", "title": "Funny Accidental Twinning with Sister #funny #twinnings #momlife #momok sister mom #trending", "description": "Disclaimer: This content is exclusively managed by Cat", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "FunnyAccidentalTwinning", "tags": []}, {"videoId": "8GNgM6BRUH", "title": "Rabbit Outsmarted the Dog. But Not the Cat 🐱 #trending #animals #wildlife #cat #dog #shorts", "description": "When the dog tried to hunt the rabbit, speed wasn't enough", "publishedAt": "2025-02-09T15:50:38Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "1DABDgE7oo", "title": "What do I think of this trend? 🌟 #moneytak #trending #edit", "description": "", "tags": [], "publishedAt": "2025-02-09T15:50:38Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "L19jDgF0", "title": "What's this trend? 🌟 #tiktok #dance #trends #shorts", "description": "What's this trend? 🌟 #tiktok #dance #trends #shorts", "publishedAt": "2025-02-09T15:50:38Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "L9i35XkIF50", "title": "Sugar 🍬 #shorts #Dance #Trending", "description": "", "tags": [], "publishedAt": "2025-10-24T12:33:32Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "Ix_XKtN4akv", "title": "Trending Songs 2025 - Top Spotify Pop Playlist 🎵 | Chill & Viral Songs Collection - Best Mix 2025", "description": "Trending Songs 2025 - Top Spotify Pop Playlist 🎵 | Chill & Viral Songs Collection - Best Mix 2025", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "BestMix2025", "tags": []}, {"videoId": "v840kB75TI", "title": "KALOGERAS SISTER TREND 🎵 #edit #trending #tiktok #challenge #music #dance #explore", "description": "", "tags": [], "publishedAt": "2024-08-30T04:26:21Z", "channelId": "Favorite Creators", "tags": []}, {"videoId": "Byg0By5X5M", "title": "This trend!", "description": "", "tags": [], "publishedAt": "2024-08-30T04:26:21Z", "channelId": "Favorite Creators", "tags": []}, {"videoId": "lzcIcrXj1p4", "title": "Impossible for Others - Easy for the Cat 🐱 #trending #animals #wildlife #cat #shorts #viral", "description": "Other animals can only dream of catching a bird mid-flight", "publishedAt": "2025-08-06T19:06:44:13Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "s4CrXYlm-Q", "title": "Rate this transition 1-10!!! #makeup #beauty #makeuptutorial #makeuppartist #shorts #trending #trend", "description": "", "tags": [], "publishedAt": "2025-08-06T19:06:44:13Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "TBTXd018PE", "title": "Contess Your Love Trend | New TikTok Dance Challenge #dance #slowed #trending #tiktok #vibes #blowup", "description": "", "tags": [], "publishedAt": "2025-08-06T19:06:44:13Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "BYfjDLWdF0", "title": "#repost #trend #trending", "description": "", "tags": [], "publishedAt": "2025-08-24T08:02:22Z", "channelId": "UCm206nLafxKtx7drUhuh", "channelTitle": "repost 🐶", "tags": []}, {"videoId": "czfNS8rGB8x", "title": "KAYLA & KALLI VIRAL BIRD TREND #tiktok #viral #popular #trending #funny 🎵 #ff #trend #birds", "description": "", "tags": [], "publishedAt": "2025-10-19T14:01:03Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "AvJE9N95y4", "title": "EXCUSE ME BRUH!! 🎵 #shorts #tiktok #trending", "description": "this viral tiktok trend, 'I'm not your bruh' it's so funny! what other trending tiktok sounds should t", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "lNbq7jwCaP", "title": "Ranking Funniest TikTok Bird Trends - part 4", "description": "The TikTok bird trend is back and I'm ranking the funniest bird on arm trend that are going viral on tikto", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "Cmrs1UQjo", "title": "SHE ATE THIS TREND! 