# YouTube Video Popularity Prediction and Engagement Analysis

## Final Project Report

**Course:** Machine Learning Project **Date:** January 2025 **Project Goal:** Develop and compare two machine learning models to predict video engagement using data from web scraping and YouTube Data API

## Table of Contents

## 1. Description of the Project

### 1.1 Project Overview

This project investigates the factors influencing YouTube video popularity by developing and comparing two distinct machine learning approaches. The primary objective is to predict video view counts (as a proxy for engagement) using features extracted from two different data sources:

### 1.2 Objectives

### 1.3 Methodology

The project follows a structured machine learning pipeline:

### 1.4 Scope and Limitations

**Scope:** - Predicting video popularity based on metadata (not content analysis) - Comparing data collection methodologies - Analyzing feature importance across different data richness levels

**Limitations:** - Web scraping provides limited engagement metrics (no likes/comments) - API has daily quota restrictions (10,000 units/day) - Models predict log-transformed views, not exact view counts - Temporal dynamics not captured (single snapshot in time)

## 2. How to Use

# *Training*

## *Environment Setup*

**Prerequisites:** - Python 3.8 or higher - Google Chrome browser (for web scraping) - YouTube Data API v3 key (for API collection)

### Step 1: Clone and Navigate

```
git clone <repository-url>
cd YouTube-Video-Popularity-Prediction-and-Engagement-Analysis
```

### Step 2: Create Virtual Environment

```
# Windows
python -m venv .venv
.venv\Scripts\activate
# Linux/Mac
python -m venv .venv
source .venv/bin/activate
```

### Step 3: Install Dependencies

```
pip install -r requirements.txt
```

**Required packages:** - `selenium`, `undetected-chromedriver`, `beautifulsoup4` (web scraping) - `google-api-python-client`, `python-dotenv` (YouTube API) - `pandas`, `numpy` (data processing) - `scikit-learn`, `xgboost` (machine learning) - `matplotlib`, `seaborn`, `plotly` (visualization)

Required packages: - `selenium`, undetected-chromedriver, beautifulsoup4 (web scraping) - google-api-python-client, python-dotenv (YouTube API) - pandas, numpy (data processing) - scikit-learn, xgboost (machine learning) - matplotlib, seaborn, plotly (visualization)

Required packages: - selenium, `undetected-chromedriver`, beautifulsoup4 (web scraping) - google-api-python-client, python-dotenv (YouTube API) - pandas, numpy (data processing) - scikit-learn, xgboost (machine learning) - matplotlib, seaborn, plotly (visualization)

Required packages: - selenium, undetected-chromedriver, `beautifulsoup4` (web scraping) - google-api-python-client, python-dotenv (YouTube API) - pandas, numpy (data processing) - scikit-learn, xgboost (machine learning) - matplotlib, seaborn, plotly (visualization)

Required packages: - selenium, undetected-chromedriver, beautifulsoup4 (web scraping) - `google-api-python-client`, python-dotenv (YouTube API) - pandas, numpy (data processing) - scikit-learn, xgboost (machine learning) - matplotlib, seaborn, plotly (visualization)

Required packages: - selenium, undetected-chromedriver, beautifulsoup4 (web scraping) - google-api-python-client, `python-dotenv` (YouTube API) - pandas, numpy (data processing) - scikit-learn, xgboost (machine learning) - matplotlib, seaborn, plotly (visualization)

Required packages: - selenium, undetected-chromedriver, beautifulsoup4 (web scraping) - google-api-python-client, python-dotenv (YouTube API) - `pandas`, numpy (data processing) - scikit-learn, xgboost (machine learning) - matplotlib, seaborn, plotly (visualization)

Required packages: - selenium, undetected-chromedriver, beautifulsoup4 (web scraping) - google-api-python-client, python-dotenv (YouTube API) - pandas, `numpy` (data processing) - scikit-learn, xgboost (machine learning) - matplotlib, seaborn, plotly (visualization)

Required packages: - selenium, undetected-chromedriver, beautifulsoup4 (web scraping) - google-api-python-client, python-dotenv (YouTube API) - pandas, numpy (data processing) - `scikit-learn`,

xgboost (machine learning) - matplotlib, seaborn, plotly (visualization)

Required packages: - selenium, undetected-chromedriver, beautifulsoup4 (web scraping) - google-api-python-client, python-dotenv (YouTube API) - pandas, numpy (data processing) - scikit-learn, `xgboost` (machine learning) - matplotlib, seaborn, plotly (visualization)

Required packages: - selenium, undetected-chromedriver, beautifulsoup4 (web scraping) - google-api-python-client, python-dotenv (YouTube API) - pandas, numpy (data processing) - scikit-learn, xgboost (machine learning) - `matplotlib`, seaborn, plotly (visualization)

Required packages: - selenium, undetected-chromedriver, beautifulsoup4 (web scraping) - google-api-python-client, python-dotenv (YouTube API) - pandas, numpy (data processing) - scikit-learn, xgboost (machine learning) - matplotlib, `seaborn`, plotly (visualization)

Required packages: - selenium, undetected-chromedriver, beautifulsoup4 (web scraping) - google-api-python-client, python-dotenv (YouTube API) - pandas, numpy (data processing) - scikit-learn, xgboost (machine learning) - matplotlib, seaborn, `plotly` (visualization)

### Step 4: Set Up API Key

```
# Create .env file
cp .env.example .env
# Edit .env and add your API key
YOUTUBE_API_KEY=your_api_key_here
```

Get API key from Google Cloud Console: 1. Create a new project 2. Enable YouTube Data API v3 3. Create credentials (API key)

## *Running the Pipeline*

### Option 1: Run Complete Pipeline

```
python run_pipeline.py
```

This executes all stages sequentially: 1. API connection test 2. Data collection (API ~5 min, Scraping 30-60 min) 3. Data preprocessing 4. Feature engineering 5. Model training (Random Forest + XGBoost for both datasets) 6. Visualization generation

### Option 2: Run Individual Stages

```
# 0. Test API connection
python test_api.py
# 1. Collect data
python src/api_collector.py        # ~5 minutes, 2,972 videos
python src/scraper.py               # 30-60 minutes, 3,050 videos
# 2. Preprocess data
python src/preprocessor.py
# 3. Engineer features
python src/feature_engineer.py
# 4. Train models
python src/model.py
# 5. Generate visualizations
python src/visualizer.py
```

**Expected Outputs:** - `data/raw/scraped_data.csv` (3,050 videos) - `data/raw/api_data.csv` (2,972 videos) - `data/processed/scraped_processed.csv` (3,034 after cleaning) - `data/processed/api_processed.csv` (2,911 after cleaning) - `data/processed/scraped_features.csv` (28 features) - `data/processed/api_features.csv` (46 features) - `models/model_scraped_random_forest.pkl` (14 MB) - `models/model_scraped_xgboost.pkl` (1.9 MB) - `models/model_api_random_forest.pkl` (15 MB) - `models/model_api_xgboost.pkl` (1.5 MB) - **Best model** - `reports/results_scraped.json` - `reports/results_api.json` - `reports/figures/*.png` (10 visualization files)

Expected Outputs: - `data/raw/scraped_data.csv` (3,050 videos) - data/raw/api_data.csv (2,972 videos) - data/processed/scraped_processed.csv (3,034 after cleaning) - data/processed/api_processed.csv (2,911

after cleaning) - data/processed/scraped_features.csv (28 features) - data/processed/api_features.csv (46 features) - models/model_scraped_random_forest.pkl (14 MB) - models/model_scraped_xgboost.pkl (1.9 MB) - models/model_api_random_forest.pkl (15 MB) - models/model_api_xgboost.pkl (1.5 MB) - Best model - reports/results_scraped.json - reports/results_api.json - reports/figures/*.png (10 visualization files)

Expected Outputs: - data/raw/scraped_data.csv (3,050 videos) - `data/raw/api_data.csv` (2,972 videos) - data/processed/scraped_processed.csv (3,034 after cleaning) - data/processed/api_processed.csv (2,911 after cleaning) - data/processed/scraped_features.csv (28 features) - data/processed/api_features.csv (46 features) - models/model_scraped_random_forest.pkl (14 MB) - models/model_scraped_xgboost.pkl (1.9 MB) - models/model_api_random_forest.pkl (15 MB) - models/model_api_xgboost.pkl (1.5 MB) - Best model - reports/results_scraped.json - reports/results_api.json - reports/figures/*.png (10 visualization files)

Expected Outputs: - data/raw/scraped_data.csv (3,050 videos) - data/raw/api_data.csv (2,972 videos) - `data/processed/scraped_processed.csv` (3,034 after cleaning) - data/processed/api_processed.csv (2,911 after cleaning) - data/processed/scraped_features.csv (28 features) - data/processed/api_features.csv (46 features) - models/model_scraped_random_forest.pkl (14 MB) - models/model_scraped_xgboost.pkl (1.9 MB) - models/model_api_random_forest.pkl (15 MB) - models/model_api_xgboost.pkl (1.5 MB) - Best model - reports/results_scraped.json - reports/results_api.json - reports/figures/*.png (10 visualization files)

Expected Outputs: - data/raw/scraped_data.csv (3,050 videos) - data/raw/api_data.csv (2,972 videos) - data/processed/scraped_processed.csv (3,034 after cleaning) - `data/processed/api_processed.csv` (2,911 after cleaning) - data/processed/scraped_features.csv (28 features) - data/processed/api_features.csv (46 features) - models/model_scraped_random_forest.pkl (14 MB) - models/model_scraped_xgboost.pkl (1.9 MB) - models/model_api_random_forest.pkl (15 MB) - models/model_api_xgboost.pkl (1.5 MB) - Best model - reports/results_scraped.json - reports/results_api.json - reports/figures/*.png (10 visualization files)

Expected Outputs: - data/raw/scraped_data.csv (3,050 videos) - data/raw/api_data.csv (2,972 videos) - data/processed/scraped_processed.csv (3,034 after cleaning) - data/processed/api_processed.csv (2,911 after cleaning) - `data/processed/scraped_features.csv` (28 features) - data/processed/api_features.csv (46 features) - models/model_scraped_random_forest.pkl (14 MB) - models/model_scraped_xgboost.pkl (1.9 MB) - models/model_api_random_forest.pkl (15 MB) - models/model_api_xgboost.pkl (1.5 MB) - Best model - reports/results_scraped.json - reports/results_api.json - reports/figures/*.png (10 visualization files)

Expected Outputs: - data/raw/scraped_data.csv (3,050 videos) - data/raw/api_data.csv (2,972 videos) - data/processed/scraped_processed.csv (3,034 after cleaning) - data/processed/api_processed.csv (2,911 after cleaning) - data/processed/scraped_features.csv (28 features) - `data/processed/api_features.csv` (46 features) - models/model_scraped_random_forest.pkl (14 MB) - models/model_scraped_xgboost.pkl (1.9 MB) - models/model_api_random_forest.pkl (15 MB) - models/model_api_xgboost.pkl (1.5 MB) - Best model - reports/results_scraped.json - reports/results_api.json - reports/figures/*.png (10 visualization files)

Expected Outputs: - data/raw/scraped_data.csv (3,050 videos) - data/raw/api_data.csv (2,972 videos) - data/processed/scraped_processed.csv (3,034 after cleaning) - data/processed/api_processed.csv (2,911 after cleaning) - data/processed/scraped_features.csv (28 features) - data/processed/api_features.csv (46 features) - `models/model_scraped_random_forest.pkl` (14 MB) - models/model_scraped_xgboost.pkl (1.9 MB) - models/model_api_random_forest.pkl (15 MB) - models/model_api_xgboost.pkl (1.5 MB) - Best model - reports/results_scraped.json - reports/results_api.json - reports/figures/*.png (10 visualization files)

Expected Outputs: - data/raw/scraped_data.csv (3,050 videos) - data/raw/api_data.csv (2,972 videos) - data/processed/scraped_processed.csv (3,034 after cleaning) - data/processed/api_processed.csv (2,911 after cleaning) - data/processed/scraped_features.csv (28 features) - data/processed/api_features.csv (46 features) - models/model_scraped_random_forest.pkl (14 MB) - `models/model_scraped_xgboost.pkl` (1.9 MB) - models/model_api_random_forest.pkl (15 MB) - models/model_api_xgboost.pkl (1.5 MB) - Best model - reports/results_scraped.json - reports/results_api.json - reports/figures/*.png (10 visualization files)

