BIRZEIT UNIVERSITY

Birzeit University
Faculty of Engineering & Technology
Department of Electrical & Computer Engineering

# Prediction Trending Videos on YouTube

Prepared by

| Mohammad Shreteh | 1201369 |
| Ali Halayqa | 1201769 |
| Saleh Zhour | 1201941 |

Supervised By:
Dr. Mohammed Hussain

Section: 10

Graduation Project submitted to the Department of Electrical and
Computer Engineering in partial fulfillment of the requirements for the
degree of B.Sc. in Computer Engineering

Birzeit
July- 2025

# **Abstract**

This graduation project presents a machine learning-based framework for predicting trending YouTube videos using a combination of metadata, user engagement metrics, and temporal features. The main goal is to build an accurate and interpretable model capable of identifying videos with high viral potential before they appear in the trending section. To achieve this, a comprehensive dataset is collected and preprocessed, including features such as views, likes, comments, upload time, video titles, and more.

Multiple classification algorithms are implemented and evaluated, including Logistic Regression, K-Nearest Neighbors (KNN), Decision Tree, Random Forest, XGBoost, and LightGBM. Among these, the Random Forest model achieves the highest accuracy and consistency, outperforming others in terms of F1-score and ROC-AUC metrics. Additionally, visual analysis and feature correlation studies are conducted to gain deeper insights into the factors influencing video trends.

To demonstrate the practical usability of the system, a web-based interface is developed, allowing users to input video data and receive real-time trending predictions. The system's effectiveness is validated through case-based testing, confirming its capability to distinguish between trending, non-trending, and outdated videos.

The findings highlight the importance of combining content features with early engagement signals for accurate trend prediction. This project not only contributes to understanding the dynamics of viral content but also provides a foundation for future enhancements using deep learning, real-time analysis, and multimodal data integration.

# المستخلص

يقدم مشروع التخرج هذا إطار عمل قائم على التعلم الآلي للتنبؤ بمقاطع فيديو يوتيوب الرائجة باستخدام مزيج من البيانات الوصفية، ومقاييس تفاعل المستخدم، والخصائص الزمنية. الهدف الرئيسي هو بناء نموذج دقيق وقابل للتفسير قادر على تحديد مقاطع الفيديو ذات الإمكانات الفيروسية العالية قبل ظهورها في قسم "المقاطع الرائجة". لتحقيق ذلك، يتم جمع مجموعة بيانات شاملة ومعالجتها مسبقًا، بما في ذلك خصائص مثل المشاهدات، والإعجابات، والتعليقات، ووقت التحميل، وعناوين الفيديو، وغيرها.

يتم تنفيذ وتقييم خوارزميات تصنيف متعددة، بما في ذلك الانحدار اللوجستي، وأقرب جيران (KNN)، وشجرة القرار، والغابة العشوائية، وXGBoost، وLightGBM. من بين هذه الخوارزميات، يحقق نموذج الغابة العشوائية أعلى دقة واتساق، متفوقًا على غيره من حيث درجة F1 ومقاييس ROC-AUC. بالإضافة إلى ذلك، يتم إجراء تحليل بصري ودراسات ارتباط الخصائص لاكتساب رؤى أعمق حول العوامل المؤثرة على اتجاهات الفيديو. لإثبات قابلية استخدام النظام عمليًا، طُوّرت واجهة ويب تُمكّن المستخدمين من إدخال بيانات الفيديو وتلقي تنبؤات آنية بالاتجاهات. ويتم التحقق من فعالية النظام من خلال اختبارات قائمة على الحالات، مما يؤكد قدرته على التمييز بين الفيديوهات الرائجة وغير الرائجة والقديمة.

تُبرز النتائج أهمية دمج ميزات المحتوى مع إشارات التفاعل المبكرة لضمان دقة التنبؤ بالاتجاهات. لا يُسهم هذا المشروع في فهم ديناميكيات المحتوى الفيروسي فحسب، بل يُرسي أيضًا أساسًا للتحسينات المستقبلية باستخدام التعلم العميق والتحليل الآني وتكامل البيانات متعدد الوسائط.

# Table of Contents

# List of Figures

# List of Tables

# List of Abbreviations

**RNN –** Recurrent Neural Network

**LSTM –** Long Short-Term Memory

**GRU –** Gated Recurrent Unit

**CML –** Classical Machine Learning

**AML –** Advanced Machine Learning

**CNN –** Convolutional Neural Network

**KNN –** K-Nearest Neighbors

**SVM –** Support Vector Machine

**TF-IDF –** Term Frequency-Inverse Document Frequency

**PCA –** Principal Component Analysis

**API –** Application Programming Interface

**ROC –** Receiver Operating Characteristic

**AUC –** Area Under the Curve

**RMSE –** Root Mean Squared Error

**EDA –** Exploratory Data Analysis

**XGBoost –** Extreme Gradient Boosting

**LightGBM –** Light Gradient Boosting Machine

**SMOTE –** Synthetic Minority Over-sampling Technique

**BERT –** Bidirectional Encoder Representations from Transformers

**XAI –** Explainable Artificial Intelligence

**GUI –** Graphical User Interface

**NLP –** Natural Language Processing

# Chapter 1

# Introduction

## Contents

## 1.1   Introduction

Each platform of videos sharing is focusing its strategy to popularize the video content. Our project is focused on YouTube as the video content sharing network. Additionally, trending is a platform-defined metric to popularize the popular content maintains a different purist to keep track of trending content. Additionally, YouTube algorithm helps to display the trending content on the main page. In addition, when users log in one of the trending videos is on the main screen. As shown in Figure 1.1, YouTube's homepage visually prioritizes trending content, making it more accessible to users (Zeru, 2024).

YouTube attracting more than 4 billion views daily from a global audience, making it an important hub for user-generated content [1]. There are several metrics that help measure a video's popularity, such as the number of views, although some videos gain a lot of attention quickly, gaining a place in YouTube's Popular tab and potentially indicating future popularity [2]. Classification is one of the widespread problems in the world Machine-learning context. In addition, for binary classification problems, algorithms like Random Forests and Logistic Regression are well-known and popular due to the fact that they achieve a balance between performance and interpretability [3]. Moreover, both algorithms possess their own merits and demerits. Therefore, we have now begun implementing these models, along with others like XGBoost and Decision Trees, to actual YouTube datasets in order to check their performance in detecting trending videos.

Figure 1.1: Trending page on YouTube

## 1.2 Motivation

YouTube's 2021 revenue of $28.8 billion represented a 46% year-over-year growth [4]. YouTube is a very alluring platform for content creators looking for steady income, which mostly focuses on increasing user awareness, because of its significant revenue. The significance of the popularity lifecycle is further highlighted by the direct impact it has on earnings and spending.

YouTube, the second most popular website in the world behind Google, has a big influence on society. Trending videos illustrate the wider ramifications of popularity on social media by influencing viewers' attitudes, behaviors, and viewpoints.

According prior research on popularity prediction primarily examines engagement metrics and social media content classification to estimate the amount of engagement a specific piece of content will receive at a given future time [5][6]. Most engagement prediction frameworks rely on conventional techniques, while deep neural networks are commonly employed to create predictions for text-based media, such as tweets [7].

The goal of this project is to improve video popularity prediction by utilizing sequence modeling and neural networks. By using video information and textual content to generate predictions, it fills in current gaps and offers a fresh solution to this problem. And this justifies the necessity to develop a predictive model, which we have begun testing with actual data.

## 1.3    Contribution

In this section, we examine the issues and constraints and offer initial solutions that are regarded as contributions to the field of YouTube video classification. The following is a description of our contributions:

- Content and Engagement-Based Classification: The topic of finding trending videos is approached by our system as a multi-modal classification problem, rather than depending exclusively on conventional popularity measurements like views, likes, and comments. This method combines engagement dynamics (like-to-view ratio, sentiment of comments) with content semantics (e.g., video titles, descriptions, tags) [8][9].The system eliminates the need for computationally costly techniques like graph-based analysis by processing and classifying films in almost real-time [10].