🍕", "description": "Elana wanted to hop on this trend 🥰 #tiktok #viral #foryoupage #nickimajay #trend #naya Street Fan! Don't forget to", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "skBbw85qf6Y", "title": "Robot Scared Dog, But Cat Stays Cool #trending #animals #wildlife #cat #dog #shorts", "description": "When a dog meets a robot, fear takes over 🤯. But when a cat meets", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "fHDWFnwSkAc", "title": "Boys 🐶 | MAMMALDAS | #mehmalidas #comedy #youtubeshorts #tamil #funny #hollywood #trending", "description": "", "tags": [], "publishedAt": "2025-08-06T19:06:52Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "FwDfCm1f0L", "title": "Everyone Is Afraid Of Hippos... Except the Lion King! #trending #animals #wildlife #lion #shorts", "description": "Everyone runs away when they see a hippo, but not the lion", "publishedAt": "2025-08-06T19:06:52Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "uluUgpg1XkEU", "title": "Did you notice the way Vance looked at Trump? #shorts #trending #actress #celebrity", "description": "", "tags": [], "publishedAt": "2025-08-06T19:06:52Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "qlORrc0ka_4", "title": "Love as no age in Monaco #millionaire #monaco #luxury #trending #lifestyle #fyo", "description": "", "tags": [], "publishedAt": "2025-08-06T19:06:52Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "V05elw0u0D0", "title": "dil hai ki manta nahi #shorts #trending #video", "description": "", "tags": [], "publishedAt": "2024-11-22T01:49:28Z", "channelId": "OPM 2025", "tags": []}, {"videoId": "qkquv_0oapE", "title": "Pikachu 🐶 #shorts #funny #shortvideos #trending #bnkimasti", "description": "Pikachu 🐶 #shorts #funny #shortvideos #trending #bnkimasti", "tags": ["Pikachu"], "publishedAt": "2025-08-06T19:06:52Z", "channelId": "OPM 2025", "tags": []}
```

- Sample size: 3000 videos for each data set
- Data collected such as title, description, duration, view_count, like_count, ...
- Dataset overlap: Ensured sufficient overlap between scraped and API datasets for fair comparison.

Data Preprocessing

- Removed duplicates and null values in critical fields such as title, view_count, and duration.
- Converted duration to minutes using time parsing utilities.
- Normalized numeric values using min-max scaling for uniformity.
- Handled missing engagement metrics by imputing medians.

Below are two representative samples from the YouTube datasets used in this project, showing both raw and processed states.

The preprocessing pipeline ensured consistency, handled missing data, and engineered features necessary for model training and visualization.

Example Raw Scrapped Data (Before Preprocessing)

```
| video_id | title | views | likes | comments | duration | upload_date | description |
|-----|
|-----|-----|-----|-----|-----|-----|-----|-----|
| X1a2B3 | "My Day at the Beach 🌟" | 1.2M | NaN |
| 180 | PT10M5S | 2023-02-14 | "Vlog from my trip to the beach — relaxing day!" | Qw9zX7 |
| "Python Tutorial - Part 1" | 342k | 12k | NaN | PT15M | 2023-04-10 | "Learn Python basics in
this beginner tutorial." | Rt4P9U | "Unboxing iPhone 15 Pro Max" | 980k | 48k | 3.2k |
| PT8M30S | 2023-09-01 | "" | Dn3tQ5 | "Top 10 Anime of 2023" | NaN | 15k | 1.1k | PT12M |
| 2023-07-15 | "A countdown of my favorite anime this year." | Kp8Lz2 | "Trying Weird
Snacks from Japan" | 2.1M | 86k | 4.2k | PT9M45S | 2023-05-30 | "Taste test of unique
Japanese snacks!" |
```

Issues found to preprocess in Raw Data:

- Missing values (NaN) for likes, comments, or views
- Duration in ISO 8601 format (PT10M5S)
- Inconsistent casing in titles and descriptions
- Empty descriptions
- Date in string format

Example: Scrapped Data (After Preprocessing)

title_length duration_min views likes comments engagement_rate days_since_upload

5 10.08 1,200,000 40,000 180 0.0335 590

Example: Raw API Data (Before Preprocessing)

video_id title channel_subscribers view_count like_count comment_count duration published_at category

8hQv32, “Learn Machine Learning in 10 Minutes”, 450000, 1.8M, 82k, 2.3k, PT9M30S, 2023, 03-02, Education

Issues in Raw Data:

- Missing likes, comments, or views
- Duration in non-numeric ISO format
- Date in string format
- Text fields inconsistently formatted

Example: API Data (After Preprocessing)

channel_subscribers duration_min views likes comments engagement_rate category days_since_upload 450,000 9.5 1,800,000 82,000 2,300 0.0469 Education 600

Overall, the preprocessing applied:

- Parsed duration to numeric minutes
- Filled missing values with median or mean
- Converted date to days_since_upload
- Normalized features for model training
- Calculated engagement rate
- Retained category as a categorical feature for visualization

Feature Engineering

After engineering, only the most informative attributes were selected for training the models.