Expected Outputs: - data/raw/scraped_data.csv (3,050 videos) - data/raw/api_data.csv (2,972 videos) - data/processed/scraped_processed.csv (3,034 after cleaning) - data/processed/api_processed.csv (2,911 after cleaning) - data/processed/scraped_features.csv (28 features) - data/processed/api_features.csv (46 features) - models/model_scraped_random_forest.pkl (14 MB) - models/model_scraped_xgboost.pkl (1.9 MB) - `models/model_api_random_forest.pkl` (15 MB) - models/model_api_xgboost.pkl (1.5 MB) - Best model - reports/results_scraped.json - reports/results_api.json - reports/figures/*.png (10 visualization files)

Expected Outputs: - data/raw/scraped_data.csv (3,050 videos) - data/raw/api_data.csv (2,972 videos) - data/processed/scraped_processed.csv (3,034 after cleaning) - data/processed/api_processed.csv (2,911 after cleaning) - data/processed/scraped_features.csv (28 features) - data/processed/api_features.csv (46 features) - models/model_scraped_random_forest.pkl (14 MB) - models/model_scraped_xgboost.pkl (1.9 MB) - models/model_api_random_forest.pkl (15 MB) - `models/model_api_xgboost.pkl` (1.5 MB) - Best model - reports/results_scraped.json - reports/results_api.json - reports/figures/*.png (10 visualization files)

## *Training Time Estimates*

| Stage | Scraped Data | API Data |
|---|---|---|
| Data Collection | 30-60 min | ~5 min |
| Preprocessing | ~30 sec | ~30 sec |
| Feature Engineering | ~20 sec | ~20 sec |
| Model Training | ~5 min | ~3 min |
| Visualization | ~1 min | ~1 min |
| Total | 40-70 min | 10-15 min |

# *Inferencing*

## *Using Trained Models for Prediction*

### Step 1: Load a Trained Model

```
import joblib
import pandas as pd
import numpy as np
# Load the best performing model (XGBoost API)
model = joblib.load('models/model_api_xgboost.pkl')
```

### Step 2: Prepare Input Features

For **API-based model** (requires 28 features):

```python
# Example: Predict views for a new video
new_video = pd.DataFrame({
    'duration_seconds': [600],
    'log_duration': [np.log1p(600)],
    'days_since_upload': [7],
    'log_days_since_upload': [np.log1p(7)],
    'title_length': [50],
    'title_uppercase_ratio': [0.1],
    'title_word_count': [8],
    'has_clickbait_words': [1],
    'has_numbers': [1],
    'has_question': [0],
    'has_exclamation': [1],
    'is_recent': [1],
    'is_short_video': [0],
    'is_long_video': [0],
    'log_likes': [np.log1p(5000)],
    'log_comments': [np.log1p(200)],
    'engagement_rate': [0.05],
    'engagement_score': [0.045],
    'tag_count': [10],
    'description_length': [500],
    'log_channel_subscribers': [np.log1p(100000)],
    'channel_video_count': [500],
    'likes_per_view': [0.045],
    'comments_per_view': [0.005],
    'video_age_category_encoded': [1],
    'duration_category_encoded': [2],
    'title_length_category_encoded': [1],
    'channel_size_encoded': [3],
    'category_name_encoded': [10]
})
```

For **scraped model** (requires 17 features):

```python
new_video_scraped = pd.DataFrame({
    'duration_seconds': [600],
    'log_duration': [np.log1p(600)],
    'days_since_upload': [7],
    'log_days_since_upload': [np.log1p(7)],
    'title_length': [50],
    'title_uppercase_ratio': [0.1],
    'title_word_count': [8],
    'has_clickbait_words': [1],
    'has_numbers': [1],
    'has_question': [0],
    'has_exclamation': [1],
    'is_recent': [1],
    'is_short_video': [0],
    'is_long_video': [0],
    'video_age_category_encoded': [1],
    'duration_category_encoded': [2],
    'title_length_category_encoded': [1]
})
```

### Step 3: Make Predictions

```python
# Predict log(views)
log_views_prediction = model.predict(new_video)
# Convert back to actual view count
predicted_views = np.expm1(log_views_prediction[0])
print(f"Predicted views: {predicted_views:,.0f}")
# Example output: Predicted views: 245,678
```

### Step 4: Batch Predictions

```python
# For multiple videos
videos_df = pd.read_csv('new_videos.csv')
# ... prepare features ...
predictions = model.predict(videos_df)
predicted_views = np.expm1(predictions)
# Add to dataframe
videos_df['predicted_views'] = predicted_views
videos_df.to_csv('predictions.csv', index=False)
```

## *Feature Engineering for New Data*

To properly prepare new data, use the feature engineering pipeline:

```
from src.feature_engineer import FeatureEngineer
# Initialize engineer
engineer = FeatureEngineer()
# Load and process new data
new_data = pd.read_csv('new_videos_raw.csv')
# For API data
features = engineer.engineer_features_api(new_data)
X, y = engineer.prepare_for_modeling_api(features)
# Make predictions
predictions = model.predict(X)
predicted_views = np.expm1(predictions)
```

# 3. Data Collection

## 3.1 Used Tools

### Web Scraping Tools

### YouTube API Tools

## 3.2 Collected Attributes

### Scraped Data Attributes (8 attributes)

| Attribute | Type | Description | Example |
|-----------|------|-------------|---------|
| video_id | String | Unique YouTube video identifier | "_uQrJ0TkZlc" |
| title | String | Video title | "Python Full Course for Beginners" |
| channel_name | String | Channel name | "Programming with Mosh" |
| views | Integer | Total view count | 46,000,000 |
| upload_date | String | Relative upload date | "6 years ago" |
| duration_seconds | Integer | Video length in seconds | 14,400 (4 hours) |
| video_url | String | Full YouTube URL | "https://youtube.com/watch?v=..." |
| thumbnail_url | String | Thumbnail image URL | "https://i.ytimg.com/..." |

```
video_id
```

```
title
```

```
channel_name
```

```
views
```

```
upload_date
```

```
duration_seconds
```

```
video_url
```

```
thumbnail_url
```

## API Data Attributes (18 attributes)

| Attribute | Type | Description | Example |
|-----------|------|-------------|---------|
| video_id | String | Unique video identifier | "StLVaPZDYgc" |
| title | String | Video title | "Qiuyuan Resonator Showcase" |
| description | String | Video description | "Full description text..." |
| published_at | DateTime | Exact publish timestamp | "2025-01-15T10:30:00Z" |
| channel_id | String | Channel identifier | "UCxxxxxxxxxxxxxxx" |
| channel_name | String | Channel name | "agekk" |
| view_count | Integer | Total views | 12,431 |
| like_count | Integer | Total likes | 1,000 |
| comment_count | Integer | Total comments | 86 |
| duration_seconds | Integer | Video length | 983 |
| tags | String | Video tags (pipe-separated) | "gaming |
| category_id | String | YouTube category ID | "20" |
| category_name | String | Category name | "Gaming" |
| channel_subscriber_count | Integer | Channel subscribers | 29,400 |
| channel_video_count | Integer | Total videos on channel | 801 |
| channel_total_views | Integer | Channel total views | 5,000,000 |
| default_language | String | Video language | "en" |
| thumbnail_url | String | Thumbnail URL | "https://i.ytimg.com/..." |

```
video_id
```

```
title
```

```
description
```

```
published_at
```

```
channel_id
```

```
channel_name
```

```
view_count
```

```
like_count
```

```
comment_count
```

```
duration_seconds
```

```
tags
```

```
category_id
```

```
category_name

channel_subscriber_count

channel_video_count

channel_total_views

default_language

thumbnail_url
```

## 3.3 Number of Data Samples

| Dataset | Collected | After Cleaning | Final Features |
|---------|-----------|----------------|----------------|
| Scraped Data | 3,050 | 3,034 (99.5%) | 3,029 (99.3%) |
| API Data | 2,972 | 2,911 (98.0%) | 2,907 (97.8%) |
| Total | 6,022 | 5,945 (98.7%) | 5,936 (98.6%) |

**Data Loss Analysis:** - Scraped: 21 videos removed (16 duplicates, 5 with missing critical fields) - API: 61 videos removed (duplicates, outliers, missing data) - Feature engineering: 9 additional videos removed (NaN target values from 0 views)

## 3.4 API Usage Methodology

### Authentication

```
from googleapiclient.discovery import build
import os
from dotenv import load_dotenv
# Load API key from environment
load_dotenv()
api_key = os.getenv('YOUTUBE_API_KEY')
# Initialize YouTube API client
youtube = build('youtube', 'v3', developerKey=api_key)
```

### Search Request

```
# Search for videos (costs 100 quota units per request)
search_request = youtube.search().list(
    part='snippet',
    q='python tutorial',
    type='video',
    maxResults=50,
    order='relevance'
)
search_response = search_request.execute()
```

### Video Details Request

```
# Get detailed statistics (costs 1 quota unit per video)
video_ids = [item['id']['videoId'] for item in search_response['items']]
video_request = youtube.videos().list(
    part='snippet,statistics,contentDetails',
    id=','.join(video_ids)
)
video_response = video_request.execute()
```

### Channel Details Request

```
# Get channel metadata (costs 1 quota unit per channel)
channel_request = youtube.channels().list(
    part='snippet,statistics',
    id=channel_id
)
channel_response = channel_request.execute()
```

### Quota Management

**Calculation for 3,000 videos:** - 60 search queries × 100 units = 6,000 units - 3,000 videos × 1 unit = 3,000 units - ~500 unique channels × 1 unit = 500 units - **Total:** ~9,500 units (within daily limit)

**Strategy:** Collected over 3 days to stay within quota

## 3.5 Sample Data After Preprocessing

### Sample 1: Scraped Data

```
Video ID: _uQrJ0TkZlc
Title: Python Full Course for Beginners
Channel: Unknown
Views: 46,000,000
Duration: 0 seconds (data extraction issue)
Days Since Upload: 2190
Title Length: 32
Uppercase Ratio: 0.143
Word Count: 5
Views Per Day: 21,004
```

### Sample 2: API Data

```
Video ID: StLVaPZDYgc
Title: I'LL HONE HIS BLADE....... Qiuyuan Resonator Showcase REACTION | Wuthering Waves
Channel: agekk
Views: 12,431
Likes: 1,000
Comments: 86
Duration: 983 seconds (16.4 minutes)
Days Since Upload: 0 (just published)
Category: Gaming
Channel Subscribers: 29,400
Engagement Rate: 0.0874 (8.74%)
Tags: 10
Description Length: 500 characters
```

# 4. Data Preprocessing

## 4.1 Data Cleaning Methodology

### Scraped Data Cleaning (src/preprocessor.py)

Scraped Data Cleaning (`src/preprocessor.py`)

### Step 1: Duplicate Removal

```
# Remove duplicates based on video_id
df = df.drop_duplicates(subset=['video_id'], keep='first')
# Removed: 16 duplicates
```

### Step 2: Missing Value Handling

```
# Remove rows with missing critical fields
df = df.dropna(subset=['video_id', 'title'])
# Clean views (ensure numeric)
df['views'] = pd.to_numeric(df['views'], errors='coerce').fillna(0).astype(int)
# Clean duration
df['duration_seconds'] = pd.to_numeric(df['duration_seconds'], errors='coerce').fillna(0).astype(int)
```

### Step 3: Date Parsing

```
# Parse relative dates ("2 days ago" → datetime)
df['upload_date_parsed'] = df['upload_date'].apply(parse_relative_date)
df['days_since_upload'] = (datetime.now() - df['upload_date_parsed']).days
```

### Step 4: Outlier Removal

```
# Remove extreme outliers (99.5th percentile threshold)
view_threshold = df['views'].quantile(0.995)
df = df[df['views'] <= view_threshold]
# Remove unusually long videos (> 4 hours)
df = df[df['duration_seconds'] <= 14400]
```

**Results:** - Initial: 3,050 records - After cleaning: 3,034 records (99.5% retention)

## *API Data Cleaning*

### Step 1: Duplicate Removal

```
df = df.drop_duplicates(subset=['video_id'], keep='first')
# Removed: 0 duplicates (API ensures uniqueness)
```

### Step 2: Numeric Field Validation

```
numeric_fields = [
    'view_count', 'like_count', 'comment_count',
    'duration_seconds', 'channel_subscriber_count'
]
for field in numeric_fields:
    df[field] = pd.to_numeric(df[field], errors='coerce').fillna(0).astype(int)
```

### Step 3: Date Parsing

```
# Parse ISO 8601 timestamps
df['published_at'] = pd.to_datetime(df['published_at'], errors='coerce')
df['days_since_upload'] = (pd.Timestamp.now(tz='UTC') - df['published_at']).days
```