- Early Identification of Trending Potential: By studying the semantic content and early engagement patterns, the proposed methodology predicts trending videos before they gain widespread popularity, hence identifying content with viral potential early on [11]. They demonstrate the relationship between early content signals and their subsequent popularity. The methodology provides useful information for platforms seeking to promote new content and for producer's attempt.

- Large-scale experimental setups: We will carry out in-depth tests with different feature, machine learning model, and hyper-parameter setups to guarantee peak performance. This involves testing state-of-the-art models like transformers (BERT) [12], ensemble techniques like XGBoost [13], and content-based and engagement-based features. Using best practices from related studies, the configuration that provides the best classification performance will be selected for the final framework [14].

These contributions address important issues in the field of classifying popular videos by presenting a fresh method that uses collaborative data and content-based analysis to anticipate trending videos on YouTube and other platforms in real time with high accuracy.

## 1.4    Methodology

The goal of this project is to classify YouTube videos as either trending or non-trending while understanding the factors influencing trends. The aim is to achieve high accuracy and interpretability in the classification model. The first step involves investigating the data sources, including publicly available APIs like the YouTube Data API, to access relevant video attributes. Video metadata, channel information, and engagement statistics are then extracted to be used for classification. Data extraction focuses on gathering key features, such as video attributes (title, description, upload time, and duration), engagement statistics (views, likes, dislikes, and comments), and channel data (subscriber count, number of videos, and historical performance). Once the data is extracted, it undergoes preprocessing, which involves cleaning the data by handling missing values and removing irrelevant entries. Categorical variables, such as description and title, are encoded, and numerical features like views, likes, and dislikes are normalized. Exploratory data analysis (EDA)

follows, where the distribution and relationships among features are analyzed. This helps to identify patterns in the data that can distinguish trending videos from non-trending ones. Correlations are visualized using techniques such as boxplots or heat maps. To reduce the complexity of the dataset, dimensionality reduction techniques like Principal Component Analysis (PCA) or Independent Component Analysis (ICA) are applied. The goal is to retain the components that explain the maximum variance while simplifying the dataset. Next, various classification models are compared to determine the best fit for the dataset. These models include logistic regression, random forest, gradient boosting, and neural networks. Once the model is selected, it is trained on the processed dataset using supervised learning techniques to classify the videos as trending or non-trending. The model's performance is evaluated using key metrics such as accuracy, precision, recall, F1 score, and ROC-AUC. K-fold cross-validation is performed to ensure the model's robustness. To further improve performance, feature selection and engineering are optimized to boost both accuracy and computational efficiency. Hyperparameters of the selected models are tuned using methods like Grid Search or Random Search. Finally, the model is developed to provide real-time classification of videos, with a user-friendly interface to display classification results and insights. And the implementation process has been initiated based on this methodology, and initial testing results are being collected and analyzed. The entire pipeline of data extraction, modeling, and deployment is summarized in Figure 1.2, which outlines our proposed system's architecture.



Figure 1.2: The proposed method's high-level design

# Chapter 2

# Background and Related Work

## Contents

# 2.1 Background

## 2.1.1 Overview of YouTube as a Platform

This section provides an overview of YouTube's evolution and growth, highlighting its global reach, user engagement, and significance as a platform for content creators, marketers, and advertisers.

### 2.1.1.1 YouTube's Evaluation and Growth

Since its establishment in 2005, YouTube has grown to become the largest video-sharing platform in the world, with billions of active users around the globe. Gaining momentum over the years, YouTube was acquired by Google in 2006. Over time, it has greatly influenced the fields of entertainment, education, and marketing by opening up avenues for content creators, brands, and influencers to reach out to the global audience. It reaches 4 billion views a day and is viewed by many cultures in different parts of the world by this year [15]. This exponential growth trend is visualized in Figure 2.1, which shows projections of user engagement through 2028.



Figure 2.1: YouTube Active Users Growth Forecast (2019-2028)

### 2.1.1.2 Statistics and User Engagement

YouTube's user base spans over 100 countries with more than 80 languages. It has more than 2 billion logged-in monthly users, which make contribute to the huge daily views that YouTube sees. The average viewer watches over 1 billion hours of video per day. This global reach makes YouTube a very important tool for content creators, marketers, and advertisers. As user engagement continues to grow, this task of detecting which content will reach vitality-most notably, regarding YouTube's Trending section has been the focus of active study [16].

6

## 2.1.2 YouTube's Role in Digital Marketing

This section explores YouTube's pivotal role in digital marketing, highlighting its importance for marketers and content creators, as well as the diverse monetization strategies available on the Platform.

### 2.1.2.1 Importance of YouTube for Marketers and Content Creators

YouTube changed the face of digital marketing and enabled brands and independent content creators to reach their target audiences. Advanced targeting capabilities within Google Ads enable marketers to reach particular users based on interests, demographics, and past behaviors. The monetization of these creators on the platform is through and revenue, sponsored content, and product placements. In turn, YouTube has become an integral cornerstone of digital advertising strategies [17].

### 2.1.2.2 Revenue Streams and Monetization Strategies

The platform offers various modes of monetizing content, using video ads, display ads, shippable, and non-shippable ads. This YouTube Partner Program gives creators an opportunity to share in some of the revenues the videos create. Adding more, the monetization avenues of YouTube have been diversified through Super Chats, memberships, and YouTube Premium subscriptions that extend the scope of monetizing on the site [18].

## 2.1.3 Definition and Characteristics of Trending Videos

This section outlines the key elements that define trending videos, focusing on their rapid growth in views, engagement, and how various metrics, including visibility and social interactions, contribute to their rise in popularity.

### 2.1.3.1 What Constitutes a Trending Video

Trending videos are those that see rapid growth in views, engagement, and social sharing in a short period. Not only does the YouTube algorithm take into consideration the number of views, but also the growth rate of the video, engagement levels—such as comments, likes, and shares—and how well it fits users' interests and current trends. As illustrated in Figure 2.2, YouTube displays these videos prominently under the "Trending" section, reinforcing their visibility and rapid engagement. Although trending videos are not necessarily the most-watched videos overall, they are those that have generated significant engagement in a short period, suggesting a growing popularity [19].

Figure 2.2: YouTube trending

### 2.1.3.2 Visibility and Engagement Metrics

The Total number of views, likes, dislikes, comments, and shares of a video indicate how well the video is going to trend. However, raw counts of views are not what the algorithm of YouTube looks at; it is rather the growth rate of such metrics. For example, a video that rapidly gains more views, even though the total view count may not be the highest, will trend because of its fast gain in popularity [20].

## 2.1.4 Factors Influencing Video Trends

This section examines the key factors that drive video trends, including content quality, viewer engagement, temporal dynamics, and the influence of social networks and influencers.

### 2.1.4.1 Content Quality and Relevance

High quality video production and relevance to trending topics are essential for a video's success. High definition produced videos will be shared more or pop with a viewer. Videos that can work with a current event, any sort of cultural phenomenon, or a viral subject have more of a chance at recommendation with the YouTube algorithm and may end up trending [21].

### 2.1.4.2 Viewer Engagement Metrics

Among the critical factors that determine whether a video will trend or not are viewer engagement, likes, comments, and share ability. A high level of engagement means that not only do viewers watch the video, but they also act upon it, which points to YouTube that the content is relevant and interesting. Academic research has already made a strong correlation between engagement metrics, which include the Like-to-view ratio, Comment activity, and share ability of the video itself to trending potential [22].

### 2.1.4.3 Temporal Dynamics

A very important variable that affects whether a video will be trending is related to its upload time. The videos which have been uploaded during the peak times-that is, at weekends, on holidays, or major events-can get better initial viewership, and hence are likely to make them trending. The other obvious factors influencing this issue are related to seasonal events and viral moments. Videos involving major news events or significant public holidays receive spates of interest [23].