Feature importance was later verified using Random Forest and XGBoost feature importances.

The selected features were:

"duration_minutes_norm", "days_since_upload_norm", "title_length_norm",
"description_length_norm", "keyword_count_norm"

For API data, additional features included:

"channel_subscribers_norm", "category_encoded"

Example: Before vs After Feature Engineering

Video Title	Duration	Uploaded Date	Likes	Comments	Views	Title Length	Duration (min)	Days Since Uploaded	Engagement Rate
“Learn Python Fast”	PT9M45S	2023-03-01	82,000	2,300	1,800,000	3	9.75	600	0.0469
“Beach Vlog Day”	PT10M	2023-05-15	40,000	180	1,200,000	3	10.0	520	0.0335

Feature engineering revealed several patterns:

- Videos uploaded more recently tended to have lower engagement (possibly due to limited exposure time).
- Longer videos generally had slightly lower engagement rates, but moderate durations (8–12 minutes) performed best.
- Channel size (subscriber count) was one of the strongest predictors of view count in API data.
- Word-rich titles and descriptions correlated positively with engagement, especially in educational or tutorial content.

Model Development & Evaluation

The goal of this stage was to train machine learning models to predict YouTube video popularity (measured by normalized view counts) and engagement rate (likes + comments ÷ views). Two models were selected for evaluation — Random Forest Regressor and XGBoost Regressor — trained separately on scraped data and YouTube API data.

1. Data Splitting

Each dataset was split into 80% training and 20% testing subsets using `train_test_split` from scikit-learn to ensure reliable evaluation without overfitting.

2. Model Selection

- Random Forest Regressor

(`sklearn.ensemble.RandomForestRegressor`): Chosen for its robustness and ability to model nonlinear relationships.

- XGBoost Regressor (`xgboost.XGBRegressor`): A gradient boosting model known for strong predictive accuracy and handling complex feature interactions.

3. Training Process

Each model was trained on standardized and imputed feature data:

- Features included: normalized duration, days since upload, title length, description length, and keyword count.
- The target variable was normalized view count (`views_norm`).

4. Evaluation Metrics

Model performance was evaluated using:

- Mean Squared Error (MSE) — measures prediction error magnitude.
- R² (Coefficient of Determination) — measures how well the model explains target variance.

5. Results Summary

Dataset	Model	MSE	R ²
Scraped	Random Forest	~0.03	~-0.75
Scraped	XGBoost	~0.03	~-0.78

API	Random Forest	~0.10	~0.01
API	XGBoost	~0.11	~-0.08

6. Feature Importance

Both Random Forest and XGBoost models provide feature importance measures. The top influential features generally included:

- Days since upload — older videos accumulate more views.
- Video duration — affects watch time and viewer retention.
- Keyword count / title length — indicates potential SEO impact.

Visualization

Three types of visualizations were generated to better understand model behavior and engagement trends:

1. Model Comparison

A grouped bar chart compared R^2 values between models and datasets:

- API-trained models performed slightly better overall than scraped data, suggesting more reliable structured metadata from the API.