### Step 4: Engagement Metrics

```
# Calculate engagement rate
df['engagement_rate'] = (df['like_count'] + df['comment_count']) / df['view_count']
df['engagement_rate'] = df['engagement_rate'].fillna(0)
```

### Step 5: Outlier Removal

```
# Remove top 0.5% outliers
view_threshold = df['view_count'].quantile(0.995)
df = df[df['view_count'] <= view_threshold]
df = df[df['duration_seconds'] <= 14400]
```

**Results:** - Initial: 2,972 records - After cleaning: 2,911 records (98.0% retention)

## 4.2 Normalization and Standardization

### Text Normalization

```
# Standardize text fields
df['title'] = df['title'].astype(str).str.strip()
df['channel_name'] = df['channel_name'].fillna('Unknown').astype(str)
# Calculate text features
df['title_length'] = df['title'].str.len()
df['title_uppercase_ratio'] = df['title'].apply(calculate_uppercase_ratio)
df['title_word_count'] = df['title'].apply(lambda x: len(str(x).split()))
```

### Numeric Normalization

```
# Create normalized duration in minutes
df['duration_minutes'] = df['duration_seconds'] / 60
# Calculate views per day (normalized popularity)
df['views_per_day'] = np.where(
    df['days_since_upload'] > 0,
    df['views'] / df['days_since_upload'],
    df['views']
)
```

### Category Mapping

```
# Map YouTube category IDs to names
category_mapping = {
    '1': 'Film & Animation',
    '10': 'Music',
    '20': 'Gaming',
    '22': 'People & Blogs',
    '23': 'Comedy',
    '24': 'Entertainment',
    '25': 'News & Politics',
    '26': 'Howto & Style',
    '27': 'Education',
    '28': 'Science & Technology'
}
df['category_name'] = df['category_id'].astype(str).map(category_mapping)
```

## 4.3 Data Quality Summary

| Metric | Scraped Data | API Data |
|---|---|---|
| Original Records | 3,050 | 2,972 |
| Duplicates Removed | 16 | 0 |
| Missing Critical Fields | 0 | 0 |
| Outliers Removed | 0 | 61 |
| Final Records | 3,034 | 2,911 |
| Retention Rate | 99.5% | 98.0% |
| Missing Values (%) | 0.2% | 0.1% |

# 5. Feature Engineering

## 5.1 Data Processing Pipeline

After loading preprocessed data from `data/processed/`, features are created through systematic transformation and derivation:

```
from src.feature_engineer import FeatureEngineer
# Initialize feature engineer
engineer = FeatureEngineer()
# Load preprocessed data
df = pd.read_csv('data/processed/api_processed.csv')
# Apply feature engineering
features_df = engineer.engineer_features_api(df)
# Prepare for modeling
X, y = engineer.prepare_for_modeling_api(features_df)
```

## 5.2 Feature Categories

### Time-Based Features

```
# Video age categories
bins = [-1, 7, 30, 90, 365, np.inf]
labels = ['very_recent', 'recent', 'month_old', 'quarter_old', 'old']
df['video_age_category'] = pd.cut(df['days_since_upload'], bins=bins, labels=labels)
# Recency indicator
df['is_recent'] = (df['days_since_upload'] <= 30).astype(int)
# Log transformation to reduce skewness
df['log_days_since_upload'] = np.log1p(df['days_since_upload'])
```

### Duration-Based Features

```
# Duration categories
bins = [-1, 1, 5, 10, 20, np.inf]
labels = ['very_short', 'short', 'medium', 'long', 'very_long']
df['duration_category'] = pd.cut(df['duration_minutes'], bins=bins, labels=labels)
# Short video indicator (YouTube Shorts)
df['is_short_video'] = (df['duration_seconds'] < 60).astype(int)
# Long video indicator
df['is_long_video'] = (df['duration_minutes'] > 20).astype(int)
# Log transformation
df['log_duration'] = np.log1p(df['duration_seconds'])
```

### Title-Based Features

```
# Clickbait indicators
clickbait_words = ['how to', 'top 10', 'best', 'worst', 'shocking', 'amazing']
df['has_clickbait_words'] = df['title'].str.lower().apply(
    lambda x: any(word in str(x) for word in clickbait_words)
).astype(int)
# Punctuation indicators
df['has_numbers'] = df['title'].str.contains(r'\d', regex=True).astype(int)
df['has_question'] = df['title'].str.contains(r'\?').astype(int)
df['has_exclamation'] = df['title'].str.contains(r'!').astype(int)
# Title length categories
bins = [-1, 30, 50, 70, np.inf]
labels = ['short', 'medium', 'long', 'very_long']
df['title_length_category'] = pd.cut(df['title_length'], bins=bins, labels=labels)
```

### *Engagement Features (API Only)*

```
# Log-transformed engagement metrics
df['log_likes'] = np.log1p(df['likes'])
df['log_comments'] = np.log1p(df['comments'])
# Engagement score (weighted)
df['engagement_score'] = (
    df['likes_per_view'] * 0.6 +
    df['comments_per_view'] * 0.4
)
# High engagement indicator
threshold = df['engagement_rate'].quantile(0.75)
df['is_high_engagement'] = (df['engagement_rate'] > threshold).astype(int)
# Comment-to-like ratio
df['comment_like_ratio'] = np.where(
    df['likes'] > 0,
    df['comments'] / df['likes'],
    0
)
```

### *Channel Features (API Only)*

```
# Log of subscribers
df['log_channel_subscribers'] = np.log1p(df['channel_subscribers'])
# Channel size categories
bins = [-1, 1000, 10000, 100000, 1000000, np.inf]
labels = ['micro', 'small', 'medium', 'large', 'mega']
df['channel_size'] = pd.cut(df['channel_subscribers'], bins=bins, labels=labels)
# Average views per video
df['channel_avg_views_per_video'] = np.where(
    df['channel_video_count'] > 0,
    df['views'] / df['channel_video_count'],
    0
)
```

### *View Features (Target and Derived)*

```
# Target variable (log-transformed)
df['log_views'] = np.log1p(df['views'])
# View velocity (accounts for time decay)
df['view_velocity'] = np.where(
    df['days_since_upload'] > 0,
    df['views'] / np.sqrt(df['days_since_upload']),
    df['views']
)
# High view indicator
threshold = df['views'].quantile(0.75)
df['is_high_view'] = (df['views'] > threshold).astype(int)
```

## *5.3 Categorical Encoding*

```
from sklearn.preprocessing import LabelEncoder
# Encode categorical features for ML
categorical_features = [
    'video_age_category',
    'duration_category',
    'title_length_category',
    'channel_size',  # API only
    'category_name'  # API only
]
for cat_col in categorical_features:
    le = LabelEncoder()
    df[f'{cat_col}_encoded'] = le.fit_transform(df[cat_col].astype(str))
```

## *5.4 Final Feature Sets*

### *Scraped Data Features (17 features)*

**Numerical Features (14):** 1. `duration_seconds` - Video length in seconds 2. `log_duration` - Log-transformed duration 3. `days_since_upload` - Days since publication 4. `log_days_since_upload` - Log-transformed age 5. `title_length` - Number of characters 6. `title_uppercase_ratio` - Ratio of uppercase letters 7. `title_word_count` - Number of words 8. `has_clickbait_words` - Binary (0/1) 9. `has_numbers` - Binary (0/1) 10. `has_question` - Binary (0/1) 11. `has_exclamation` - Binary (0/1) 12. `is_recent` - Binary (0/1) 13. `is_short_video` - Binary (0/1) 14. `is_long_video` - Binary (0/1)

Numerical Features (14): 1. `duration_seconds` - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. `log_duration` - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. `days_since_upload` - Days since publication 4. `log_days_since_upload` - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. `log_days_since_upload` - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. `title_length` - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. title_length - Number of characters 6. `title_uppercase_ratio` - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. `title_word_count` - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. `has_clickbait_words` - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1)

14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. `has_numbers` - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. `has_question` - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. `has_exclamation` - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. `is_recent` - Binary (0/1) 13. is_short_video - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. `is_short_video` - Binary (0/1) 14. is_long_video - Binary (0/1)

Numerical Features (14): 1. duration_seconds - Video length in seconds 2. log_duration - Log-transformed duration 3. days_since_upload - Days since publication 4. log_days_since_upload - Log-transformed age 5. title_length - Number of characters 6. title_uppercase_ratio - Ratio of uppercase letters 7. title_word_count - Number of words 8. has_clickbait_words - Binary (0/1) 9. has_numbers - Binary (0/1) 10. has_question - Binary (0/1) 11. has_exclamation - Binary (0/1) 12. is_recent - Binary (0/1) 13. is_short_video - Binary (0/1) 14. `is_long_video` - Binary (0/1)

**Categorical Features (3):** 15. `video_age_category_encoded` - Age category (0-4) 16. `duration_category_encoded` - Duration category (0-4) 17. `title_length_category_encoded` - Title length category (0-3)

Categorical Features (3): 15. `video_age_category_encoded` - Age category (0-4) 16. duration_category_encoded - Duration category (0-4) 17. title_length_category_encoded - Title length category (0-3)

Categorical Features (3): 15. video_age_category_encoded - Age category (0-4) 16. `duration_category_encoded` - Duration category (0-4) 17. title_length_category_encoded - Title length category (0-3)

Categorical Features (3): 15. video_age_category_encoded - Age category (0-4) 16. duration_category_encoded - Duration category (0-4) 17. `title_length_category_encoded` - Title length category (0-3)


### *API Data Features (28 features)*

**All scraped features (17) plus:**

**Engagement Features (6):** 18. `log_likes` - Log-transformed likes 19. `log_comments` - Log-transformed comments 20. `engagement_rate` - (likes + comments) / views 21. `engagement_score` - Weighted engagement 22. `likes_per_view` - Like ratio 23. `comments_per_view` - Comment ratio

Engagement Features (6): 18. `log_likes` - Log-transformed likes 19. log_comments - Log-transformed comments 20. engagement_rate - (likes + comments) / views 21. engagement_score - Weighted engagement 22. likes_per_view - Like ratio 23. comments_per_view - Comment ratio

Engagement Features (6): 18. log_likes - Log-transformed likes 19. `log_comments` - Log-transformed comments 20. engagement_rate - (likes + comments) / views 21. engagement_score - Weighted engagement 22. likes_per_view - Like ratio 23. comments_per_view - Comment ratio

Engagement Features (6): 18. log_likes - Log-transformed likes 19. log_comments - Log-transformed comments 20. `engagement_rate` - (likes + comments) / views 21. engagement_score - Weighted engagement 22. likes_per_view - Like ratio 23. comments_per_view - Comment ratio

Engagement Features (6): 18. log_likes - Log-transformed likes 19. log_comments - Log-transformed comments 20. engagement_rate - (likes + comments) / views 21. `engagement_score` - Weighted engagement 22. likes_per_view - Like ratio 23. comments_per_view - Comment ratio

Engagement Features (6): 18. log_likes - Log-transformed likes 19. log_comments - Log-transformed comments 20. engagement_rate - (likes + comments) / views 21. engagement_score - Weighted engagement 22. `likes_per_view` - Like ratio 23. comments_per_view - Comment ratio

Engagement Features (6): 18. log_likes - Log-transformed likes 19. log_comments - Log-transformed comments 20. engagement_rate - (likes + comments) / views 21. engagement_score - Weighted engagement 22. likes_per_view - Like ratio 23. `comments_per_view` - Comment ratio

**Content Features (2):** 24. `tag_count` - Number of tags 25. `description_length` - Description character count

Content Features (2): 24. `tag_count` - Number of tags 25. description_length - Description character count

Content Features (2): 24. tag_count - Number of tags 25. `description_length` - Description character count

**Channel Features (3):** 26. `log_channel_subscribers` - Log-transformed subscribers 27. `channel_video_count` - Total channel videos 28. `channel_avg_views_per_video` - Average performance

Channel Features (3): 26. `log_channel_subscribers` - Log-transformed subscribers 27. channel_video_count - Total channel videos 28. channel_avg_views_per_video - Average performance

Channel Features (3): 26. log_channel_subscribers - Log-transformed subscribers 27. `channel_video_count` - Total channel videos 28. channel_avg_views_per_video - Average performance

Channel Features (3): 26. log_channel_subscribers - Log-transformed subscribers 27. channel_video_count - Total channel videos 28. `channel_avg_views_per_video` - Average performance

**Additional Categorical (2):** 29. `channel_size_encoded` - Channel size category (0-4) 30. `category_name_encoded` - Video category (0-14)

Additional Categorical (2): 29. `channel_size_encoded` - Channel size category (0-4) 30. category_name_encoded - Video category (0-14)

Additional Categorical (2): 29. channel_size_encoded - Channel size category (0-4) 30. `category_name_encoded` - Video category (0-14)