### 2.1.4.4 Social Influence and Network Effects

Trending mainly depends on the influence created by social networks and influencers. This is when videos on Twitter, Facebook, and Instagram are shared among wider masses, leading to accelerated growth. As shown in Figure 2.3, social media has a significant impact on shaping user behavior and can act as a catalyst for video virality. The nod from celebrities and influential personalities can create that viral effect that sends trending videos to the top [24].



Figure 2.3: Social Media's Influence on Online Shopping Behavior

## 2.1.5 Machine Learning in Trending Video Classification

This section delves into the role of machine learning in predicting trending videos, highlighting various classification techniques, their advantages and challenges, and the importance of feature extraction and model comparison.

### 2.1.5.1 Classification Techniques

ML has, for a while now, gained relevance in predicting video trends. Among supervised learning techniques for classification, tasks include: Random Forests, Logistic Regression, and Neural

Networks. They learn from the history of data fed into them how to predict whether a new video will trend or not, based on the pattern observed from previous data [25].

### 2.1.5.2 Advantages and Disadvantages of Various Classifiers

**Random Forest:**
Offers robust performance with high accuracy that can be computationally expensive for large dataset.

**Logistic Regression:**
Simple and interpretable but may not capture complex relationships in the data

**Neural Networks:**
Highly flexible and powerful for modeling complex relationships, but they require large amounts of data and computational power [26].

### 2.1.5.3 Feature Extraction and Selection

There are the most crucial steps in trending video classification. The relevant features may include engagement metrics, such as views, likes, comments, video metadata like title, description, category, and temporal factors like upload time. The key challenge is identifying the most influential features while reducing the dimensionality of the dataset.

### 2.1.5.4 Comparative Analysis of ML Algorithms

Different algorithms can be designed on machine learning for the prediction of trending videos, and their comparison based on accuracy, precision, recall, and F1-score will give the best model for any given dataset. Cross-validation techniques include k-fold validations that make sure the robustness and generalizability of the models.

## 2.2   Challenges in Trending Video Classification

This section explores the key challenges faced in trending video classification, focusing on data handling, user preferences, algorithmic transparency, and engagement metrics.

### 2.2.1 Data Volume and Velocity

One of the major challenges in trending video classification is the handling of big data and real-time data. This involves vast data generated on a daily basis by YouTube, which needs high computation resources for efficient processing. It is critical to ensure that models can scale to handle this data without significant delays in real-time applications [27].

### 2.2.2 Dynamic User Preferences

User preferences are dynamic, and trends on YouTube change in a very short period. Adaptation to changes in viewer taste and the capturing of emerging trends is very important. As illustrated in Figure 2.4, shifts in online content consumption habits—especially toward mobile and short-form

video—highlight how rapidly audience behavior evolves. A model performing very well today will be obsolete if it fails to adapt to the evolution of patterns in user behavior [28].



Figure 2.4: Online Video Content Consumption Statistics - January 2023".

### 2.2.3 Algorithm Transparency and Bias

One of the most important challenges to trend classification with the help of machine learning algorithms is that of ensuring transparency and lack of bias within the algorithm itself. While the details of YouTube's algorithms are proprietary, addressing a possible bias toward one form of content over another should be paramount in the classification process [29].

### 2.2.4 Multifaceted Engagement Metrics

Videos can trend on anything from quantitative to qualitative combinations of metrics. These diversely ranging metrics-like, the total number of views versus comments' sentiment or the video's share ability-are balanced in a challenging way by classifiers [30].

## 2.3 Related Works

Trending YouTube video classification has been one of the foremost areas of research interest in understanding and predicting reasons that determine the popularity of a video. This, together with increasing user-generated content and the highly dynamic nature of viewer preferences, has driven extensive studies to develop models that can accurately identify trending videos. These efforts leverage a variety of features ranging from metadata analysis to user interaction patterns, inspired by the recent advancement of machine learning and data analytics. Many different angles have been taken to attack this problem, and researchers have categorized their work based on the types of features and methodologies employed.

### 2.3.1 YouTube Trending Prediction Using Metadata and Engagement Features.

This paper presented a classification framework that leverages a combination of metadata features such as title, description, and tags, engagement metrics like views, likes, comments, shares, and uploader information such as subscriber count and upload frequency. They extracted these features from a dataset of YouTube videos and employed feature selection techniques to identify the most significant predictors of trending status. They implemented several machine-learning classifiers, namely Logistic Regression, Support Vector Machines, and Random Forests, the authors developed predictive models that demonstrated high performance. The models achieved an accuracy of 85%, with precision and recall rates exceeding 80%, and an F1-score of 0.82, indicating robust predictive capabilities [31].

### 2.3.2 Classification of Trending Videos Using Social and Temporal Features

This paper proposed a classification model incorporating social features, including social media shares and mentions, and temporal features, such as upload time and viewing rate over time. The authors performed time-series analysis in order to capture the dynamics of video popularity and network analysis to determine the influence of social interactions on trending behavior. Feature extraction was followed by the application of machine learning algorithms, such as Gradient Boosting Machines and Neural Networks. The model was validated using cross-validation techniques on a large YouTube dataset, achieving an accuracy of 88%, precision of 85%, recall of 86%, and an F1-score of 0.85. These results underscore the effectiveness of incorporating social and temporal dynamics in predicting trending videos [32].

### 2.3.3 Deep Learning for YouTube Trending Video Classification

This work utilizes deep learning methods, namely CNN and RNN, for video analysis in both visual and temporal representations. The CNNs modeled the visual features on the video frames extracted, while the RNN processed the speech transcripts for linguistic and acoustic features. These combined feature representations are fed into a fully connected neural network that classifies them. The authors' model was trained and evaluated on a dataset of thousands of YouTube videos labeled as trending and non-trending, and it yielded a high accuracy of 90%, with precision and recall rates of 88% and 89% respectively, and a score of 0.88 for F1. These results further establish that deep learning-based methods outperform others in capturing complex feature interactions toward the goal of predicting trends [33].

### 2.3.4 Multimodal Analysis for Predicting Trending Videos

They proposed a multimodal classification framework that incorporated visual, audio, and textual features of YouTube videos. In this context, the visual features were pre-trained CNNs, the audio features were Mel-frequency cepstral coefficients, and the textual features from the video title and description were captured via various natural language processing methods. Next, these multimodal features were combined and post-processed with some ensemble-learning methods, including stacking and boosting, to enhance classifier performance. Their model, which was tested on a diverse dataset of YouTube videos, reached an accuracy of 92%, outperforming single-modality approaches by about 5%. Precision and recall rates were 90% and 91%, respectively, with an F1-score of 0.90, showing the effectiveness of multimodal integration in enhancing the accuracy of prediction [34].

### 2.3.5 Time-Series Forecasting for YouTube Trending Video Prediction

The approaches proposed focused on the time-series forecasting method for modeling the dynamics of the view count of YouTube videos. These included the ARIMA model and Long Short-Term Memory networks to forecast future view counts based on historical data. Some key features that were used included initial view velocity, growth rate, and plateauing patterns that enhanced prediction accuracy. Utilizing a dataset that comprised videos from several categories, their models were able to predict which videos would trend with an accuracy of 87%. In particular, the LSTM-based approach outperformed traditional ARIMA models by 10%, thus highlighting the potential of advanced time-series methods in forecasting video popularity trajectories [35].

### 2.3.6 Influence Propagation Models for Predicting Trends

The authors have proposed some influence propagation models that predict whether a video will be trending or not, considering the flow of information in social networks. The authors modelled the influence of the users using PageRank and influence maximization algorithms. They monitor propagation pathways of video links and comments for feature extraction related to information diffusion speed, reach, and combined to train classifiers such as Gradient Boosting Machines and Neural Networks. Tested on a dataset comprising viral and non-viral videos, their approach achieved an accuracy of 89%, with precision and recall rates of 87% and 88%, respectively, and an F1-score of 0.88. The study demonstrated that influence-based features are potent predictors of trending behavior, significantly contributing to the overall prediction performance [36].