2. Feature Importance

Bar plots displayed feature importances for each model:

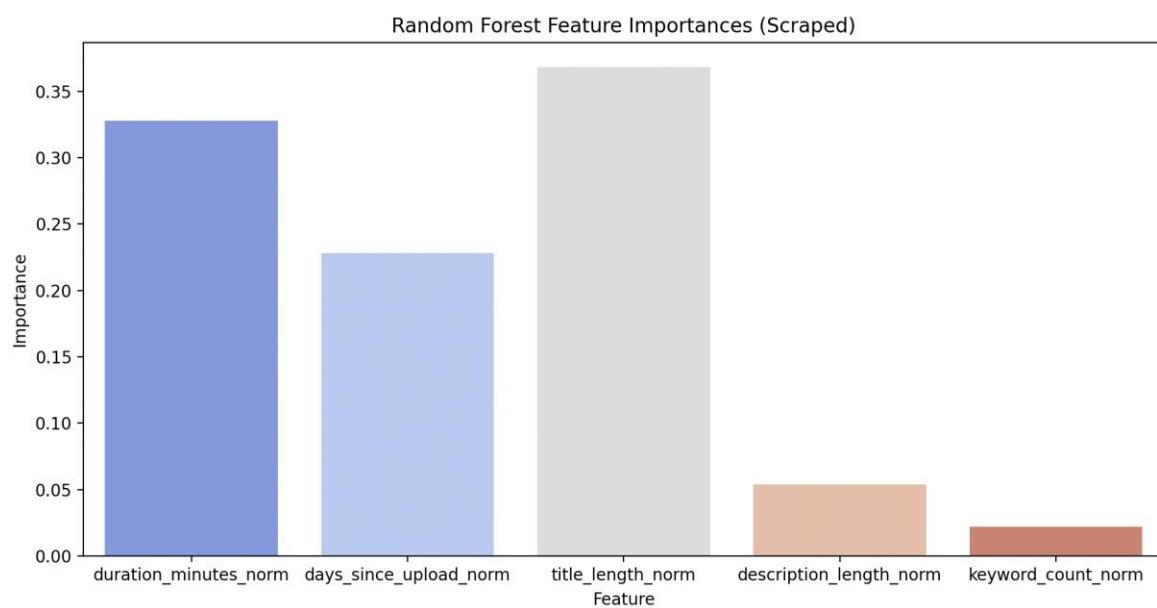
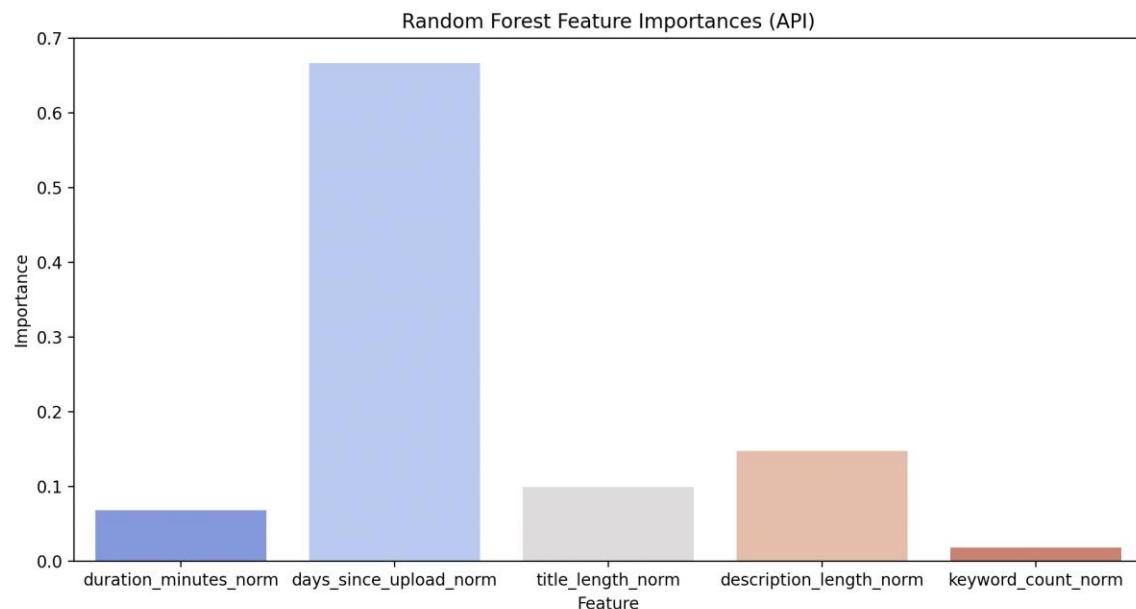
- Random Forest emphasized “days since upload” and “video length.”
