

Republic of Iraq
Ministry of Higher Education and Scientific Research
University of Babylon
College of Information Technology
Software Department



Video Popularity Prediction Based on Youtube Metadata and Thumbnail Images

A Thesis

Submitted To the Council of The College of Information Technology For
Postgraduate Studies at University of Babylon in Partial Fulfilment of The
Requirements For the Degree of Master in Information Technology/ Software

By

Heba Hussein Abd-Alabbas Naif

Supervised By

Lecturer. Dr. Wadhah Razooqi Abood Hassan Baiee

2023 A.C.

1445 A.H.

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

أَلَمْ نَشْرَحْ لَكَ صَدْرَكَ ❖
وَوَضَعْنَا عَنكَ وِزْرَكَ ❖
الَّذِي أَقْبَضَ ظَهْرَكَ ❖
وَرَفَعْنَا لَكَ ذِكْرَكَ ❖
فَإِنَّ مَعَ الْعُسْرِ يُسْرًا ❖
إِن مَّعَ الْعُسْرِ يُسْرًا ❖
فَإِذَا فَرَغْتَ فَانصَبْ ❖
وَإِلَىٰ رَبِّكَ فَارْغَبْ ❖

بِسْمِ اللَّهِ
الرَّحْمَنِ الرَّحِيمِ

Dedications

To the Prophet of Mercy, Muhammad, may Allah bless him and grant
him peace:

I dedicate this thesis with deep reverence and respect to the Prophet Muhammad, whose teachings of compassion, wisdom, and guidance continue to inspire countless lives around the world. His legacy serves as a beacon of light and an enduring source of spiritual strength for all of humanity.

To the Imam, the Awaited Mahdi:

Whose profound spiritual presence offers hope and guidance to those who await his appearance.

To My Family:

You have consistently been the spark that reignited my light when it dimmed. Thank you for your unwavering love and support along this journey I have undertaken. I love you all, now and for all time.

Heba Hussien

Declaration

I hereby declare that this dissertation entitled (Video Popularity Prediction Based on Youtube Metadata and Thumbnail Images), submitted to University of Babylon in partial fulfilment of requirements for the degree of Master in Information Technology \ Software, has not been submitted as an exercise for a similar degree at any other University. I also certify that this work described here is entirely my own except for experts and summaries whose source are appropriately cited in the references

Signature:

Name: Heba Hussein Abd-Alabbas

Date: / /2023

Supervisor Certification

I certify that this thesis (Video Popularity Prediction Based on Youtube Metadata and Thumbnail Images) was prepared under my supervision at the Department of Software / Collage of Information Technology / University of Babylon, by **Heba Hussien** as a partial fulfilment of the requirements for the degree of **Master in Information Technology**.

Signature:

Name: Lecturer. Dr. Wadhah Razooqi Baiee

Date: / / 2023

The Head of the Department Certification

In view of the available recommendation, we forward this Thesis for debate by the examining committee.

Signature:

Name: Asst. Prof. Dr. Sura Zaki Alrashid

Date: / / 2023

Certification of the Examination Committee

We, the undersigned, certify that (Heba Hussein Abd-Alabbas) candidate for the degree of Master in Information Technology - Software, has presented his thesis of the following title (Video Popularity Prediction Based on Youtube Metadata and Thumbnail Images) as it appears on the title page and front cover of the thesis that the said thesis is acceptable in form and content and displays a satisfactory knowledge of the field of study as demonstrated by the candidate through an oral examination held on:

Signature:

Name: Israa Hadi Ali

Title: Professor

Date: / /2023

(Chairman)

Signature:

Name: Ayad Hameed Mousa

Title: Assistant Professor

Date: / /2023

(Member)

Signature:

Name: Safa Saad Abbas

Title: Assistant Professor

Date: / /2023

(Member)

Signature:

Name: Wadhah Razooqi Abood

Title: Lecturer

Date: / /2023

(supervisor)

Signature:

Name: Wesam Sameer Bhaya

Title: Professor

Date: / /2023

(Dean of Collage of Information Technology)

Acknowledgements

First and foremost, I offer my praise and thanks to the Almighty, Allah, for His abundant blessings that sustained me throughout the successful completion of my research project.