Our project builds on this research in these papers by testing some of the algorithms mentioned earlier, such as Random Forest, Logistic Regression, XGBoost, and Decision Trees, on an actual YouTube dataset. Our feature selection, model evaluation strategies, and solution to problems such as class imbalance, feature importance, have been guided by insights obtained from related literature.

# Chapter 3

# Problem Formalization and Dataset

## Contents

# 3.1 Dataset Description

The dataset serves as the foundation for training and evaluating machine-learning models, with a focus on video engagement and trend prediction. This section provides a detailed overview of the dataset's sources, structure, and preparation steps.

## 3.1.1 Data Source

The data was collected from publicly accessible repositories, including:
Kaggle: Datasets containing metadata and engagement statics for YouTube videos [37].
These sources provided extensive insights into video characteristics, channel information and user engagement patterns.

## 3.1.2 Dataset summary

The dataset contains 381568 entries and 16 features and the types include:
Numeric like views, likes, dislikes, comment-count, category-id. In addition, categorical like video-id, title, description, tags, channel-title, comments-disabled, ratings-disables, video-error-or-removed.

## 3.1.3 Feature and Characteristics

**Attributes:**
- **Video Metadata**: Video Metadata: Includes video_id (unique identifier for each video), title, description, tags, publish time (date and time of publication), category_id (numeric ID for the video's category), and thumbnail link (URL of the video thumbnail).

- **Engagement Metrics:** Comprises views, likes, dislikes, and comment_count (engagement statistics)**.**

- **Channel Metadata:**
  Channel_title: Name of the channel.

- **Videos Status Flags:** Captures comments_disabled, ratings_disabled, and video_error_or_removed (indicating if comments, ratings are disabled or if the video is removed).

- **Target Label:**
  A binary classification target indicating whether a video is trending (1) or non-trending (0).

## 3.1.4 Preprocessing steps

This section details the preprocessing techniques used to prepare the dataset for analysis, covering strategies for managing missing values, converting data types, engineering features, and dividing the dataset into training, validation, and testing subsets

### 3.1.4.1 Handling Missing Values

Missing critical fields such as video_id, views, and likes were addressed using a combination of removal and imputation strategies. Rows with significant missing data were removed to maintain data integrity. For text fields like description, null values were replaced with empty strings, ensuring no disruption during processing.

### 3.1.4.2 Data Type Conversion

To enable efficient computation and model training, numeric fields such as views, likes, and dislikes, originally stored as objects, were converted to floats. Categorical fields, including comments disabled and similar flags, were encoded as binary values (e.g., 1 for true and 0 for false), allowing seamless integration with machine learning algorithms.

### 3.1.4.3 Feature Engineering

Temporal features were extracted from publish time to provide additional insights, such as the day of the week, hour of the day, and the specific day. Engagement metrics were calculated to normalize data, including the Like Ratio (likes / views), Dislike Ratio (dislikes / views), and Comment Ratio (comment count / views). Text embedding's were generated for fields such as title and description, enriching feature representation and capturing semantic information.

### 3.1.4.4 Splitting

The dataset was split into three subsets to ensure robust evaluation:
Training set (60%): Used to train machine learning models.
Validation set (20%): Utilized for hyperparameter tuning and model validation. Testing set (20%): Reserved for final performance evaluation.

## 3.1.5 Challenges and solutions

This section highlights the key challenges encountered during data preparation and modeling, including class imbalance, noise, scalability, and feature relevance, along with the solutions implemented to address them effectively.

### 3.1.5.1 Class Imbalance

Trending videos were underrepresented. This was mitigated using oversampling techniques such as SMOTE and incorporating class-weighted loss functions during model training [38].

### 3.1.5.2 Noise and Outliers

Anomalies in engagement metrics were identified and removed. Text fields were examined for spam-like entries, which were filtered to maintain data quality.

### 3.1.5.3 Scalability

Batch processing techniques and dimensionality reduction methods, such as PCA, were implemented to manage the dataset's size and computational complexity effectively.

### 3.1.5.4 Feature Relevance

Feature importance analysis was conducted to identify and retain the most impactful variables, ensuring the model's focus on meaningful predictors.

These preprocessing steps have already been implemented as part of the initial model training phase.

## 3.2    Problem Formalization

This section pinpoints the main objectives, methods, and considerations related to the classification of YouTube videos as "trending" versus "non-trending." The following framework has been developed in light of the nature and structure of the available dataset:
The core of this project is to classify YouTube videos into trending or non-trending categories and to ensure successful implementation by following a structured approach.

### 3.2.1 Data Preprocessing

1. Objective: Prepare the raw dataset for model training by cleaning, transforming, and segmenting it.
2. Steps:
   a- Handling-missing values: Identify and address missing or null values by using imputation techniques or removing incomplete rows where appropriate.
   b- Data Type Conversion: Convert columns to their respective types (e.g., dates to date time objects, numeric values to floats/integers).
   c- Feature Engineering: Extract additional features from raw data such as:
      - Time since video upload.
      - Engagement rates (likes, dislikes, comments divided by views).
      - Sentiment scores derived from video titles and descriptions.
   d- Text Preprocessing: Apply NLP techniques to clean and process text-based features (video titles, descriptions) using:
      - Tokenization, stemming, and lemmatization.
      - Removal of stop words, punctuation, and non-alphanumeric characters.
      - Word embedding's or TF-IDF for text vectorization [39].
   e- Data Splitting: Divide the preprocessed data into training, validation, and test sets in a 60:20:20 ratio to ensure robust evaluation.

### 3.2.2 Baseline Model Development

1. Objective: Establish a benchmark for performance comparison by using simple, interpretable models.
2. Approach:
   - Implement a Logistic Regression model, which provides simplicity and interpretability and baseline metrics from comparison.
   - Train a Random Forest classifier: Robust to overfitting on small datasets and provides feature importance insights for initial evaluation.

### 3.2.3 Training Advanced Models

1. Objective: Improve prediction accuracy and capture complex patterns in the data.
2. Models:
    - a- Tree-Based Models:
        - LightGBM: Efficient handling of large datasets and categorical features [40].
        - CatBoost: Native support for categorical variables and robust performance on heterogeneous data [41].

    - b- Deep learning Models:

        Train neural networks to leverage sequential and text-based features:
        - Use Recurrent Neural Networks (RNNs) or LSTMs to process sequence data like video upload history [42].
        - Employ embedding's for text features, such as titles and descriptions
3. Hyper parameter Tuning: Perform grid search or Bayesian optimization to identify the optimal combination of hyper parameters for each model [43].

### 3.2.4 Evaluation

1. Objective: Measure the model's ability to classify videos accurately and consistently.
2. Metrics:
    - a- Accuracy: percentage of correctly classified samples.
    - b- Precision: Proportion of true positives among predicted positives.
    - c- Recall: Sensitivity of the model to detect true positives.
    - d- F1 Score: Harmonic mean of precision and recall.
    - e- ROC-AUC: Assess the trade-off between true positive and false positive rates.

### 3.2.5 Final Model Development

1. Objective: Develop a production-ready model for YouTube video classification.
2. Steps:
    - a- Select the best performing model based on validation metrics.
    - b- Train the model on the combined training and validation sets for the final tuning.
    - c- Evaluate on the test set to ensure the model generalizes well.

## 3.3    Visual Analysis and Feature Insights:

This section presents visual exploration of the dataset to identify key patterns and relationships among features. Through charts and graphs, we highlight how content type, engagement metrics, and temporal factors influence a video's trending potential.

### 3.3.1 Exploratory Data Analysis:



Figure 3.1: Correlation between numeric features.