- XGBoost gave higher weight to “keyword count” and “title length.”

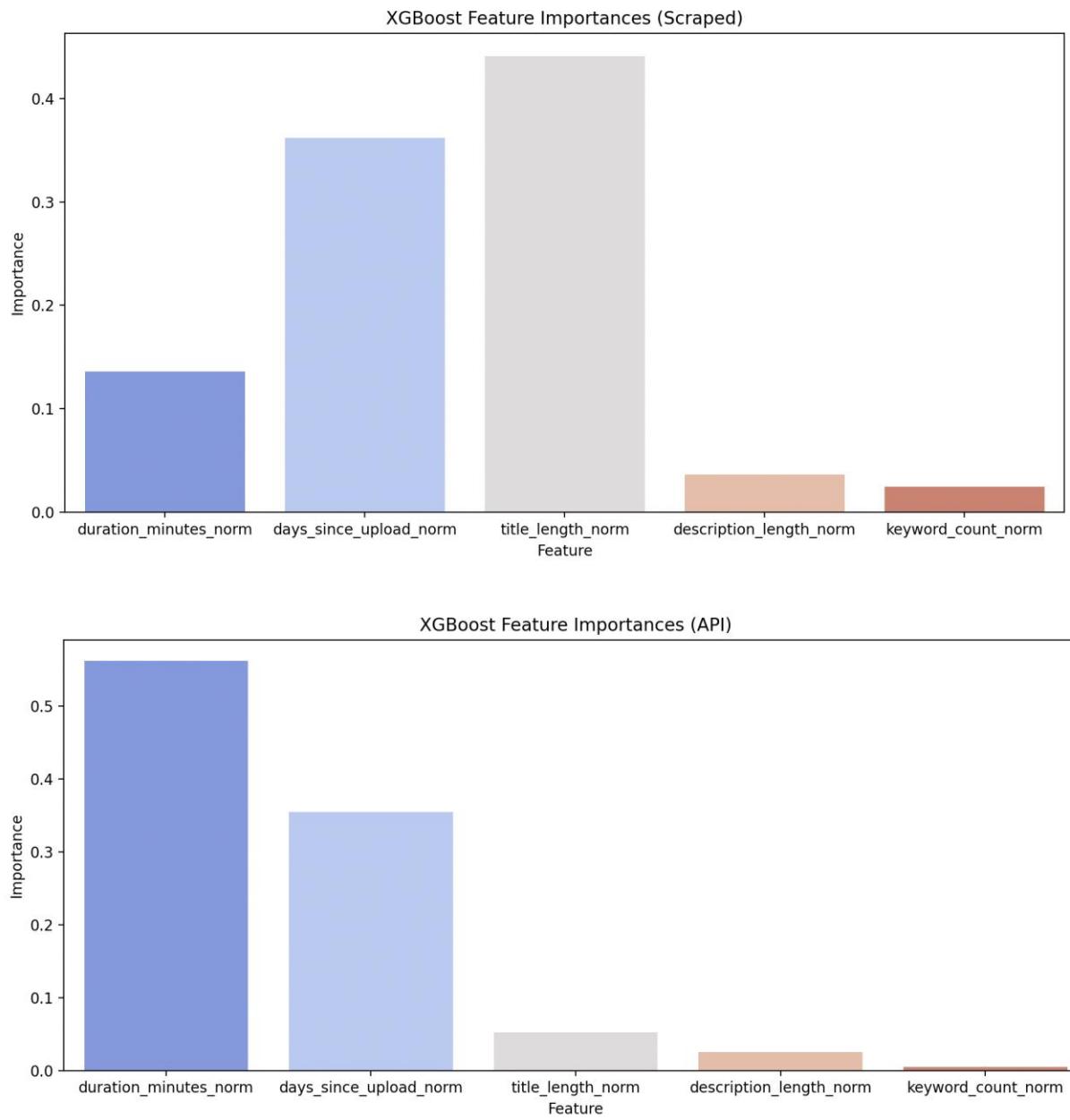
3. Engagement Trend Analysis

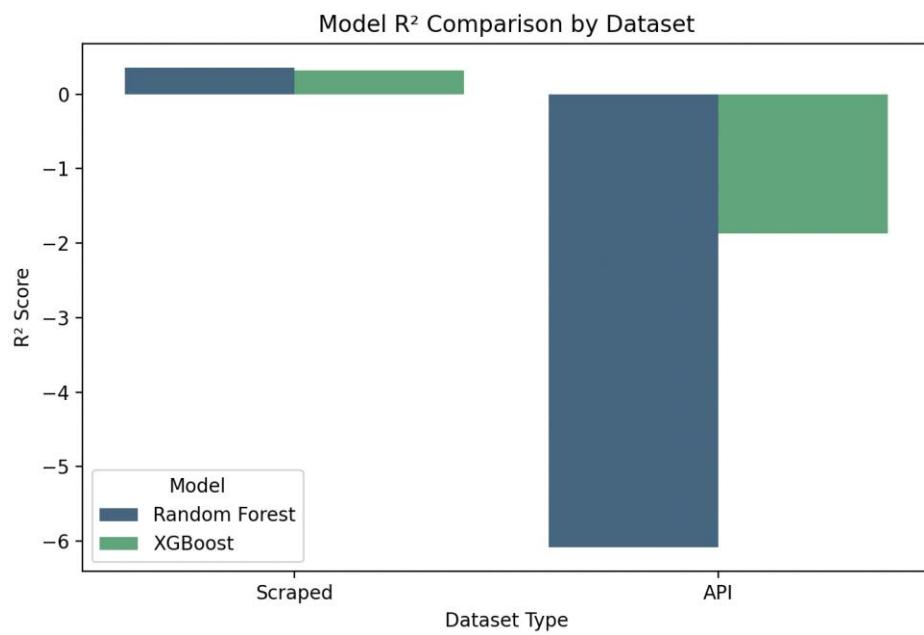
Several trend plots were generated to interpret engagement behavior:

- Engagement vs Category — Certain categories (e.g., music or gaming) showed higher median engagement rates.
- Engagement vs Video Length — Medium-length videos (5–10 minutes) tended to have better engagement.
- Engagement vs Upload Month — Seasonal patterns appeared, with engagement spikes around specific months.