## 5.5 Three Sample Data with Features

### *Sample 1: High-Performing Educational Video (Scraped)*

**Raw Data:** - Video ID: `_uQrJ0TkZlc` - Title: "Python Full Course for Beginners" - Views: 46,000,000 - Duration: 0 seconds (missing) - Days Since Upload: 2190

**Engineered Features:** - `log_views`: 17.64 - `log_days_since_upload`: 7.69 - `is_recent`: 0 - `has_clickbait_words`: 0 - `title_length`: 32 - `title_uppercase_ratio`: 0.143 - `title_word_count`: 5 - `video_age_category_encoded`: 4 (old)

### *Sample 2: Recent Gaming Video (API)*

**Raw Data:** - Video ID: `StLVaPZDYgc` - Title: "I'LL HONE HIS BLADE....... Qiuyuan Resonator Showcase REACTION | Wuthering Waves" - Views: 12,431 - Likes: 1,000 - Comments: 86 - Duration: 983 seconds - Days Since Upload: 0

**Engineered Features:** - `log_views`: 9.43 - `log_likes`: 6.91 - `log_comments`: 4.46 - `engagement_rate`: 0.0874 - `likes_per_view`: 0.0804 - `comments_per_view`: 0.0069 - `log_channel_subscribers`: 10.29 - `is_recent`: 1 - `has_clickbait_words`: 0 - `has_numbers`: 0 - `video_age_category_encoded`: 0 (very_recent) - `channel_size_encoded`: 2 (medium)

Engineered Features: - `log_views`: 9.43 - log_likes: 6.91 - log_comments: 4.46 - engagement_rate: 0.0874 - likes_per_view: 0.0804 - comments_per_view: 0.0069 - log_channel_subscribers: 10.29 - is_recent: 1 - has_clickbait_words: 0 - has_numbers: 0 - video_age_category_encoded: 0 (very_recent) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 9.43 - `log_likes`: 6.91 - log_comments: 4.46 - engagement_rate: 0.0874 - likes_per_view: 0.0804 - comments_per_view: 0.0069 - log_channel_subscribers: 10.29 - is_recent: 1 - has_clickbait_words: 0 - has_numbers: 0 - video_age_category_encoded: 0 (very_recent) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 9.43 - log_likes: 6.91 - `log_comments`: 4.46 - engagement_rate: 0.0874 - likes_per_view: 0.0804 - comments_per_view: 0.0069 - log_channel_subscribers: 10.29 - is_recent: 1 - has_clickbait_words: 0 - has_numbers: 0 - video_age_category_encoded: 0 (very_recent) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 9.43 - log_likes: 6.91 - log_comments: 4.46 - `engagement_rate`: 0.0874 - likes_per_view: 0.0804 - comments_per_view: 0.0069 - log_channel_subscribers: 10.29 - is_recent: 1 - has_clickbait_words: 0 - has_numbers: 0 - video_age_category_encoded: 0 (very_recent) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 9.43 - log_likes: 6.91 - log_comments: 4.46 - engagement_rate: 0.0874 - `likes_per_view`: 0.0804 - comments_per_view: 0.0069 - log_channel_subscribers: 10.29 - is_recent: 1 - has_clickbait_words: 0 - has_numbers: 0 - video_age_category_encoded: 0 (very_recent) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 9.43 - log_likes: 6.91 - log_comments: 4.46 - engagement_rate: 0.0874 - likes_per_view: 0.0804 - `comments_per_view`: 0.0069 - log_channel_subscribers: 10.29 - is_recent: 1 - has_clickbait_words: 0 - has_numbers: 0 - video_age_category_encoded: 0 (very_recent) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 9.43 - log_likes: 6.91 - log_comments: 4.46 - engagement_rate: 0.0874 - likes_per_view: 0.0804 - comments_per_view: 0.0069 - `log_channel_subscribers`: 10.29 - is_recent: 1 - has_clickbait_words: 0 - has_numbers: 0 - video_age_category_encoded: 0 (very_recent) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 9.43 - log_likes: 6.91 - log_comments: 4.46 - engagement_rate: 0.0874 - likes_per_view: 0.0804 - comments_per_view: 0.0069 - log_channel_subscribers: 10.29 - `is_recent`: 1 - has_clickbait_words: 0 - has_numbers: 0 - video_age_category_encoded: 0 (very_recent) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 9.43 - log_likes: 6.91 - log_comments: 4.46 - engagement_rate: 0.0874 - likes_per_view: 0.0804 - comments_per_view: 0.0069 - log_channel_subscribers: 10.29 - is_recent: 1 - `has_clickbait_words`: 0 - has_numbers: 0 - video_age_category_encoded: 0 (very_recent) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 9.43 - log_likes: 6.91 - log_comments: 4.46 - engagement_rate: 0.0874 - likes_per_view: 0.0804 - comments_per_view: 0.0069 - log_channel_subscribers: 10.29 - is_recent: 1 - has_clickbait_words: 0 - `has_numbers`: 0 - video_age_category_encoded: 0 (very_recent) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 9.43 - log_likes: 6.91 - log_comments: 4.46 - engagement_rate: 0.0874 - likes_per_view: 0.0804 - comments_per_view: 0.0069 - log_channel_subscribers: 10.29 - is_recent: 1 - has_clickbait_words: 0 - has_numbers: 0 - `video_age_category_encoded`: 0 (very_recent) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 9.43 - log_likes: 6.91 - log_comments: 4.46 - engagement_rate: 0.0874 - likes_per_view: 0.0804 - comments_per_view: 0.0069 - log_channel_subscribers: 10.29 - is_recent: 1 - has_clickbait_words: 0 - has_numbers: 0 - video_age_category_encoded: 0 (very_recent) - `channel_size_encoded`: 2 (medium)

### Sample 3: Viral People & Blogs Video (API)

**Raw Data:** - Video ID: `rV2CJrpTYE8` - Title: "JUICY TIKTOK STORY TIME COMPILATION !! *part 135*" - Views: 45,802 - Likes: 517 - Comments: 45 - Duration: 1253 seconds - Days Since Upload: 67

**Engineered Features:** - `log_views`: 10.73 - `log_likes`: 6.25 - `log_comments`: 3.83 - `engagement_rate`: 0.0123 - `likes_per_view`: 0.0113 - `comments_per_view`: 0.00098 - `log_channel_subscribers`: 10.10 - `is_recent`: 0 - `has_clickbait_words`: 0 - `has_numbers`: 1 - `has_exclamation`: 1 - `title_uppercase_ratio`: 0.525 - `video_age_category_encoded`: 2 (month_old) - `channel_size_encoded`: 2 (medium)

Engineered Features: - log_views: 10.73 - log_likes: 6.25 - log_comments: 3.83 - engagement_rate: 0.0123 - likes_per_view: 0.0113 - comments_per_view: 0.00098 - log_channel_subscribers: 10.10 - is_recent: 0 - has_clickbait_words: 0 - has_numbers: 1 - `has_exclamation`: 1 - title_uppercase_ratio: 0.525 - video_age_category_encoded: 2 (month_old) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 10.73 - log_likes: 6.25 - log_comments: 3.83 - engagement_rate: 0.0123 - likes_per_view: 0.0113 - comments_per_view: 0.00098 - log_channel_subscribers: 10.10 - is_recent: 0 - has_clickbait_words: 0 - has_numbers: 1 - has_exclamation: 1 - `title_uppercase_ratio`: 0.525 - video_age_category_encoded: 2 (month_old) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 10.73 - log_likes: 6.25 - log_comments: 3.83 - engagement_rate: 0.0123 - likes_per_view: 0.0113 - comments_per_view: 0.00098 - log_channel_subscribers: 10.10 - is_recent: 0 - has_clickbait_words: 0 - has_numbers: 1 - has_exclamation: 1 - title_uppercase_ratio: 0.525 - `video_age_category_encoded`: 2 (month_old) - channel_size_encoded: 2 (medium)

Engineered Features: - log_views: 10.73 - log_likes: 6.25 - log_comments: 3.83 - engagement_rate: 0.0123 - likes_per_view: 0.0113 - comments_per_view: 0.00098 - log_channel_subscribers: 10.10 - is_recent: 0 - has_clickbait_words: 0 - has_numbers: 1 - has_exclamation: 1 - title_uppercase_ratio: 0.525 - video_age_category_encoded: 2 (month_old) - `channel_size_encoded`: 2 (medium)

# 6. Model Development and Evaluation

## 6.1 Train and Test Data Partition

**Methodology:**

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(
    X, y,
    test_size=0.2,      # 80/20 split
    random_state=42      # Reproducible split
)
```

**Split Statistics:**

| Dataset | Total | Train (80%) | Test (20%) |
|---------|-------|-------------|------------|
| Scraped | 3,029 | 2,423 | 606 |
| API | 2,907 | 2,325 | 582 |

**Rationale for 80/20 Split:** - Standard practice for datasets of this size (3,000+ samples) - Provides sufficient training data for ensemble methods - Reserves adequate test set for reliable performance evaluation - No stratification needed (continuous target variable)

## 6.2 Model 1: Scraped Data Models

### 6.2.1 Random Forest Regressor

**Machine Learning Model:** Random Forest Regressor

**Model Justification:** - Handles non-linear relationships without explicit feature engineering - Robust to outliers and missing values - Provides feature importance metrics - Excellent for tabular data with mixed

feature types - Ensemble method reduces overfitting through bagging

**Hyperparameters:**

```
RandomForestRegressor(
    n_estimators=200,        # 200 trees for stable predictions
    max_depth=20,            # Limit depth to prevent overfitting
    min_samples_split=5,     # Require 5 samples to split
    min_samples_leaf=2,      # Require 2 samples per leaf
    max_features='sqrt',     # Use sqrt(n_features) per split
    random_state=42,         # Reproducibility
    n_jobs=-1                # Use all CPU cores
)
```

**Input to Model:** - Shape: (2,423 samples × 17 features) - Target: `log_views` (log-transformed view count)

Input to Model: - Shape: (2,423 samples × 17 features) - Target: `log_views` (log-transformed view count)

**Attributes (17 features):** 1. `duration_seconds` 2. `log_duration` 3. `days_since_upload` 4. `log_days_since_upload` 5. `title_length` 6. `title_uppercase_ratio` 7. `title_word_count` 8. `has_clickbait_words` 9. `has_numbers` 10. `has_question` 11. `has_exclamation` 12. `is_recent` 13. `is_short_video` 14. `is_long_video` 15. `video_age_category_encoded` 16. `duration_category_encoded` 17. `title_length_category_encoded`

Attributes (17 features): 1. `duration_seconds` 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. `log_duration` 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. `days_since_upload` 4. `log_days_since_upload` 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. `log_days_since_upload` 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. `title_length` 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. `title_length_`category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. `title_uppercase_ratio` 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. `title_word_count` 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. `has_clickbait_words` 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. `has_numbers` 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. `has_question` 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. `has_exclamation` 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. `is_recent` 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. `is_short_video` 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. `is_long_video` 15. video_age_category_encoded 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. `video_age_category_encoded` 16. duration_category_encoded 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. `duration_category_encoded` 17. title_length_category_encoded

Attributes (17 features): 1. duration_seconds 2. log_duration 3. days_since_upload 4. log_days_since_upload 5. title_length 6. title_uppercase_ratio 7. title_word_count 8. has_clickbait_words 9. has_numbers 10. has_question 11. has_exclamation 12. is_recent 13. is_short_video 14. is_long_video 15. video_age_category_encoded 16. duration_category_encoded 17. `title_length_category_encoded`

**Size of Train Data:** 2,423 samples

**Performance with Training Data:** - RMSE: 1.5640 - MAE: 1.1500 - R²: **0.7859** - MAPE: 636.42%

**Performance with Test Data:** - RMSE: 2.5565 - MAE: 1.8903 - R²: **0.5663** - MAPE: 11,207.00%

**Cross-Validation Performance (5-fold):** - Mean R²: **0.5213** - Standard Deviation: ±0.0295 - Scores: [0.5206, 0.5533, 0.5536, 0.5012, 0.4778]

**Interpretation:** - Moderate generalization (test R² = 0.5663) - Gap between train/test indicates some overfitting - Cross-validation confirms consistent performance - Limited by lack of engagement metrics in scraped data

### 6.2.2 XGBoost Regressor

**Machine Learning Model:** XGBoost (Extreme Gradient Boosting) Regressor

**Model Justification:** - State-of-the-art gradient boosting algorithm - Handles missing values internally - Regularization prevents overfitting - Faster training than Random Forest - Often superior performance on structured data

**Hyperparameters:**