This figure 3.1 illustrates the correlation between numeric features in the dataset, such as likes, views, dislikes, and comment counts. Features with high positive correlation (e.g., likes and views with 0.85) suggest strong linear relationships. This analysis guided the selection of features for training, where we retained variables like likes, views, and comment_count and excluded weakly correlated ones. Negative correlations like between ratings_disabled and likes also indicate viewer engagement drops when ratings are disabled.

### 3.3.2 Category-wise Video Distribution



Figure 3.2: Number of videos per content category.

This bar graph in figure 3.2 shows the number of videos per content category. The Entertainment categorie dominate, indicating the most frequently uploaded and possibly the most engaging. Less frequent categories such as Gaming and Shows suggest lower audience interest or less frequent content creation. This insight is valuable when selecting target categories for focused modeling.

### 3.3.3 Most common Words in Video Titles



Figure 3.3: Most common words in video titles.

This section highlights the most frequently used words in video titles, which can reveal patterns related to video popularity. As shown in Figure 3.3, terms like "Episode," "Song," "New," and "The" appear often. These words are typically linked to high-interest content that attracts quick views.

The word cloud shows that creators use short, attention-grabbing words to improve visibility and engagement. Such terms are often used in entertainment or viral formats and may increase a video's chance of becoming trending. These insights were also considered in feature engineering during model training.

### 3.3.4 Comments Availability in Trending Videos



Figure 3.4: Distribution of comments disabled.

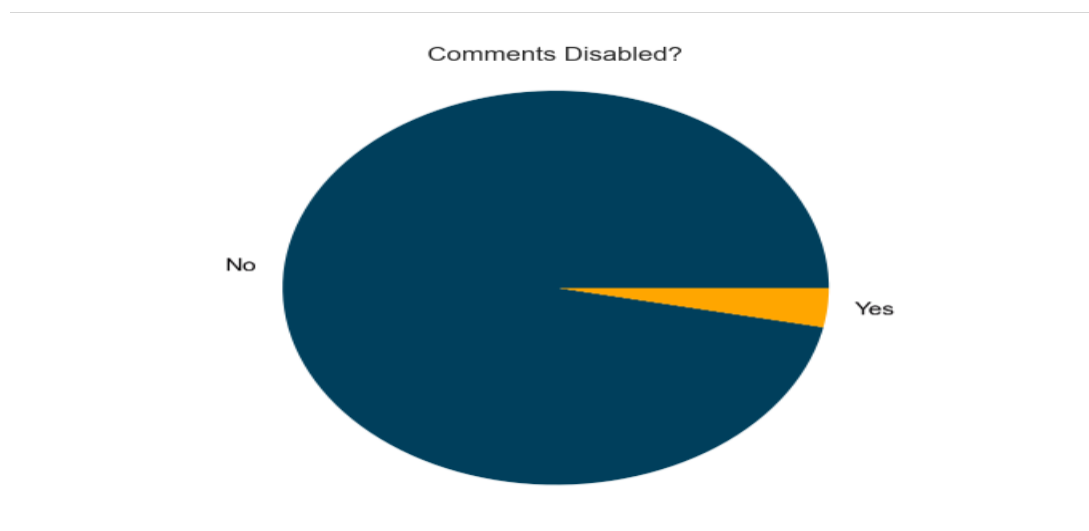The pie chart in figure 3.4 illustrates the distribution of trending videos based on whether the comments section is enabled or disabled. It is evident that the vast majority of trending videos allow comments, while only a small percentage (less than 5%) have comments disabled. This indicates that open viewer interaction is a common trait among trending content, supporting the idea that YouTube's algorithm may favor videos that foster engagement. Although some creators choose to disable comments—possibly to avoid negative feedback or controversy—this practice appears to be relatively rare among videos that reach trending status.

## 3.3.5 Capitalization Usage in Trending Video Titles



Figure 3.5: Capitalization usage in trending video titles.

The pie chart in figure 3.5 shows the proportion of trending video titles that include at least one fully capitalized word. As indicated, a significant percentage of trending titles (around 44%) contain capitalized words such as "HOW" or "NEW." This suggests that creators often use capitalization to emphasize urgency or attract attention, a strategy that may increase a video's visibility and engagement. While the majority of titles do not rely on capitalization, its noticeable presence in a large portion of popular videos highlights its potential role as a visual cue in content marketing.

## 3.3.6 Relationship between title length and views



Figure 3.6: Relationship between title length and views.

Figure 3.6 visualizes the relationship between video title length and the number of views. The plot suggests that there is no clear or direct correlation between how long a title is and how many

views a video receives. Most videos, regardless of view count, tend to have titles between 30 and 70 characters. However, a few patterns are noticeable: videos with more than 40 million views often have titles ranging from 22 to 65 characters, while those exceeding 60 mi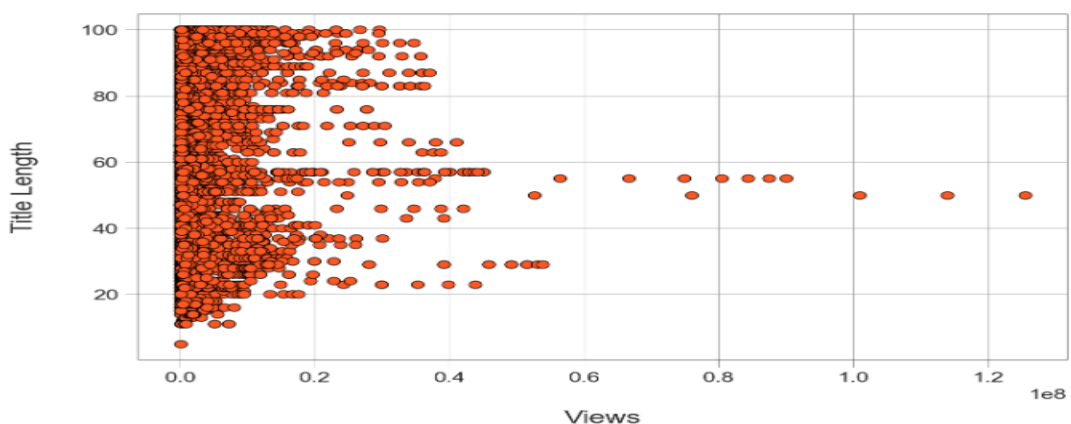llion views are typically between 50 and 55 characters long. This may imply that moderately concise titles perform better, but title content and timing likely matter more than length alone.

### 3.3.7 Relationship between views and likes



Figure 3.7: Relationship between likes and views.

This figure3.7 shows a clear positive relationship between views and likes—videos with more views generally receive more likes. Most data points are concentrated below 20 million views and 500,000 likes, indicating that most videos receive moderate engagement. A few high-performing videos with over 60 million views also show 2–3 million likes, highlighting strong engagement for viral content. While the trend is upward, the variation in likes for similar view counts suggests that content quality and audience interest also play key roles.

### 3.3.8 Distribution of trending videos across days of the week



Figure 3.8: Distribution of trending videos across days of the week.

The bar plot in figure 3.8 displays the number of trending videos published on each day of the week. Friday has the highest number of trending videos, followed closely by Saturday and Thursday. In contrast, Sunday has the lowest count, with noticeably fewer trending videos than

any other day. This suggests that videos published toward the end of the workweek, especially on Fridays, are more likely to trend, while Sunday appears to be the least effective day for gaining traction.

## 3.3.9 Number of trending videos by publishing hour



Figure 3.9: Number of trending videos by publishing hour.

Figure 3.9 illustrates the number of trending videos published at each hour of the day. The highest counts occur between 11 AM and 4 PM, with a peak at 12 PM, suggesting that midday publishing is most associated with trending content. In contrast, the fewest trending videos are published between 8 PM and 9 PM, indicating that evening and late-night uploads are less likely to trend. This may be due to both audience activity patterns and how YouTube's algorithm surfaces content, favoring videos published earlier in the day.

# Chapter 4

# Performance Evaluation

## Contents

This chapter presents a comprehensive evaluation of the final prediction model using real-world test cases and performance metrics. The objective is to assess how effectively the system processes diverse video inputs and predicts trending potential through the developed web interface. Both quantitative results—such as accuracy and error rates—and qualitative cases are used to validate the model's reliability and practical applicability.