I would like to extend my profound and heartfelt appreciation to my supervisor, Dr. Wadhah Baiee, for his invaluable guidance and unwavering support during the course of this research.

To my dear family, especially my mother and Aunt Maryam.

I want to take a moment to express my sincere appreciation for your unwavering love and support. Your presence in my life has been my greatest source of strength and encouragement. Thank you for always being there for me, through every up and down. I am truly blessed to have you in my life.

I also extend my gratitude to my brothers and sisters for their endless support. It would not have been possible to write this thesis without their unwavering encouragement. Their love and support have played a significant role in my development and success.

Sincerely, thank you, my dear friend Shahad Hussein, for your inspiring friendship, understanding, and unwavering belief in my abilities. Your presence and support have proven invaluable throughout this challenging endeavor.

Finally, I express my heartfelt gratitude to all those who stood by me and provided moral support during this research journey.

Abstract

The exponential growth of social media platforms and online video content has revolutionized the way information is communicated, interacted with, and consumed. The predicting of video popularity has become a crucial aspect for influencers, platform administrators, and marketers, with millions of videos being shared daily on platforms like YouTube. This thesis delves into the challenge of accurately predicting video popularity on social media platforms, specifically focusing on the YouTube platform. The predicting of video popularity is considered a crucial aspect for influencers, platform administrators, and marketers, with millions of videos being shared daily on platforms like YouTube. The factors contributing to a video's popularity are understood in this thesis, significantly impacting content marketing and YouTube platform strategies.

To achieve the goal of enhancing the precision of forecasting video popularity, an approach has been formulated that utilizes both metadata analysis (including selected impact features) and thumbnail image analysis to create a more effective prediction model. The precision of forecasting video popularity is aimed to be enhanced by extracting relevant features from these sources.

A range of powerful classification algorithms, including Support Vector Machine, Random Forest, Extreme Gradient Boosting, and K-nearest Neighbor, is employed to achieve this goal. Preprocessing operations were conducted on the data to make it suitable for machine learning algorithms. Subsequently, features were extracted in three categories: textual features, visual features, and time-related features. The dataset was enriched with additional impactful attributes capable of enhancing the model's accuracy by combining these newly added features

with the existing ones. Substantial accuracy was achieved by training the model using these algorithms and implementing feature extraction techniques, particularly with the Extreme Gradient Boosting and Random Forest algorithms.

The predicting process was conducted in two stages: predicting both original and extracted features, and predicting only the extracted features. The model's performance was evaluated using various metrics, including accuracy, resulting in rates of 96.33% and 96.94% for the Extreme Gradient Boosting and Random Forest algorithms using all features, respectively, while accuracy reached 84.65% and 93.00% using only the extracted features. Valuable insights into the factors influencing video popularity are provided by combining metadata and thumbnail analysis. This enables influencers and marketers to better personalize their content, resonating with their target audience and increasing the likelihood of success on social media platforms.

Table of Contents

CHAPTER ONE GENERAL INTRODUCTION.....	1
1.1 INTRODUCTION.....	1
1.2 THESIS PROBLEM	3
1.3 THESIS QUESTIONS	3
1.4 AIMS OF THESIS.....	4
1.5 RELATED WORKS	4
1.6 THESIS OUTLINE	9
CHAPTER TWO THEORETICAL BACKGROUND	10
2.1 OVERVIEW	10
2.2 SOCIAL MEDIA AND POPULARITY PREDICATION	10
2.2.1 <i>Predicting Content Popularity</i>	11
2.2.2 <i>Factors Affecting Video Popularity</i>	12
2.3 DATASET	14
2.4 FEATURE EXTRACTION TECHNICAL	16
2.5 DATA PREPARATION METHODS	19
2.5.1 <i>Removal of Row Duplications</i>	19
2.5.2 <i>Data Labelling[1]</i>	20
2.5.3 <i>Handling Missing Values</i>	20
2.5.4 <i>Processing and Extracting Text Features</i>	21
2.5.5 <i>Encoding Features</i>	24
2.6 DATA MINING	24
2.7 EVALUATION OF THE PREDICTION MODELS	34
CHAPTER THREE THE PROPOSED SYSTEM	36
3.1 OVERVIEW.....	36
3.2 THE PROPOSED SYSTEM ARCHITECTURE	36
3.2.1 <i>Extraction Features from Images</i>	38
3.2.2 <i>Data Preprocessing</i>	41
3.2.3 <i>Machine Learning Models</i>	47
3.2.4 <i>Evaluating the Performance of the Proposed Model</i>	47
CHAPTER FOUR RESULTS AND DISCUSSION	49