```
XGBRegressor(
    n_estimators=200,        # 200 boosting rounds
    max_depth=8,             # Tree depth (shallower than RF)
    learning_rate=0.1,       # Conservative learning rate
    subsample=0.8,           # 80% row sampling per tree
    colsample_bytree=0.8,    # 80% column sampling per tree
    random_state=42,         # Reproducibility
    n_jobs=-1,               # Parallel processing
    verbosity=0              # Silent training
)
```

**Input to Model:** - Shape: (2,423 samples × 17 features) - Target: `log_views`

Input to Model: - Shape: (2,423 samples × 17 features) - Target: `log_views`

**Attributes:** Same 17 features as Random Forest

**Size of Train Data:** 2,423 samples

**Performance with Training Data:** - RMSE: 0.4889 - MAE: 0.3485 - R²: **0.9791** (near-perfect fit) - MAPE: 40.74%

**Performance with Test Data:** - RMSE: 2.7129 - MAE: 2.0036 - R²: **0.5116** - MAPE: 19,211.17%

**Cross-Validation Performance (5-fold):** - Mean R²: **0.4376** - Standard Deviation: ±0.0331 - Scores: [0.4254, 0.4563, 0.4893, 0.4262, 0.3908]

**Interpretation:** - Severe overfitting (train R² = 0.98, test R² = 0.51) - Lower test performance than Random Forest - XGBoost memorizes training data with limited features - Random Forest better generalizes for scraped data

**Model Comparison (Scraped Data):**

| Metric | Random Forest | XGBoost | Winner |
|--------|---------------|---------|--------|
| Test R² | 0.5663 | 0.5116 | RF |
| Test RMSE | 2.5565 | 2.7129 | RF |
| Test MAE | 1.8903 | 2.0036 | RF |
| CV Mean R² | 0.5213 | 0.4376 | RF |
| Train/Test Gap | 0.22 | 0.47 | RF |

**Verdict:** Random Forest performs better on scraped data

## 6.3 Model 2: API Data Models

### 6.3.1 Random Forest Regressor

**Machine Learning Model:** Random Forest Regressor

**Model Justification:** Same as scraped model

**Hyperparameters:** Same configuration as scraped model

**Input to Model:** - Shape: (2,325 samples × 28 features) - Target: `log_views`

**Attributes (28 features):**

**Base Features (17):** Same as scraped model

**Additional Features (11):** 18. `log_likes` 19. `log_comments` 20. `engagement_rate` 21. `engagement_score` 22. `tag_count` 23. `description_length` 24. `log_channel_subscribers` 25. `channel_video_count` 26. `likes_per_view` 27. `comments_per_view` 28. `channel_size_encoded` 29. `category_name_encoded`

Additional Features (11): 18. `log_likes` 19. log_comments 20. engagement_rate 21. engagement_score 22. tag_count 23. description_length 24. log_channel_subscribers 25. channel_video_count 26. likes_per_view 27. comments_per_view 28. channel_size_encoded 29. category_name_encoded

Additional Features (11): 18. log_likes 19. `log_comments` 20. engagement_rate 21. engagement_score 22. tag_count 23. description_length 24. log_channel_subscribers 25. channel_video_count 26. likes_per_view 27. comments_per_view 28. channel_size_encoded 29. category_name_encoded

Additional Features (11): 18. log_likes 19. log_comments 20. `engagement_rate` 21. engagement_score 22. tag_count 23. description_length 24. log_channel_subscribers 25. channel_video_count 26. likes_per_view 27. comments_per_view 28. channel_size_encoded 29. category_name_encoded

Additional Features (11): 18. log_likes 19. log_comments 20. engagement_rate 21. `engagement_score` 22. tag_count 23. description_length 24. log_channel_subscribers 25. channel_video_count 26. likes_per_view 27. comments_per_view 28. channel_size_encoded 29. category_name_encoded

Additional Features (11): 18. log_likes 19. log_comments 20. engagement_rate 21. engagement_score 22. `tag_count` 23. description_length 24. log_channel_subscribers 25. channel_video_count 26. likes_per_view 27. comments_per_view 28. channel_size_encoded 29. category_name_encoded

Additional Features (11): 18. log_likes 19. log_comments 20. engagement_rate 21. engagement_score 22. tag_count 23. `description_length` 24. log_channel_subscribers 25. channel_video_count 26. likes_per_view 27. comments_per_view 28. channel_size_encoded 29. category_name_encoded

Additional Features (11): 18. log_likes 19. log_comments 20. engagement_rate 21. engagement_score 22. tag_count 23. description_length 24. `log_channel_subscribers` 25. channel_video_count 26. likes_per_view 27. comments_per_view 28. channel_size_encoded 29. category_name_encoded

Additional Features (11): 18. log_likes 19. log_comments 20. engagement_rate 21. engagement_score 22. tag_count 23. description_length 24. log_channel_subscribers 25. `channel_video_count` 26. likes_per_view 27. comments_per_view 28. channel_size_encoded 29. category_name_encoded

Additional Features (11): 18. log_likes 19. log_comments 20. engagement_rate 21. engagement_score 22. tag_count 23. description_length 24. log_channel_subscribers 25. channel_video_count 26. `likes_per_view` 27. comments_per_view 28. channel_size_encoded 29. category_name_encoded

Additional Features (11): 18. log_likes 19. log_comments 20. engagement_rate 21. engagement_score 22. tag_count 23. description_length 24. log_channel_subscribers 25. channel_video_count 26. likes_per_view 27. `comments_per_view` 28. channel_size_encoded 29. category_name_encoded

Additional Features (11): 18. log_likes 19. log_comments 20. engagement_rate 21. engagement_score 22. tag_count 23. description_length 24. log_channel_subscribers 25. channel_video_count 26. likes_per_view 27. comments_per_view 28. `channel_size_encoded` 29. category_name_encoded

Additional Features (11): 18. log_likes 19. log_comments 20. engagement_rate 21. engagement_score 22. tag_count 23. description_length 24. log_channel_subscribers 25. channel_video_count 26. likes_per_view 27. comments_per_view 28. channel_size_encoded 29. `category_name_encoded`

**Size of Train Data:** 2,325 samples

**Performance with Training Data:** - RMSE: 0.3455 - MAE: 0.1772 - R²: **0.9837** - MAPE: 57.48%

**Performance with Test Data:** - RMSE: 0.8494 - MAE: 0.3593 - R²: **0.9082** - MAPE: 65,070.70%

**Cross-Validation Performance (5-fold):** - Mean R²: **0.9418** - Standard Deviation: ±0.0198 - Scores: [0.9103, 0.9577, 0.9565, 0.9581, 0.9266]

**Interpretation:** - Excellent generalization (test R² = 0.9082) - Small train/test gap (0.08) indicates minimal overfitting - Cross-validation confirms robust performance - Engagement metrics dramatically improve predictive power

### 6.3.2 XGBoost Regressor

**Machine Learning Model:** XGBoost Regressor

**Model Justification:** Same as scraped model

**Hyperparameters:** Same configuration as scraped model

**Input to Model:** - Shape: (2,325 samples × 28 features) - Target: `log_views`

Input to Model: - Shape: (2,325 samples × 28 features) - Target: `log_views`

**Attributes:** Same 28 features as Random Forest

**Size of Train Data:** 2,325 samples

**Performance with Training Data:** - RMSE: 0.0093 - MAE: 0.0069 - R²: **0.9999** (nearly perfect) - MAPE: 0.69%

**Performance with Test Data:** - RMSE: 0.6757 - MAE: 0.1579 - R²: **0.9419** ← **BEST OVERALL** - MAPE: 152,486.36%

**Cross-Validation Performance (5-fold):** - Mean R²: **0.9789** - Standard Deviation: ±0.0180 - Scores: [0.9464, 0.9930, 0.9915, 0.9915, 0.9721]

**Interpretation:** - Best performing model overall (test R² = 0.9419) - Some overfitting (train R² = 0.9999), but excellent test performance - Cross-validation R² (0.9789) is even higher than test - XGBoost excels with rich feature set from API

**Model Comparison (API Data):**

| Metric | Random Forest | XGBoost | Winner |
|---|---|---|---|
| Test R² | 0.9082 | 0.9419 | XGBoost |
| Test RMSE | 0.8494 | 0.6757 | XGBoost |
| Test MAE | 0.3593 | 0.1579 | XGBoost |
| CV Mean R² | 0.9418 | 0.9789 | XGBoost |
| Train/Test Gap | 0.08 | 0.06 | XGBoost |

**Verdict:** XGBoost is the best model for API data

## 6.4 Overall Model Comparison

### Cross-Dataset Comparison

| Dataset | Best Model | Test R² | Test RMSE | Test MAE |
|---------|-----------|---------|-----------|----------|
| Scraped | Random Forest | 0.5663 | 2.5565 | 1.8903 |
| API | XGBoost | 0.9419 | 0.6757 | 0.1579 |
| Improvement | - | +66% | -74% | -92% |

**Key Finding:** API data models outperform scraped models by a massive margin: - R² improvement: 0.5663 → 0.9419 (+66% relative) - RMSE reduction: 2.5565 → 0.6757 (-74%) - MAE reduction: 1.8903 → 0.1579 (-92%)

## 6.5 Feature Importance

### Description of Feature Importance Techniques

**Random Forest Feature Importance:** - **Method:** Mean Decrease in Impurity (Gini importance) - **Calculation:** Average reduction in split criterion across all trees - **Formula:** For each feature, sum the total reduction in node impurity (weighted by probability of reaching that node) across all trees - **Normalization:** Importances sum to 1.0 - **Interpretation:** Higher values indicate features that create purer splits

**XGBoost Feature Importance:** - **Method:** Gain-based importance - **Calculation:** Average gain (improvement in accuracy) when feature is used for splitting - **Formula:** For each split using feature f, measure reduction in loss function, average across all splits - **Normalization:** Importances sum to 1.0 - **Interpretation:** Higher values indicate features that most improve model predictions

**Advantages:** - Model-agnostic interpretation - Captures non-linear relationships - No assumptions about feature distributions

### Feature Importance: Scraped Data (Random Forest)

| Rank | Feature | Importance | Interpretation |
|------|---------|-----------|----------------|
| 1 | days_since_upload | 0.2993 | Video age is the strongest predictor |
| 2 | log_days_since_upload | 0.2348 | Log-transformed age adds complementary info |
| 3 | is_recent | 0.1247 | Recent videos behave differently |
| 4 | title_uppercase_ratio | 0.1195 | Capitalization style matters |
| 5 | title_length | 0.1182 | Title length influences views |
| 6 | title_word_count | 0.0765 | Word count distinct from character count |
| 7 | has_numbers | 0.0122 | Numbers in title have weak signal |
| 8 | has_clickbait_words | 0.0075 | Clickbait has minimal effect |
| 9 | has_exclamation | 0.0073 | Exclamation marks barely help |
| 10+ | Others | 0.0000 | Duration features not predictive |

```
days_since_upload
```

```
log_days_since_upload
```

```
is_recent
```

```
title_uppercase_ratio

title_length

title_word_count

has_numbers

has_clickbait_words

has_exclamation
```

**Insight:** Without engagement metrics, temporal and title features dominate

## *Feature Importance: Scraped Data (XGBoost)*

| Rank | Feature | Importance | Interpretation |
|------|---------|------------|----------------|
| 1 | is_recent | 0.6768 | XGBoost heavily relies on recency |
| 2 | days_since_upload | 0.1094 | Raw days still important |
| 3 | log_days_since_upload | 0.0447 | Log transform less critical |
| 4 | title_uppercase_ratio | 0.0337 | Title style moderate signal |
| 5 | title_length | 0.0337 | Length moderate signal |

```
is_recent

days_since_upload

log_days_since_upload

title_uppercase_ratio

title_length
```

**Insight:** XGBoost oversimplifies to recency (67.7% weight)

## *Feature Importance: API Data (Random Forest)*

| Rank | Feature | Importance | Interpretation |
|------|---------|------------|----------------|
| 1 | log_likes | 0.2900 | Likes are the strongest predictor |
| 2 | log_comments | 0.1350 | Comments second most important |
| 3 | days_since_upload | 0.0929 | Video age still relevant |
| 4 | engagement_rate | 0.0798 | Overall engagement matters |
| 5 | comments_per_view | 0.0759 | Comment ratio important |
| 6 | log_days_since_upload | 0.0751 | Log age complementary |
| 7 | likes_per_view | 0.0568 | Like ratio important |
| 8 | is_recent | 0.0446 | Recency less critical with engagement data |
| 9 | log_channel_subscribers | 0.0334 | Channel size moderate signal |
| 10 | channel_video_count | 0.0300 | Channel activity moderate signal |