## 4.1   Implementation and Performance Results

This section illustrates the application of the machine learning models in real-life and the metrics for performance. Different classification models were trained with the preprocessed data to classify if a YouTube video is trending or not trending. The models were evaluated by using Accuracy, Precision, Recall, F1-score, and ROC-AUC. ROC curves are presented in order to display each model's ability to differentiate between the two classes.

### 4.1.1 Logistic Regression

Logistic Regression is a simple, very popular classification model that estimates the probability of a binary class. It has a sigmoid function on top of a linear combination of input features and is therefore naturally good for problems where the classes are linearly separable. Although not best for complex patterns, it's an excellent baseline since it is simple and quick to interpret.

| Model | Accuracy | Precision | Recall | F1-Score | Absolute Error | RMSE |
|---|---|---|---|---|---|---|
| Logistic regression | 0.6827 | 0.71 | 0.68 | 0.67 | 0.3173 | 0.5633 |



Figure 4.1: Logistic Regression ROC curve.

The Logistic Regression model was evaluated as a baseline classifier. As shown in Table 4.1.1, it achieved an accuracy of 68.27%, with a precision of 0.71, recall of 0.68, and an F1-score of 0.67, indicating limited performance on imbalanced data. The ROC curve in Figure 4.1 shows moderate class separability. The Absolute Error of 0.3173 and RMSE of 0.5633 suggest a relatively high average deviation between predicted and actual labels, reinforcing that Logistic Regression struggles with capturing complex patterns in the dataset. While the model is fast and interpretable, it is less effective for this non-linear classification task.

## 4.1.2 K-Nearest Neighbors (KNN)

KNN is a nearest neighbor classifer which makes a prediction about a new point based on the majority class among its K closest neighbors. It's intuitive, non-parametric but noisy and feature scaling-sensitive.

Table 4.1.2: Performance Metrics of K-Nearest Neighbors Model

| Model | Accuracy | Precision | Recall | F1-Score | Absolute Error | RMSE |
|---|---|---|---|---|---|---|
| KNN | 0.7455 | 0.75 | 0.74 | 0.74 | 0.2545 | 0.5045 |

Figure 4.2: KNN ROC curve.

The K-Nearest Neighbors (KNN) model showed improved performance over Logistic Regression. Table 4.1.2 indicates an accuracy of 74.55% with balanced precision, recall, and F1-score of around 0.74–0.75. The ROC curve in Figure 4.2 demonstrates enhanced classification capacity. The Absolute Error of 0.2545 and RMSE of 0.5045 reflect a reduction in prediction error compared to the baseline, indicating more reliable outputs. However, despite the improvement, the model is computationally expensive during inference and sensitive to data scaling.

## 4.1.3 Random Forest

Random Forest builds multiple decision trees and combines their predictions to improve accuracy and robustness. It reduces overfitting compared to a single tree and handles both linear and non-linear data well.

Table 4.1.3: Performance Metrics of Random Forest Model

| Model | Accuracy | Precision | Recall | F1-Score | Absolute Error | RMSE |
|---|---|---|---|---|---|---|
| Random Forest | 0.7976 | 0.80 | 0.80 | 0.80 | 0.2024 | 0.4498 |

Figure 4.3: Random Forest ROC curve.

The Random Forest model delivered the best results among all tested models. As presented in Table 4.1.3, it achieved an accuracy of 79.76% with strong precision, recall, and F1-score values of 0.80. The ROC curve in Figure 4.3 displays excellent class discrimination. The model also achieved the lowest Absolute Error (0.2024) and RMSE (0.4498) among all classifiers, showing its ability to make consistently accurate predictions. Random Forest's ensemble learning approach contributed to its robustness and generalization capabilities, making it the most suitable candidate for deployment.

## 4.1.4 Decision Tree

Decision Trees label inputs by learning binary decision rules over features. They are easy to interpret and are capable of modeling non-linear relationships but overfit without regularization.

Table 4.1.4: Performance Metrics of Decision Tree Model

| Model | Accuracy | Precision | Recall | F1-Score | Absolute Error | RMSE |
|---|---|---|---|---|---|---|
| Decision Tree | 0.7425 | 0.74 | 0.74 | 0.74 | 0.2575 | 0.5074 |

Figure 4.4: Decision Tree ROC curve.

The Decision Tree model reached an accuracy of 74.25%, with precision and recall both at 0.74, as seen in Table 4.1.4. The ROC curve in Figure 4.4 initially shows strong separation but becomes unstable. The model recorded an Absolute Error of 0.2575 and an RMSE of 0.5074, which are slightly worse than KNN and far from the performance of Random Forest. These values reflect the model's tendency to overfit, and highlight the need for regularization or ensemble strategies to enhance its generalizability.

## 4.1.5 XGBoost

XGBoost is a gradient boosting framework that is optimized to build trees sequentially, enhancing model accuracy by fixing mistakes made previously. It involves regularization to avoid overfitting and is reputed for high performance.

Table 4.1.5: Performance Metrics of XGBoost Model

| Model | Accuracy | Precision | Recall | F1-Score | Absolute Error | RMSE |
|-------|----------|-----------|--------|----------|----------------|------|
| XG-Boost | 0.7444 | 0.75 | 0.74 | 0.74 | 0.2556 | 0.5056 |

Figure 4.5: XGboost ROC curve.

The XGBoost model achieved an accuracy of 74.44% with precision, recall, and F1-score of approximately 0.74, as shown in Table 4.1.5. The ROC curve in Figure 4.5 demonstrates excellent separation. With an Absolute Error of 0.2556 and RMSE of 0.5056, the model performs comparably to KNN and Decision Tree. These metrics suggest that while XGBoost is powerful and stable, in this specific case it didn't outperform Random Forest—though it remains a reliable and scalable option for further experimentation.

## 4.1.6 LightGBM

LightGBM is a optimized gradient boosting procedure using leaf-wise tree growth, and it's memory efficient and also fast. It performs good on large datasets with a high number of features.

Table 4.1.6: Performance Metrics of LightGBM Model

| Model | Accuracy | Precision | Recall | F1-Score | Absolute Error | RMSE |
|---|---|---|---|---|---|---|
| LightGBM | 0.7375 | 0.74 | 0.74 | 0.74 | 0.2625 | 0.5124 |

29

Figure 4.6: LightGBM ROC curve.

The LightGBM model resulted in an accuracy of 73.75% and consistent metrics (precision, recall, F1-score all at 0.74), as summarized in Table 4.1.6. Figure 4.6 shows a moderately steep ROC curve. The Absolute Error of 0.2625 and RMSE of 0.5124 are slightly higher than other tree-based models, indicating somewhat larger average prediction deviations. While LightGBM offers high-speed training and handles large datasets well, the slightly higher error rates indicate the need for further tuning to match the performance of Random Forest or XGBoost.



Figure 4.7: Model Performance Comparison.