4.1 INTRODUCTION	49
4.2 RESEARCH ENVIRONMENT	49
4.3 PROPOSED SYSTEM RESULTS	49
4.3.1 <i>Extract Visual Features from Images</i>	49
4.3.2 <i>Data Preparation Results</i>	52
4.3.3 <i>Results of the Classification Methods</i>	63
CHAPTER FIVE CONCLUSIONS AND FUTURE WORKS	69
5.1 CONCLUSIONS.....	69
5.2 FUTUREWORK	70
REFERENCES	72
الخلاصة.....	81

List Of Figures

Figure	Caption	Page Number
Figure 2.1	Factors That Effect Video Popularity	12
Figure 2.2	Sample Of Dataset	15
Figure 2.3	The YOLO Architecture	17
Figure 2.4	Yolo For Object Detection	18
Figure 2.5	Application Of NLP	22
Figure 2.6	Example Process of Lexicon-Based Sentiment Analysis.	23
Figure 2.7	Data Mining Goals	25
Figure 2.8	Classification Of Data by Support Vector Machine (SVM)	27
Figure 2.9	KNN Classification with Large K.	29
Figure 2.10	The Idea of Random Forest	30
Figure 2.11	Boosting Method.	33
Figure 3.1	The Proposed System Architecture	37
Figure 3.2	Sample Of Images Dataset	38
Figure 3.3	Steps Visual Feature Extraction and Merging with Original Dataset.	43
Figure 3.4	The Confusion Matrix of a Four Class.	48
Figure 4.1	The Result of Using Pre-Trained Yolo V5 Model	50
Figure 4.2	Distribution Sentiments Across the Description Features.	59
Figure 4.3	Distribution Sentiments Across the Video Title Features.	60
Figure 4.4	Distribution Sentiments Across the Channel Title Features.	61
Figure 4.5	The Ratio of Positive (+), Negative (-), And Neutral (Zero) Sentiments in The Video Title Feature Description	61

Figure 4.6	The Ratio of Positive (+), Negative (-), And Neutral (Zero) Sentiments in The Video Title Feature	62
Figure 4.7	The Ratio of Positive (+), Negative (-), And Neutral (Zero) Sentiments in the Channel Title Feature	62
Figure 4.8	Separation of Data into Training and Testing Data	64
Figure 4.9	Accuracy Comparison of Classifiers Using all features	66
Figure 4.10	Confusion Matrix of Random Forest Classifier	68

List Of Tables

Table	Caption	Page Number
Table 1.1	Summary of Previous Related Works	7
Table 2.1	Illustration of The Dataset Features	15
Table 4.1	Feature Set Extracted from Images (Before Processing Step)	50
Table 4.2	Feature Set Extracted from Images (After Processing Step).	51
Table 4.3	Shows the Number and What It Represents (Category)	52
Table 4.4	The Result of Labelling the Data into Four Classes	52
Table 4.5	The Results Before of Removing Duplicate Records	52
Table 4.6	The Result After of Removing Duplicate Records	53
Table 4.7	Description Feature Before Handling the Missing Values	53
Table 4.8	Description Feature Before Handling the Missing Values	53
Table 4.9	The Results of Each Step on the Sample Video Description	54
Table 4.10	The Result Before Handling Non-Numerical Features	57
Table 4.11	The Result After Handling Non-Numerical Features	58
Table 4.12	The Textual Features Extraction	58
Table 4.13	Time interval Feature for Each Video in the Dataset	63
Table 4.14	Accuracy Comparison of Classifiers Using all features	65