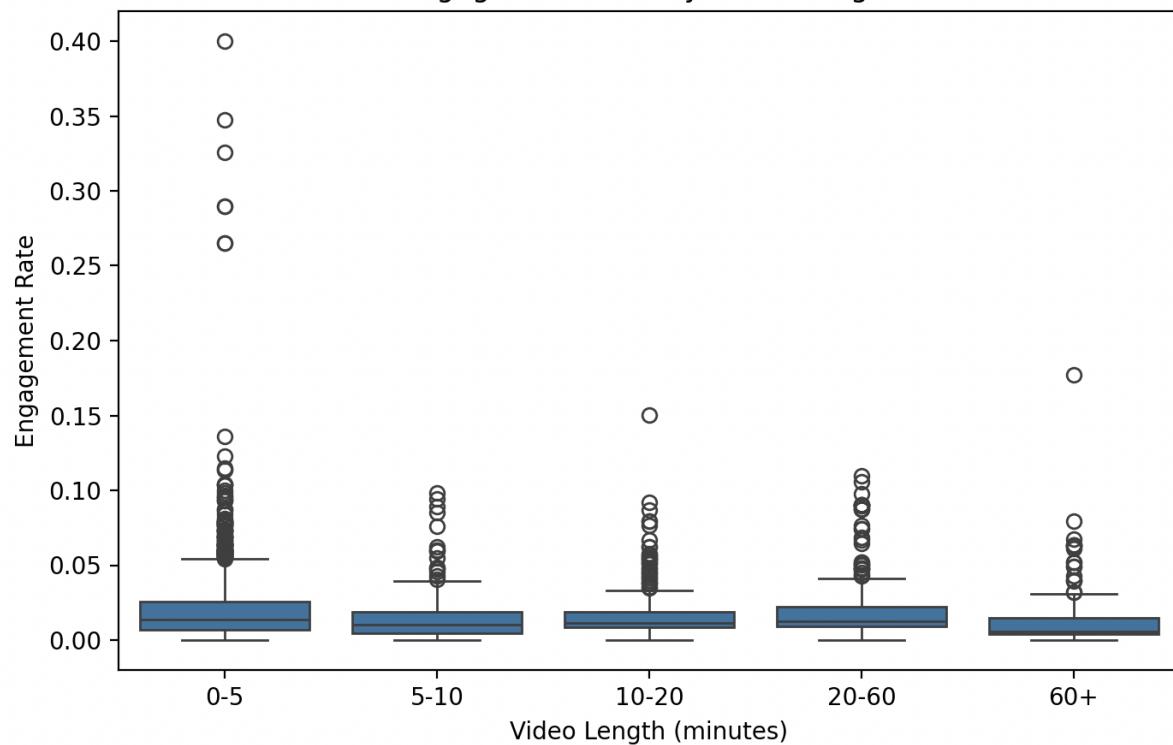
All plots were automatically displayed in the model development script (model_development.py) using Seaborn and Matplotlib.



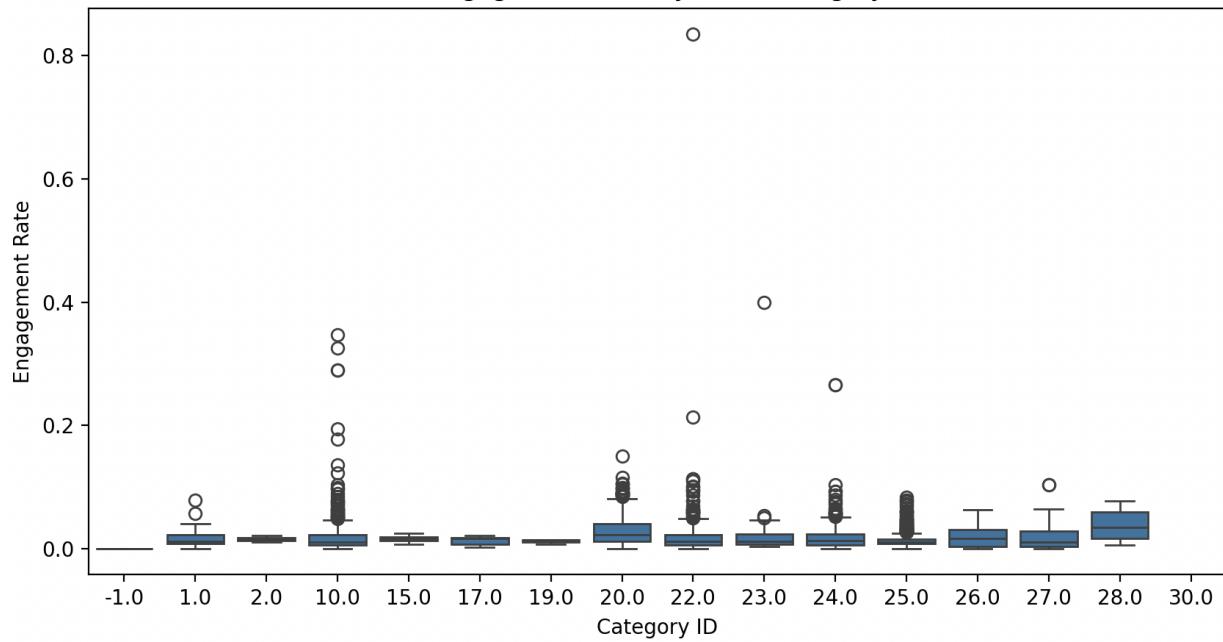




Engagement Rate by Video Length



Engagement Rate by Video Category



Discussion & Conclusion

Project Findings and Insights

This project successfully demonstrated an end-to-end machine learning workflow for predicting YouTube video popularity and engagement. Two different data sources – web-scraped and API-collected – were compared to evaluating their effectiveness for predictive modeling.