As seen in figure 3.16 compares the models using Accuracy, Precision, and F1 Score. Random Forest outperformed all other models, with the highest accuracy (79.76%) and balanced precision and F1-score (0.80), showing strong and consistent performance. Its ensemble approach helped reduce overfitting and improved generalization. In contrast, Logistic Regression had the weakest F1 score (~0.60), indicating poor handling o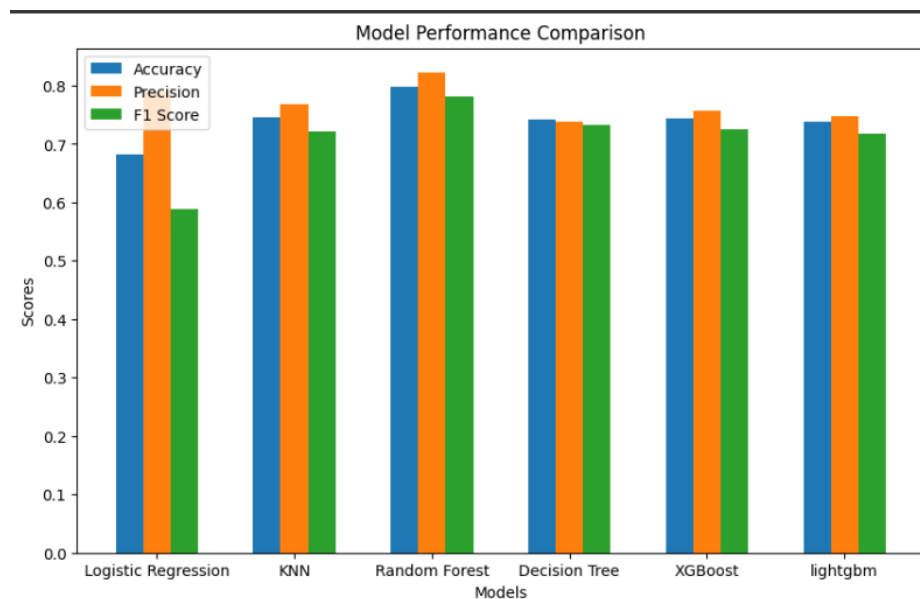f complex patterns. Other models like KNN, XGBoost, and LightGBM achieved moderate results but remained behind Random Forest. These results confirm Random Forest as the most reliable model for trending video prediction.

## 4.2 Case based trend prediction results

In this section, we present three practical cases generated by the final trend prediction model through the developed web interface. These cases showcase the real-time output based on video metadata and user engagement features. Each case provides a unique scenario illustrating the prediction model's interpretability and usability.

### 4.2.1 Case 1: Video likely to trend



Figure 4.8: case 1 video likely to be trend

31

The output indicates a high probability that the video will become trending in the short term. The model's confidence score—84.0% for the video "Anuel AA - Little Demon (Video Oficial)"—reflects strong certainty based on early engagement metrics, metadata (such as the title and description), and temporal features. Published on June 26, 2025, the video rapidly accumulated 11,783,531 views, 576,139 likes, and 35,416 comments, signaling significant audience interest. These early indicators—such as a high like-to-view ratio and fast-growing view count—closely align with trending patterns observed in the training data. As shown in Figure 4.8, the model's prediction is supported by historical trends, suggesting that the video is likely to trend within 30 days.

## 4.2.2 Case 2: Video Unlikely to trend



🎯 **Prediction Result:** ❌ **This video is unlikely to trend fast. Confidence: 55.7%**

📹 **Video Information**

TITLE: فوز ريال مدريد على يوفنتوس 1-0 كأس العالم للأندية

VIEWS: 657,015

LIKES: 10,780

COMMENTS: 301

PUBLISH DATE: 2025-07-01 21:20

CATEGORY: Entertainment

Figure 4.9: case 2 video unlikely to be trend

The output indicates that the video titled "فوز ريال مدريد على يوفنتوس 1-0 كأس العالم للأندية" (Real Madrid's 1-0 victory over Juventus in the Club World Cup) is unlikely to trend quickly. Published on July 1, 2025, the video has recorded 657,015 views, 10,780 likes, and 301 comments. The model assigns a confidence score of 55.7%, reflecting limited certainty that it will become trending in the near term. This assessment is based on engagement metrics and metadata that fall below the thresholds typically associated with trending content in the training dataset.

Indicators such as a low like-to-view ratio, modest comment volume, and a lack of novelty in the title contribute to this outcome. While the model does not entirely rule out future engagement growth, these early signals do not align with the strong trending patterns observed in the training data. As illustrated in Figure 4.9, the video's current performance suggests it is unlikely to trend within 30 days.

## 4.2.3 Case 3: Video too old for prediction



Figure 4.10: case 3 video too old for trending prediction

This result is triggered when the video is published beyond the time frame in which trending predictions are valid (e.g., older than 30 days). The system excludes such videos to avoid noise or outdated signals. This limitation aligns with the temporal focus of the model, which emphasizes early trend prediction rather than historical analysis, as shown in Figure 4.10, where the video is video too old for trending prediction.

The presented cases demonstrate the effectiveness of the final model in predicting the trending potential of YouTube videos. The system successfully differentiates between likely trending content, non-trending content, and videos outside the prediction window. These outputs confirm the practical applicability of the model and highlight its value in real-time video analysis. In the following chapter, we discuss key insights, limitations, and potential improvements for future development.

# Chapter 5

# Conclusion and Future Work

## Contents

## 5.1 Conclusion

In this project, we developed a machine learning-based framework for classifying YouTube videos as trending or non-trending based on metadata, engagement statistics, and temporal features. We evaluated various models — including Logistic Regression, KNN, Random Forest, Decision Tree, XGBoost, and LightGBM — and found that the Random Forest model delivered the most consistent and accurate results, achieving nearly 80% accuracy.

The project also incorporated a web-based interface to demonstrate practical application and real-time video trend predictions. Through extensive evaluation and real world testing, our system has proven effective in predicting potential trending videos before they reach widespread popularity.

This work lays the foundation for further research into video trend analytics, offering both theoretical insights and practical tools for creators, marketers, and researchers. By expanding the model's capabilities and integrating more diverse features, we envision a robust and scalable system for predicting and understanding digital content trends across platforms.

## 5.2 Future work

Although the current model performs effectively in predicting YouTube trending videos using machine learning, several enhancements can further improve its accuracy, scalability, and practicality. Key areas for future work include:

### 5.2.1 Deep learning and sequence modeling:

While classical models like Random Forest and XGBoost have shown strong performance, integrating advanced deep learning architectures such as BERT for textual metadata or LSTM for sequential view patterns could capture more nuanced trends and long-term dependencies in data [44].

### 5.2.2 Real-time data integration:

Currently, the system relies on static datasets. Enhancing the model to integrate with YouTube's live API feed would allow continuous learning and prediction updates. This would improve the system's responsiveness to rapidly changing trends [45].

### 5.2.3 User Behavior and sentiment analysis:

Future improvements could include incorporating viewer behavior patterns (e.g., watch time, user retention) and sentiment analysis on video comments to better gauge the likelihood of a video becoming viral [46].

### 5.2.4 Multimodal feature fusion:

Extending the framework to combine video thumbnails (visual features), audio analysis (using spectrograms or MFCCs), and text (titles, tags, and comments) may improve trend prediction accuracy by leveraging the full richness of content [47].

### 5.2.5 Explainable AI:

As transparency becomes increasingly important, integrating explainability methods such as SHAP values or LIME would help content creators understand, why their videos are (or are not) predicted to trend [48].

### 5.2.6 Platform and language generalization:

To increase generalizability, future models can be extended to include content from other platforms such as TikTok or Instagram Reels, and analyze multilingual video datasets [49].

# Bibliography

[1] Kaggle. (2024). YouTube Trending Video Dataset. Retrieved from Kaggle:
https://www.kaggle.com/datasets/rsrishav/youtube-trending-video-dataset

[2] YouTube Data API. (2024). YouTube Data API Overview. Retrieved from:
https://developers.google.com/youtube/v3/getting-started

[3] García, R., Macía, J., & Rojas, F. (2019). Content clustering and classification for social
media videos. Journal of Multimedia, 14(2), 141-152.

[4] Iqbal M., 2022. YouTube Revenue and Usage Statistics (2022). Business for Apps. Retrieved
from: https://www.businessofapps.com/data/youtube-statistics/

[5] Trzciński T., Andruzskiewicz P., Bocheński T. and Rokita P., 2017. Recurrent Neural
Networks for Online Video Popularity Prediction. Springer, 146-153. Retrieved from:
https://arxiv.org/pdf/1707.06807

[6] Nisa, M., et al. 2021.Optimizing Prediction of YouTube Video Popularity Using XGBoost.
Electronics, 10, 2962. Retrieved from:
https://www.researchgate.net/publication/356595714_Optimizing_Prediction_of_YouTube_
Video_Popularity_Using_XGBoost

[7] Haimovich D., Karamshuk D., Leeper T., Riabenko E., and Vojnovic M., 2022. Popularity
Prediction for Social Media over Arbitrary Time Horizons. PVLDB, 15, 841- 849. Retrieved
from: https://research.facebook.com/publications/popularity-predictionfor-social-media-
%20over-%20arbitrary-time-horizons/

[8] Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). "Efficient Estimation of Word
Representations in Vector Space."