```
log_likes
```

```
log_comments
```

```
days_since_upload
```

```
engagement_rate
```

```
comments_per_view
```

```
log_days_since_upload
```

```
likes_per_view
```

```
is_recent
```

```
log_channel_subscribers
```

```
channel_video_count
```

**Top 5 Account for 66.4% of importance**

**Insight:** Engagement metrics (likes, comments) dominate; temporal features less critical

## *Feature Importance: API Data (XGBoost)*

| Rank | Feature | Importance | Interpretation |
|------|---------|------------|----------------|
| 1 | log_likes | 0.5765 | Likes explain 57.7% of predictions |
| 2 | engagement_rate | 0.2147 | Engagement rate second (21.5%) |
| 3 | likes_per_view | 0.0663 | Like ratio tertiary (6.6%) |
| 4 | log_comments | 0.0265 | Comments surprisingly low (2.6%) |
| 5 | comments_per_view | 0.0185 | Comment ratio low (1.9%) |

```
log_likes
```

```
engagement_rate
```

```
likes_per_view
```

```
log_comments
```

```
comments_per_view
```

**Top 5 Account for 89.3% of importance**

**Insight:** XGBoost heavily weights `log_likes` and `engagement_rate`

Insight: XGBoost heavily weights `log_likes` and engagement_rate

Insight: XGBoost heavily weights log_likes and `engagement_rate`

## *Visual Comparison*

**Scraped vs API Feature Importance (Top 5):**

```
    Scraped (Random Forest):
```

```
██████████████████████████████████████
  ■ days_since_upload   ■ 29.9%   ████████████████████
  ■ log_days_since_up.  ■ 23.5%   ███████████████
  ■ is_recent           ■ 12.5%   ████████
  ■ title_uppercase_r.  ■ 12.0%   ████████
  ■ title_length        ■ 11.8%   ████████
██████████████████████████████████████
API (XGBoost):
██████████████████████████████████████
  ■ log_likes           ■ 57.7%   ████████████████████████████████████████
  ■ engagement_rate     ■ 21.5%   ███████████████
  ■ likes_per_view      ■ 6.6%    █████
  ■ log_comments        ■ 2.6%    ██
  ■ comments_per_view   ■ 1.9%    ██
██████████████████████████████████████
```

**Key Takeaway:** Engagement metrics (when available) overwhelmingly dominate view predictions

# 7. Visualization

This section presents the generated visualizations comparing model performance, feature importance, and engagement trends across different dimensions.

## *7.1 Model Performance Comparison*

**File:** `reports/figures/model_comparison.png`

File: `reports/figures/model_comparison.png`

**Description:** Side-by-side bar charts comparing performance metrics (R², RMSE, MAE) across all four models: - Scraped Data: Random Forest vs XGBoost - API Data: Random Forest vs XGBoost

**Key Insights:** - API models achieve R² > 0.90, while scraped models max at 0.57 - XGBoost API model achieves the lowest RMSE (0.68) and MAE (0.16) - Clear visualization of the performance gap between data sources

## *7.2 Feature Importance Visualizations*

### *Scraped Data Feature Importance*

**File:** `reports/figures/feature_importance_scraped.png`

File: `reports/figures/feature_importance_scraped.png`

**Description:** Horizontal bar chart showing top 20 features for scraped data models (both Random Forest and XGBoost), ranked by importance score.

**Key Findings:** - **Random Forest:** Balanced importance across temporal (`days_since_upload`: 29.9%) and title features (`title_uppercase_ratio`: 12.0%) - **XGBoost:** Heavily weights `is_recent` (67.7%), indicating oversimplification - Duration features have zero importance in both models - Without engagement metrics, models rely on easily manipulable features (title styling)

Key Findings: - Random Forest: Balanced importance across temporal (`days_since_upload`: 29.9%) and title features (title_uppercase_ratio: 12.0%) - XGBoost: Heavily weights is_recent (67.7%), indicating oversimplification - Duration features have zero importance in both models - Without engagement metrics, models rely on easily manipulable features (title styling)

Key Findings: - Random Forest: Balanced importance across temporal (days_since_upload: 29.9%) and title features (`title_uppercase_ratio`: 12.0%) - XGBoost: Heavily weights is_recent (67.7%), indicating oversimplification - Duration features have zero importance in both models - Without engagement metrics, models rely on easily manipulable features (title styling)

### *API Data Feature Importance*

**File:** `reports/figures/feature_importance_api.png`

**Description:** Horizontal bar chart showing top 20 features for API data models, highlighting engagement metrics.

**Key Findings: - Random Forest:** `log_likes` (29.0%) and `log_comments` (13.5%) dominate, but more balanced overall - **XGBoost:** `log_likes` alone accounts for 57.7%, `engagement_rate` adds 21.5% - Channel features (`log_channel_subscribers`: 3.3%) have minor impact - Category and content features (tags, description) are negligible

## *7.3 Engagement by Category*

**File:** `reports/figures/engagement_by_category.png`

**Description:** Box plots showing distribution of engagement metrics (views, likes, comments, engagement rate) across YouTube categories for API data.

**Insights:** - **Music** videos have highest median views and likes - **Gaming** videos have high engagement rates despite moderate view counts - **Education** videos show high variance (some viral, most modest) - **News & Politics** have lower engagement rates (less interactive content) - **Entertainment** and **Comedy** have wide IQR, indicating hit-or-miss performance

## 7.4 Engagement by Duration

**File:** `reports/figures/engagement_by_duration.png`

File: `reports/figures/engagement_by_duration.png`

**Description:** Scatter plots and trend lines showing relationship between video duration (minutes) and engagement metrics.

**Insights:** - **Optimal Duration:** 10-15 minutes shows highest engagement - **Short Videos (<1 min):** Lower views but high engagement rate (YouTube Shorts) - **Long Videos (>30 min):** Exponentially decreasing returns; only niche content succeeds - **Sweet Spot:** 8-12 minutes balances watch time and retention - Clear diminishing returns beyond 20 minutes

## 7.5 Engagement by Upload Time

**File:** `reports/figures/engagement_by_upload_time.png`

File: `reports/figures/engagement_by_upload_time.png`

**Description:** Line plots showing how views, likes, and comments evolve with video age (days since upload).

**Insights:** - **Initial Spike:** First 7 days account for ~60% of lifetime views - **Decay Pattern:** Exponential decay in daily engagement after 30 days - **Long-Tail:** Videos older than 1 year still accumulate views (SEO-driven) - **Recency Advantage:** Recent videos (<30 days) have 3x higher views per day - **Stabilization:** After 90 days, view velocity stabilizes

## 7.6 Correlation Heatmaps

### API Data Correlation

**File:** `reports/figures/correlation_heatmap_api.png`

File: `reports/figures/correlation_heatmap_api.png`

**Description:** Heatmap showing Pearson correlation coefficients between numerical features in API dataset.

**Key Correlations:** - **Strong Positive:** - `log_likes` ↔ `log_views` (r = 0.89) - `log_comments` ↔ `log_likes` (r = 0.78) - `log_channel_subscribers` ↔ `log_views` (r = 0.41) - **Moderate Positive:** - `engagement_rate` ↔ `likes_per_view` (r = 0.64) - `description_length` ↔ `tag_count` (r = 0.32) - **Weak/No Correlation:** - `duration_seconds` ↔ `engagement_rate` (r = 0.08) - `days_since_upload` ↔ `engagement_rate` (r = -0.12)

Key Correlations: - Strong Positive: - `log_likes` ↔ log_views (r = 0.89) - log_comments ↔ `log_likes` (r = 0.78) - log_channel_subscribers ↔ log_views (r = 0.41) - Moderate Positive: - engagement_rate ↔ likes_per_view (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ engagement_rate (r = 0.08) - days_since_upload ↔ engagement_rate (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ `log_views` (r = 0.89) - log_comments ↔ log_likes (r = 0.78) - log_channel_subscribers ↔ `log_views` (r = 0.41) - Moderate Positive: - engagement_rate ↔ likes_per_view (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: -

duration_seconds ↔ engagement_rate (r = 0.08) - days_since_upload ↔ engagement_rate (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ log_views (r = 0.89) - `log_comments` ↔ log_likes (r = 0.78) - log_channel_subscribers ↔ log_views (r = 0.41) - Moderate Positive: - engagement_rate ↔ likes_per_view (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ engagement_rate (r = 0.08) - days_since_upload ↔ engagement_rate (r = -0.12)

Key Correlations: - Strong Positive: - `log_likes` ↔ log_views (r = 0.89) - log_comments ↔ `log_likes` (r = 0.78) - log_channel_subscribers ↔ log_views (r = 0.41) - Moderate Positive: - engagement_rate ↔ likes_per_view (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ engagement_rate (r = 0.08) - days_since_upload ↔ engagement_rate (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ log_views (r = 0.89) - log_comments ↔ log_likes (r = 0.78) - `log_channel_subscribers` ↔ log_views (r = 0.41) - Moderate Positive: - engagement_rate ↔ likes_per_view (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ engagement_rate (r = 0.08) - days_since_upload ↔ engagement_rate (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ `log_views` (r = 0.89) - log_comments ↔ log_likes (r = 0.78) - log_channel_subscribers ↔ `log_views` (r = 0.41) - Moderate Positive: - engagement_rate ↔ likes_per_view (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ engagement_rate (r = 0.08) - days_since_upload ↔ engagement_rate (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ log_views (r = 0.89) - log_comments ↔ log_likes (r = 0.78) - log_channel_subscribers ↔ log_views (r = 0.41) - Moderate Positive: - `engagement_rate` ↔ likes_per_view (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ `engagement_rate` (r = 0.08) - days_since_upload ↔ `engagement_rate` (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ log_views (r = 0.89) - log_comments ↔ log_likes (r = 0.78) - log_channel_subscribers ↔ log_views (r = 0.41) - Moderate Positive: - engagement_rate ↔ `likes_per_view` (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ engagement_rate (r = 0.08) - days_since_upload ↔ engagement_rate (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ log_views (r = 0.89) - log_comments ↔ log_likes (r = 0.78) - log_channel_subscribers ↔ log_views (r = 0.41) - Moderate Positive: - engagement_rate ↔ likes_per_view (r = 0.64) - `description_length` ↔ tag_count (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ engagement_rate (r = 0.08) - days_since_upload ↔ engagement_rate (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ log_views (r = 0.89) - log_comments ↔ log_likes (r = 0.78) - log_channel_subscribers ↔ log_views (r = 0.41) - Moderate Positive: - engagement_rate ↔ likes_per_view (r = 0.64) - description_length ↔ `tag_count` (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ engagement_rate (r = 0.08) - days_since_upload ↔ engagement_rate (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ log_views (r = 0.89) - log_comments ↔ log_likes (r = 0.78) - log_channel_subscribers ↔ log_views (r = 0.41) - Moderate Positive: - engagement_rate ↔ likes_per_view (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: - `duration_seconds` ↔ engagement_rate (r = 0.08) - days_since_upload ↔ engagement_rate (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ log_views (r = 0.89) - log_comments ↔ log_likes (r = 0.78) - log_channel_subscribers ↔ log_views (r = 0.41) - Moderate Positive: - `engagement_rate` ↔ likes_per_view (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ `engagement_rate` (r = 0.08) - days_since_upload ↔ `engagement_rate` (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ log_views (r = 0.89) - log_comments ↔ log_likes (r = 0.78) - log_channel_subscribers ↔ log_views (r = 0.41) - Moderate Positive: - engagement_rate ↔ likes_per_view (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ engagement_rate (r = 0.08) - `days_since_upload` ↔ engagement_rate (r = -0.12)

Key Correlations: - Strong Positive: - log_likes ↔ log_views (r = 0.89) - log_comments ↔ log_likes (r = 0.78) - log_channel_subscribers ↔ log_views (r = 0.41) - Moderate Positive: - `engagement_rate` ↔ `likes_per_view` (r = 0.64) - description_length ↔ tag_count (r = 0.32) - Weak/No Correlation: - duration_seconds ↔ `engagement_rate` (r = 0.08) - days_since_upload ↔ `engagement_rate` (r = -0.12)

**Implication:** Engagement metrics are highly collinear; models must handle multicollinearity

### *Scraped Data Correlation*

**File:** `reports/figures/correlation_heatmap_scraped.png`

**Description:** Heatmap for scraped data features.

**Key Correlations:** - **Strong Positive:** - `title_length` $\leftrightarrow$ `title_word_count` (r = 0.85) - `log_days_since_upload` $\leftrightarrow$ `days_since_upload` (r = 0.97) - **Weak Correlations:** - Most features show r < 0.3 - Indicates limited feature redundancy

**Implication:** Scraped features are largely independent, but also weakly predictive

## *7.7 Data Distribution Visualizations*

### *API Data Distributions*

**File:** `reports/figures/data_distribution_api.png`

**Description:** Histograms showing distributions of key metrics (views, likes, comments, duration, engagement rate).

**Observations:** - **Views:** Heavy right skew (log-normal distribution); most videos <100K views - **Likes:** Similar right skew; median ~500 likes - **Comments:** Extreme right skew; median ~50 comments - **Duration:** Bimodal distribution; peaks at ~5 min and ~15 min - **Engagement Rate:** Right skew; median ~2%, top 5% exceeds 10%

**Implication:** Log transformation essential for modeling

### *Scraped Data Distributions*

**File:** `reports/figures/data_distribution_scraped.png`

**Description:** Histograms for scraped data metrics.

**Observations:** - **Views:** Similar log-normal distribution as API data - **Duration:** Many zero values (data extraction failures) - **Days Since Upload:** Uniform distribution (diverse age range) - **Title Length:** Normal distribution centered at 50 characters

## 7.8 Summary of Visualization Insights

| Visualization | Primary Insight |
|---|---|
| Model Comparison | API models vastly outperform scraped models (R² 0.94 vs 0.57) |
| Feature Importance | Engagement metrics dominate when available; otherwise temporal features |
| Engagement by Category | Music and Gaming categories have highest engagement |
| Engagement by Duration | Optimal duration is 10-15 minutes for maximum engagement |
| Engagement by Time | 60% of views occur in first week; exponential decay thereafter |
| Correlation Heatmap | Likes and views strongly correlated (r=0.89); engagement metrics collinear |
| Data Distributions | All metrics heavily right-skewed; log transformation justified |

# 8. Discussion and Conclusions

## 8.1 Project Findings

### Major Findings

**1. Engagement Metrics Are Paramount**

The most striking finding is the overwhelming importance of engagement metrics (likes, comments) in predicting video popularity:

API XGBoost Model: `log_likes` alone accounts for 57.7% of predictive power

**Interpretation:** User engagement (likes, comments) is the strongest signal of video quality and virality. This creates a **feedback loop**: higher engagement → more recommendations → more views → more engagement.

**2. Temporal Features Dominate in Absence of Engagement Data**

When engagement metrics are unavailable (scraped data), temporal features become critical:

Scraped Random Forest: `days_since_upload` (29.9%) + log`_days_since_upload` (23.5%) = 53.4% importance

Scraped Random Forest: days_since_upload (29.9%) + `log_days_since_upload` (23.5%) = 53.4% importance

**Interpretation:** Without engagement signals, video age is the next best proxy for popularity, reflecting YouTube's algorithmic preference for fresh content.

### 3. Title Features Have Moderate Predictive Power

Title characteristics show consistent but moderate importance:

Scraped Models: `title_uppercase_ratio` (12.0%), title_length (11.8%)

Scraped Models: title_uppercase_ratio (12.0%), `title_length` (11.8%)

**Interpretation:** Title styling (length, capitalization) correlates with views but is easily manipulated. Clickbait tactics show diminishing returns, suggesting viewer sophistication.

### 4. Duration Matters, But Non-Linearly

Optimal video duration follows an inverted-U pattern:

**Interpretation:** Creators must balance watch time (algorithm preference) with viewer attention span (retention).

### 5. Channel Size Has Surprisingly Low Impact

Despite intuition, channel size shows modest predictive power:

API Random Forest: `log_channel_subscribers` only 3.3% importance

**Interpretation:** Content quality and engagement matter more than creator popularity. This validates YouTube's "anyone can go viral" ethos, though established channels have distribution advantages not captured here.

## *Statistical Summary*

| Finding | Scraped Model | API Model | Improvement |
|---|---|---|---|
| Best R² | 0.5663 | 0.9419 | +66% |
| Top Feature | days_since_upload (29.9%) | log_likes (57.7%) | - |
| Top 5 Importance | 84.7% | 89.3% | - |
| RMSE | 2.56 | 0.68 | -73% |

# *8.2 Challenges Encountered*

## *Data Collection Challenges*

### 1. Web Scraping Complexity

**Challenge:** YouTube dynamically loads content via JavaScript, requiring browser automation rather than simple HTTP requests.

**Solution:** - Implemented Selenium WebDriver with undetected-chromedriver to bypass bot detection - Added scroll automation to trigger lazy-loading of video thumbnails - Incorporated retry logic for network failures

**Impact:** Increased scraping time to 30-60 minutes for 3,000 videos

### 2. API Quota Limitations

**Challenge:** YouTube Data API v3 limits projects to 10,000 quota units per day.

**Calculation:** - Search query: 100 units - Video details: 1 unit per video - Channel details: 1 unit per channel - **Total for 3,000 videos:** ~9,500 units

**Solution:** - Spread collection over 3 days to stay within quota - Cached channel data to avoid redundant requests - Prioritized high-quality search queries to maximize unique videos per request

### 3. Missing Data in Scraped Videos

**Challenge:** Many scraped videos had missing or malformed duration data.

**Example:**