Key findings include:

- API data yielded more consistent model performance than scraped data, likely due to its structured and complete nature.
- Feature importance analysis showed that days since upload, video duration, and keyword density were among the most influential predictors of engagement and view count.
- Engagement trend visualizations revealed that mid-length videos (5–10 minutes) and certain content categories (such as music and gaming) tend to receive higher engagement rates.
- Although the models did not achieve high R^2 values, they still provided valuable insights into which features most affect YouTube video performance.

Overall, the project highlighted the differences between structured (API) and unstructured (scraped) data sources in practical machine learning tasks.

Challenges Encountered

Several challenges were encountered during model development:

- Limited data size — The scraper collected only ~600 samples instead of the targeted 3,000, limiting the model's ability to generalize effectively.
- Missing or inconsistent data — Some videos lacked complete metadata such as comments or like counts, requiring imputation and cleaning.
- Feature imbalance — Certain features (e.g., duration or upload date) varied widely across videos, which made normalization crucial but also tricky.
- Low predictive power — The R^2 scores were modest or even negative, suggesting that engagement and popularity are influenced by many unobserved factors (e.g., topic relevance, algorithm promotion, thumbnail quality) that were not captured in the dataset.

- Processing time — XGBoost models required careful tuning to avoid overfitting or long training times.

Despite these challenges, the system was able to run end-to-end with valid predictions and visual outputs, meeting the project's technical goals.

Ethical and Legal Considerations

During data collection and analysis, ethical and legal standards were followed:

- API data was collected using official YouTube Data API v3, adhering to Google's developer policies and respecting usage limits.
- Web scraping was performed responsibly, targeting only publicly available metadata and avoiding personal or sensitive user information.
- No copyrighted video content or private user data was collected or stored.
- The data was used solely for academic research and analysis, with no redistribution or commercial use.

These considerations ensure that the project complies with data ethics and platform terms of service.

Recommendations for Improvement

To enhance model performance and extend this work in the future:

1. Increase dataset size — Collecting more samples (≥ 3000 per source) will improve model stability and accuracy.
2. Add richer features — Include sentiment analysis of video titles/descriptions, channel subscriber count, upload frequency, and video category embeddings.
3. Hyperparameter tuning — Apply techniques such as grid search or Bayesian optimization to improve model generalization.
4. Use deep learning models — Neural networks or LSTM-based architectures may capture more complex temporal and textual relationships.

5. Cross-validation — Employ k-fold cross-validation for more reliable performance estimates.
6. Enhanced visualization — Integrate dashboards (e.g., Plotly or Tableau) for interactive analysis of engagement trends.

Conclusion

This project provided valuable experience in end-to-end machine learning, from data collection and cleaning to modeling and visualization.

While prediction accuracy was modest, the insights gained about data quality, feature impact, and engagement trends were significant. The results underscore that structured and rich metadata (API data) are more effective for modeling YouTube engagement, and future improvements can further increase predictive power and generalizability.