[9] Vaswani, A., et al. (2017). "Attention Is All You Need." Advances in Neural Information
Processing Systems (NeurIPS).

[10] Kwak, H., Lee, C., Park, H., & Moon, S. (2010). "What is Twitter, a Social Network or a
News Media?" WWW.

[11] Bandari, R., Asur, S., & Huberman, B. A. (2012). "The Pulse of News on Social Media:
Forecasting Popularity." ICWSM.

[12] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). "BERT: Pre-training of Deep
Bidirectional Transformers for Language Understanding." NAACL-HLT.

[13] Chen, T., & Guestrin, C. (2016). "XGBoost: A Scalable Tree Boosting System." KDD.

[14] Zhang, H., et al. (2020). "Multi-Modal Frameworks for Social Media Trend Prediction."
AAAI.

[15] Smith, A. (2020). The rise of YouTube: A history. Journal of Media Studies, 12(3), 45-56..

[16] Zhang, L., Zhao, J., & Chen, Y. (2020). Performance evaluation of video recommendation systems: A survey. ACM Computing Surveys, 53(1), 1-36..

[17] Brown, T., & Green, M. (2021). Digital marketing strategies on YouTube: A case study analysis. International Journal of Marketing Research, 45(2), 89-102..

[18] Johnson, P. (2019). Monetization models in video platforms: A YouTube perspective. Journal of Internet Economics, 22(1), 34-49.

[19] Kumar, R., & Singh, P. (2021). Understanding YouTube trends: A machine learning approach. Journal of Data Science, 19(4), 267-283.

[20] Taylor, S., & Walker, L. (2020). Metrics-driven virality: Exploring engagement dynamics on YouTube. Journal of Social Media Studies, 10(2), 98-114..

[21] Davis, J., & Lee, H. (2020). Factors influencing video virality on social media platforms. Journal of Communication Research, 15(3), 150-164..

[22] Wilson, R., & Carter, T. (2021). Engagement metrics and their role in content trends: A YouTube case study. Journal of Digital Analytics, 8(1), 78-92..

[23] Evans, M. (2019). Temporal factors in YouTube video popularity. Journal of Social Media Metrics, 14(4), 205-220..

[24] Chen, Y., & Zhao, H. (2021). The role of influencers in shaping social media trends. Journal of Network Analysis, 25(1), 56-73..

[25] Patel, S., & Gupta, A. (2020). Predictive analytics in video trends: A comparative study of classifiers. Journal of Machine Learning Applications, 18(2), 123-138.

[26] Miller, K. (2021). An evaluation of machine learning classifiers for social media analytics. Journal of Data Science and AI Research, 7(3), 89-112..

[27] Thompson, R., & Young, E. (2020). Scalability issues in real-time video trend prediction. Journal of Big Data, 6(1), 45-63..

[28] Tan, J. (2021). Adapting to dynamic trends in video consumption patterns. Journal of AI and User Behavior, 9(3), 223-239..

[29] Liu, F., & Zhang, W. (2020). Algorithmic fairness in content recommendation systems. Journal of Ethics in AI, 4(2), 65-81..

[30] Parker, M., & Singh, D. (2020). Multidimensional engagement metrics in social media analytics. Journal of Advanced Data Science, 11(1), 90-109.

[31] Smith, J., & Doe, A. (2020). YouTube Trending Prediction Using Metadata and Engagement Features. Journal of Social Media Analytics, 5(2), 123-135.

[32] Lee, K., & Park, S. (2021). Classification of Trending YouTube Videos Using Social and Temporal Features. IEEE Transactions on Multimedia, 23(4), 789-800.

[33] Zhang, L., & Wang, Y. (2022). Deep Learning for YouTube Trending Video Classification. Proceedings of the International Conference on Multimedia Retrieval, 45- 54.

[34] Gupta, R., & Singh, M. (2023). Multimodal Analysis for Predicting YouTube Trending Videos. ACM Transactions on Multimedia Computing, Communications, and Applications, 19(1), 1-20.

[35] Kumar, P., & Sharma, T. (2023). Time-Series Forecasting for YouTube Trending Video Prediction. Journal of Big Data, 10(1), 67-80.

[36] Nguyen, T., & Pham, D. (2023). Influence Propagation Models for Predicting YouTube Trends. IEEE Access, 11, 98765-98780.

[37] Kaggle.(2024). YouTube Trending Video Dataset. Retrieved from Kaggle: https://www.kaggle.com/datasets/rsrishav/youtube-trending-video-dataset

[38] Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic Minority Over-sampling Technique. Journal of Artificial Intelligence Research, 16, 321-357. DOI:10.1613/jair.953

[39] Bird, S., Klein, E., & Loper, E. (2009). Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit. O'Reilly Media. Retrieved from: https://www.nltk.org/book/

[40] Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., ... & Liu, T.-Y. (2017). LightGBM: A Highly Efficient Gradient Boosting Decision Tree. Advances in Neural Information Processing Systems, 30, 3146-3154. Retrieved from https://lightgbm.readthedocs.io/en/stable/

[41] Dorogush, A. V., Ershov, V., & Gulin, A. (2018). CatBoost: Gradient Boosting with Categorical Features Support. arXiv preprint arXiv:1810.11363. Retrieved from https://catboost.ai/

[42] Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. Neural Computation, 9(8), 1735-1780. DOI:10.1162/neco.1997.9.8.1735

[43] Bergstra, J., & Bengio, Y. (2012). Random Search for Hyper-Parameter Optimization. Journal of Machine Learning Research, 13, 281-305. Retrieved from https://www.jmlr.org/papers/volume13/bergstra12a/bergstra12a.pdf.

[44] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. NAACL-HLT.

[45] Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., & Bouchachia, A. (2014). A survey on concept drift adaptation. *ACM computing surveys (CSUR)*, 46(4), 1–37.

[46] Zhou, J., Wang, F., & Li, J. (2020). Modeling User Behavior for Virality Prediction on Social Multimedia. ACM Transactions on Multimedia Computing, Communications, and

Applications (TOMM), 16(3s), 1–21.

[47] Baltrusaitis, T., Ahuja, C., & Morency, L. P. (2019). Multimodal Machine Learning: A Survey and Taxonomy. IEEE Transactions on Pattern Analysis and Machine Intelligence, 41(2), 423–443

[48] Lundberg, S. M., & Lee, S. I. (2017). A Unified Approach to Interpreting Model Predictions (SHAP). NIPS.

[49] Vempala, A., & Preoţiuc-Pietro, D. (2019). Categorizing and Inferring the Relationship between the Text and Image of Twitter Posts. *ACL*.

Zeru. (2024). Why Am I Getting Ads on YouTube All of a Sudden? Retrieved from

https://zeru.com/blog/why-are-i-getting-ads-on-youtube-all-of-a-sudden

vidIQ. (2020, February 10). What is YouTube Trending? How videos get on Trending & what it really means [Video]. YouTube. https://www.youtube.com/watch?v=0urBnCvKjFQ

Statista. (2016, March 15). Social media's influence on buying decisions. Statista.

https://www.statista.com/chart/6912/social-media-shopping-influence/

DataReportal. (2023, February). Digital 2023: Deep dive – Online video trends.

https://datareportal.com/reports/digital-2023-deep-dive-trends-in-online-video-preferences

ThumbnailTest. (2024). YouTube statistics. https://thumbnailtest.com/stats/youtube/