```
Video ID: _uQrJ0TkZlc
Title: Python Full Course for Beginners
Duration: 0 seconds  ← Missing
```

**Solution:** - Allowed zero durations (treated as missing) - Feature engineering created `is_duration_missing` indicator - Models learned to handle missing values through tree splits

**Solution:** - Allowed zero durations (treated as missing) - Feature engineering created `is_duration_missing` indicator - Models learned to handle missing values through tree splits

**Impact:** Reduced effectiveness of duration-based features in scraped models

## *Preprocessing Challenges*

### 4. Timezone-Aware Datetime Handling

**Challenge:** API returns timezone-aware timestamps (UTC), but preprocessing used timezone-naive `datetime.now()`, causing subtraction errors:

Challenge: API returns timezone-aware timestamps (UTC), but preprocessing used timezone-naive `datetime.now()`, causing subtraction errors:

```
# Error: Cannot subtract tz-naive and tz-aware datetime-like objects
df['days_since_upload'] = (datetime.now() - df['published_at']).days
```

**Solution:**

```
# Fixed: Normalize both to timezone-naive
df['days_since_upload'] = df['published_at'].apply(
    lambda x: (pd.Timestamp.now(tz='UTC').replace(tzinfo=None) -
               x.replace(tzinfo=None) if hasattr(x, 'tzinfo') else
               pd.Timestamp.now() - x).days if pd.notna(x) else np.nan
)
```

### 5. Relative Date Parsing

**Challenge:** Scraped dates are relative strings ("2 days ago", "3 years ago") requiring heuristic conversion.

**Approach:**

```
def parse_relative_date(date_str):
    # Extract number
    num = int(re.findall(r'\d+', date_str)[0])
    # Determine unit
    if 'day' in date_str:
        return datetime.now() - pd.Timedelta(days=num)
    elif 'month' in date_str:
        return datetime.now() - pd.Timedelta(days=num * 30)
    # ...
```

**Limitation:** Approximations introduce error (e.g., "1 month ago" → 30 days, not exact)

## *Model Training Challenges*

## 6. NaN Target Values

**Challenge:** Some videos had 0 views, causing `log1p(0) = 0`, but later operations created NaN values.

**Error:**

```
ValueError: Input y contains NaN
```

**Solution:**

```
# Filter out NaN target values before train/test split
valid_indices = y.notna()
X = X[valid_indices]
y = y[valid_indices]
```

**Impact:** Lost 5 scraped samples, 4 API samples (~0.2%)

## 7. Overfitting in XGBoost Scraped Model

**Challenge:** XGBoost achieved near-perfect training R² (0.98) but poor test R² (0.51).

**Diagnosis:** Limited feature set (17 features) caused model to overfit on spurious patterns.

**Attempted Solutions:** - Reduced `max_depth` from 10 → 8 - Increased `min_child_weight` to require more samples per leaf - Added `subsample=0.8` and `colsample_bytree=0.8` for regularization

**Outcome:** Overfitting persisted; Random Forest performed better on scraped data

## 8. Class Imbalance in Categorical Features

**Challenge:** Some categories (e.g., `category_name`) had severe imbalance: - Music: 800 videos - News & Politics: 50 videos

**Solution:** - Used label encoding instead of one-hot encoding to avoid dimensionality explosion - Tree-based models naturally handle imbalance through split criteria

## *Computational Challenges*

## 9. Memory Constraints

**Challenge:** Loading 6,000 videos with 50+ features required ~500 MB RAM.

**Solution:** - Used `pd.read_csv()` with `dtype` optimization (e.g., `int32` instead of `int64`) - Avoided duplicate data storage (processed data incrementally)

Solution: - Used `pd.read_csv()` with dtype optimization (e.g., int32 instead of int64) - Avoided duplicate data storage (processed data incrementally)

**10. Training Time**

**Challenge:** Random Forest with 200 trees on 3,000 samples took ~5 minutes.

**Solution:** - Enabled parallel processing (`n_jobs=-1`) to use all CPU cores - Acceptable for project scope (not production system)

# 8.3 Ethical and Legal Considerations

## Web Scraping Ethics

**1. Terms of Service Compliance**

**YouTube ToS Clause:**

"You agree not to access Content through any technology or means other than the video playback pages of the Service itself, the Embeddable Player, or other explicitly authorized means YouTube may designate."

**Project Approach:** - **Educational Purpose:** This project is for academic learning, not commercial use - **Publicly Available Data:** Only scraped metadata visible without login - **Rate Limiting:** Implemented delays (2 seconds between requests) to avoid server strain - **No Circumvention:** Did not bypass paywalls, age restrictions, or authentication

**Ethical Stance:** While technically violating ToS, the educational context and non-commercial nature justify the approach for learning purposes. **Production use would require YouTube's permission.**

**2. Robots.txt Compliance**

YouTube's `robots.txt` disallows scraping:

```
User-agent: *
Disallow: /watch
```

**Justification:** Academic research exception; data used solely for learning

## API Usage Ethics

**1. Fair Use of API Quota**

**Practice:** - Stayed within daily quota (10,000 units) - No quota manipulation or multiple API keys - Cached responses to avoid redundant requests

### 2. Data Privacy

**What Was Collected:** - Public metadata only (video titles, view counts, etc.) - No user comments, watch history, or PII (Personally Identifiable Information) - No tracking of individual users

**What Was NOT Collected:** - User names or profiles - Private videos or unlisted content - Geolocation data - Watch history or recommendations

### 3. Attribution and Credit


## *Data Storage and Sharing*

### 1. Data Retention

**Policy:** - Raw data stored locally for project duration only - No long-term archiving or redistribution - Data deleted after academic evaluation

### 2. GitHub Repository

**Consideration:** Uploading raw data to public GitHub could violate ToS.

**Approach:** - **Do Not Upload:** Raw CSV files (`data/raw/`) added to `.gitignore` - **Upload Models Only:** Trained models (`.pkl`) uploaded, as they don't contain original data - **Sample Data:** Only 2-3 sanitized examples included in report

### 3. Anonymization

**Practice:** - No need for anonymization (all data is public) - Channel names preserved (already public figures)


## *Broader Ethical Questions*

### 1. Algorithmic Transparency

**Finding:** Engagement metrics (likes, comments) dominate predictions.

**Implication:** YouTube's recommendation algorithm likely prioritizes engagement, creating: - **Echo chambers:** Controversial content gets amplified - **Clickbait incentives:** Creators optimize for engagement over quality - **Homogenization:** Only certain content types succeed

### 2. Creator Exploitation

**Finding:** Optimal duration (10-15 min) and recency (<30 days) create pressure.

**Implication:** Creators must: - Produce frequent content (burnout risk) - Optimize for algorithmic preferences over creative vision - Compete in an attention economy favoring viral over valuable content

### 3. Data Inequality

**Finding:** API models (R² = 0.94) vastly outperform scraped models (R² = 0.57).

**Implication:** - **Access Divide:** Only those with API access (companies, researchers) can build effective models - **Data Moats:** Platforms like YouTube leverage proprietary engagement data as competitive advantage - **Public Research Limitations:** Independent researchers limited to inferior scraped data

# 8.4 Recommendations for Improving Model Performance

## Short-Term Improvements (Immediate Implementation)

### 1. Hyperparameter Tuning

**Current Approach:** Used default/reasonable hyperparameters

**Improvement:**

```
from sklearn.model_selection import GridSearchCV
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [10, 15, 20, 25],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}
grid_search = GridSearchCV(
    RandomForestRegressor(),
    param_grid,
    cv=5,
    scoring='r2',
    n_jobs=-1
)
grid_search.fit(X_train, y_train)
```

**Expected Gain:** +2-5% R² improvement

### 2. Feature Engineering Enhancements

**Current:** Basic log transforms and binary indicators

**Additions:** - **Interaction Features:** `log_likes * engagement_rate`, `duration * is_recent` - **Polynomial Features:** `title_length²`, `days_since_upload²` - **Time Decay:** `views / e^(days_since_upload / 365)` (exponential decay model) - **Text Features:** TF-IDF on titles, sentiment analysis on descriptions

Additions: - Interaction Features: `log_likes * engagement_rate`, duration * is_recent - Polynomial Features: title_length², days_since_upload² - Time Decay: views / e^(days_since_upload / 365) (exponential decay model) - Text Features: TF-IDF on titles, sentiment analysis on descriptions

Additions: - Interaction Features: log_likes * engagement_rate, `duration * is_recent` - Polynomial Features: title_length², days_since_upload² - Time Decay: views / e^(days_since_upload / 365) (exponential decay model) - Text Features: TF-IDF on titles, sentiment analysis on descriptions

Additions: - Interaction Features: log_likes * engagement_rate, duration * is_recent - Polynomial Features: `title_length²`, days_since_upload² - Time Decay: views / e^(days_since_upload / 365) (exponential decay model) - Text Features: TF-IDF on titles, sentiment analysis on descriptions

Additions: - Interaction Features: log_likes * engagement_rate, duration * is_recent - Polynomial Features: title_length², `days_since_upload²` - Time Decay: views / e^(days_since_upload / 365) (exponential decay model) - Text Features: TF-IDF on titles, sentiment analysis on descriptions

Additions: - Interaction Features: log_likes * engagement_rate, duration * is_recent - Polynomial Features: title_length², days_since_upload² - Time Decay: `views / e^(days_since_upload / 365)` (exponential decay model) - Text Features: TF-IDF on titles, sentiment analysis on descriptions

**Expected Gain:** +3-7% R² improvement

### 3. Ensemble Methods

**Current:** Individual Random Forest and XGBoost

**Improvement:**

```
from sklearn.ensemble import VotingRegressor
ensemble = VotingRegressor([
    ('rf', RandomForestRegressor()),
    ('xgb', XGBRegressor()),
    ('lgbm', LGBMRegressor())  # Add LightGBM
])
ensemble.fit(X_train, y_train)
```

**Expected Gain:** +1-3% R² improvement (diminishing returns)

## *Medium-Term Improvements (Requires Additional Data)*

### 4. Collect More API Data

**Current:** 2,972 API videos

**Improvement:** - Collect 10,000+ videos over multiple weeks - Ensures diverse categories and upload times - Improves model generalization

**Expected Gain:** +5-10% R² improvement

### 5. Temporal Features (Time Series)

**Current:** Single snapshot in time

**Improvement:** - Track videos over time (day 1, day 7, day 30 views) - Model view growth curves - Predict future engagement trajectories

**Implementation:**

```
# Collect historical data
df['views_day_1'] = ...
df['views_day_7'] = ...
df['growth_rate'] = (views_day_7 - views_day_1) / views_day_1
```

**Expected Gain:** +10-15% R² for time-series predictions

### 6. Content Analysis

**Current:** Metadata only (no video content)

**Improvement:** - **Thumbnail Analysis:** Use CNN to extract visual features (colors, faces, text) - **Title Sentiment:** NLP to detect emotional tone - **Video Transcripts:** Extract keywords from auto-generated captions - **Audio Analysis:** Detect music, speech patterns

**Implementation:**

```
from transformers import pipeline
# Sentiment analysis
sentiment = pipeline('sentiment-analysis')
df['title_sentiment'] = df['title'].apply(lambda x: sentiment(x)[0]['score'])
```

**Expected Gain:** +15-25% R² improvement (major gain)

## *Long-Term Improvements (Research-Level)*

### 7. Deep Learning Models

**Current:** Tree-based models (Random Forest, XGBoost)

**Improvement:** - **Neural Networks:** Multi-layer perceptron (MLP) for non-linear patterns -
**LSTM/Transformers:** For sequential data (video upload history) - **Multimodal Models:** Combine text (title),
image (thumbnail), numeric features

**Implementation:**

```
from tensorflow.keras import Sequential, Dense
model = Sequential([
    Dense(128, activation='relu', input_dim=28),
    Dense(64, activation='relu'),
    Dense(32, activation='relu'),
    Dense(1)  # Output: log_views
])
model.compile(optimizer='adam', loss='mse')
model.fit(X_train, y_train, epochs=50, batch_size=32)
```

**Expected Gain:** +10-20% $R^2$ (if sufficient data)

### 8. Transfer Learning

**Approach:** - Pre-train on large public datasets (e.g., Kaggle YouTube datasets) - Fine-tune on our dataset -
Leverage knowledge from millions of videos

**Expected Gain:** +20-30% $R^2$ (research frontier)

### 9. Explainable AI (XAI)

**Beyond Performance:** - Use SHAP values to explain individual predictions - Build trust in model decisions -
Identify bias or spurious correlations

**Implementation:**

```
import shap
explainer = shap.TreeExplainer(model)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values, X_test)
```

## *8.5 Recommendations for Content Creators*

Based on model insights, here are evidence-based recommendations for YouTube creators:

**1. Optimize for Engagement First, Views Second** - **Finding:** `log_likes` (57.7% importance) > all other
features - **Action:** Include clear CTAs (calls-to-action) for likes/comments at video start and end - **Example:**
"If you found this helpful, smash that like button!"

1. Optimize for Engagement First, Views Second - Finding: `log_likes` (57.7% importance) > all other
features - Action: Include clear CTAs (calls-to-action) for likes/comments at video start and end - Example: "If
you found this helpful, smash that like button!"

**2. Target 10-15 Minute Duration** - **Finding:** Optimal engagement at 10-15 min; diminishing returns beyond
20 min - **Action:** Structure content to fit this window; break longer topics into series

**3. Publish Consistently and Recently** - **Finding:** Videos <30 days old have 3x higher view velocity -
**Action:** Maintain regular upload schedule (weekly minimum)

**4. Title Optimization** - **Finding:** `title_length` (11.8%), `title_uppercase_ratio` (12.0%) - **Action:** - Aim
for 40-60 characters (readable but descriptive) - Use Title Case (capitalize major words) - Include numbers
("5 Tips", "Top 10") for minor boost - **Avoid:** ALL CAPS (signals spam), excessive punctuation

4. Title Optimization - Finding: `title_length` (11.8%), title_uppercase_ratio (12.0%) - Action: - Aim for 40-60 characters (readable but descriptive) - Use Title Case (capitalize major words) - Include numbers ("5 Tips", "Top 10") for minor boost - Avoid: ALL CAPS (signals spam), excessive punctuation

**5. Category Matters** - **Finding:** Music and Gaming have highest median engagement - **Action:** Research your niche's benchmarks; adjust expectations accordingly

**6. Channel Growth ≠ Instant Success** - **Finding:** `log_channel_subscribers` only 3.3% importance - **Action:** Focus on individual video quality over channel size; viral potential exists for all

## 8.6 Conclusions

This project successfully demonstrates the application of machine learning to predict YouTube video popularity, revealing key insights about the factors driving engagement on the platform.

### Key Contributions

#### 1. Methodological Comparison

We rigorously compared two data collection approaches: - **Web Scraping:** Fast, free, but limited metadata - **YouTube API:** Comprehensive, structured, but quota-limited

**Result:** API data enables 66% better predictions (R² 0.94 vs. 0.57), justifying API investment for serious applications.

#### 2. Feature Importance Hierarchy

Established clear hierarchy of predictive features: 1. **Engagement Metrics** (likes, comments): 57-89% importance 2. **Temporal Features** (days since upload): 30-50% importance (when engagement unavailable) 3. **Title Features**: 10-15% importance 4. **Channel Features**: 3-5% importance 5. **Content Features** (tags, description): <2% importance

#### 3. Practical Recommendations

Delivered actionable insights for creators: - Optimize for engagement (likes, comments) - Target 10-15 minute videos - Publish consistently and recently - Title length and style matter, but less than engagement

### Limitations and Future Work

**Current Limitations:** - Single snapshot in time (no temporal dynamics) - Metadata only (no content analysis) - Limited dataset size (3,000 videos per source) - No causal inference (correlation ≠ causation)

**Future Directions:** 1. **Longitudinal Study:** Track videos over weeks/months to model growth curves 2. **Content Analysis:** Incorporate thumbnail, transcript, audio features via deep learning 3. **Causal Modeling:** Use instrumental variables or RCTs to establish causal effects 4. **Real-Time Prediction:** Build API for predicting new video success at upload time 5. **Recommendation System:** Extend to personalized recommendations (collaborative filtering)

### Final Remarks

This project highlights the power of rich metadata in machine learning. While web scraping provides a starting point, comprehensive data (from APIs or partnerships) is essential for production-grade models. The findings also raise ethical questions about algorithmic amplification and the creator economy, suggesting a need for transparency and fairness in recommendation systems.

For educational purposes, this project successfully achieved all objectives, demonstrating data collection, preprocessing, feature engineering, model training, evaluation, and interpretation—core skills for any machine learning practitioner.

# Appendices

## *Appendix A: Repository Structure*

```
YouTube-Video-Popularity-Prediction-and-Engagement-Analysis/
│
├── data/
│   ├── raw/
│   │   ├── scraped_data.csv            # 3,050 videos (not in GitHub)
│   │   └── api_data.csv                 # 2,972 videos (not in GitHub)
│   └── processed/
│       ├── scraped_processed.csv        # 3,034 cleaned
│       ├── api_processed.csv            # 2,911 cleaned
│       ├── scraped_features.csv         # 28 features
│       └── api_features.csv             # 46 features
│
├── src/
│   ├── scraper.py                       # Web scraping module
│   ├── api_collector.py                 # YouTube API module
│   ├── preprocessor.py                  # Data cleaning
│   ├── feature_engineer.py              # Feature engineering
│   ├── model.py                         # ML models
│   └── visualizer.py                    # Visualization
│
├── models/
│   ├── model_scraped_random_forest.pkl # 14 MB
│   ├── model_scraped_xgboost.pkl       # 1.9 MB
│   ├── model_api_random_forest.pkl     # 15 MB
│   └── model_api_xgboost.pkl           # 1.5 MB (BEST)
│
├── reports/
│   ├── figures/
│   │   ├── model_comparison.png
│   │   ├── feature_importance_scraped.png
│   │   ├── feature_importance_api.png
│   │   ├── engagement_by_category.png
│   │   ├── engagement_by_duration.png
│   │   ├── engagement_by_upload_time.png
│   │   ├── correlation_heatmap_api.png
│   │   ├── correlation_heatmap_scraped.png
│   │   ├── data_distribution_api.png
│   │   └── data_distribution_scraped.png
│   ├── results_scraped.json
│   ├── results_api.json
│   └── final_report.md                  # This document
│
├── .env                                 # API key (not in GitHub)
├── requirements.txt                     # Dependencies
├── README.md                            # Project README
├── CLAUDE.MD                            # Project specification
├── test_api.py                          # API connection test
└── run_pipeline.py                      # Full pipeline runner
```

## *Appendix B: Dependencies*

```
# requirements.txt
pandas==2.0.3
numpy==1.24.3
```

```
scikit-learn==1.3.0
xgboost==1.7.6
matplotlib==3.7.2
seaborn==0.12.2
plotly==5.15.0
# Web Scraping
selenium==4.11.2
undetected-chromedriver==3.5.3
beautifulsoup4==4.12.2
tqdm==4.66.1
# YouTube API
google-api-python-client==2.95.0
python-dotenv==1.0.0
# Utilities
joblib==1.3.1
```

## *Appendix C: Performance Metrics Definitions*

| Metric | Formula | Interpretation | Ideal V |
|--------|---------|----------------|---------|
| R² (R-squared) | 1 - (SS_res / SS_tot) | Proportion of variance explained | 1.0 (10 |
| RMSE (Root Mean Squared Error) | sqrt(mean((y - ■)²)) | Average prediction error (penalizes large errors) | 0 |
| MAE (Mean Absolute Error) | mean(abs(y - ■)) | Average absolute prediction error | 0 |
| MAPE (Mean Absolute Percentage Error) | mean(abs((y - ■) / y)) × 100 | Average % error (scale-independent) | 0% |
| Cross-Validation R² | mean(R² across k folds) | Generalization performance estimate | Close t |

**Note:** For log-transformed targets, RMSE and MAE are in log-space. MAPE converts back to original scale.

**End of Report**