Table 4.15	Shows That Random Forest (RF) And XGBoost (XGB) Outperformed the Previous Research on the Same Dataset Using Extracted Features Only.	67
------------	---	----

List Of Algorithms

Algorithm	Caption	Page Number
Algorithm 2.1	Pre-Trained YOLO For Object Detection	18
Algorithm 2.2	Lexicon-Based Sentiment Analysis	23
Algorithm 2.3	Random Forest	32
Algorithm 3.1	Retrieve Video Thumbnails.	38
Algorithm 3.2	Extract Features from Images	39
Algorithm 3.3	Labelling Video	41
Algorithm 3.4	The Missing value in Description feature	43
Algorithm 3.5	Convert Non numerical Feature into numerical Feature	44
Algorithm 3.6	Time Interval Feature Extraction	46

List of Abbreviations

Abbreviation	Meaning
AL	Artificial Intelligent
API	Application Programming Interface
CNN	Convolutional Neural Network
Doc2Vec	Document-To-Vector
GBM	Gradient Boosting Machines
GloVe	Global Vectors for Word Representation
https	Hypertext Transfer Protocol Secure
ID	Identification
IOU	Intersection over Union
KNN	K-Nearest Neighbours
LTRCN	Long-Term Recurrent Convolutional Network
MAE	Mean Absolute Error
MRBF	Multivariate Radial Basis Function
MSE	Mean Squared Error
NLP	Natural Language Processing
NMS	Non-maximum suppression
Popularity-SVR	Popularity- Support Vector Regression
ResNet	Residual Network
RF	Random Forest
SMHP	Social Media Headline Prediction
SVM	Support Vector Machine
TF-IDF	Term Frequency-Inverse Document Frequency
URLs	Uniform Resource Locators
VGG	Visual Geometry Group
WL	Weak Learner
Word2Vec	Word-To-Vector
WWW	World Wide Web
XGB	Extreme Gradient Boosting
YOLO v5	You Only Look Once Version 5

Chapter One

General Introduction

CHAPTER ONE

General Introduction

1.1 Introduction

At present, the advent of online video-sharing platforms has revolutionized how media content is consumed and engaged with. Among these platforms, the YouTube has emerged as a dominant player, with billions of users worldwide and an immense library of videos covering a vast range of topics. As YouTube continues to grow, influencers and marketers strive to understand the factors that contribute to the popularity of videos on the platform. Predicting the popularity of YouTube videos has become a subject of great interest, as it offers valuable insights for influencers, marketers, and platform administrators[1,2].

Understanding what makes a video popular on YouTube is not only intriguing from a social perspective but also carries significant practical implications. For influencers, accurately predicting video popularity can guide decisions regarding content creation, title optimization, thumbnail design, and promotional strategies[3,4]. Marketers can leverage predictive models to identify potential viral videos and allocate their advertising budgets effectively[5]. Moreover, YouTube itself can benefit from predictive analytics by enhancing user experience, optimizing recommendations, and attracting more creators to the platform.

In recent years, advancements in data science, machine learning, and natural language processing techniques have opened up exciting possibilities for predicting the popularity of YouTube videos. By analyzing various features of

a video, such as view count, likes, comments, video duration, tags, and channel characteristics, predictive models can be trained to estimate the likelihood of a video becoming popular. These models can uncover patterns and relationships that humans may overlook, enabling more accurate predictions of video performance.

Predicting the popularity of YouTube videos is a multifaceted task with several challenges[3]. The considerable volume of accessible data, the platform's dynamic nature, and the impact of external factors like trending topics and algorithmic changes all contribute to these obstacles[4]. Additionally, the subjective nature of popularity and the inherent unpredictability of viral phenomena add further complexity to the prediction process[5].

This thesis encompasses a thorough analysis of various features that contribute to a video's popularity. Specifically, it investigates the impact of different video metadata, such as the ratio of likes to dislikes and the number of comments, as well as examines the role of video title, description, and thumbnails, among other factors. By carefully analyzing these features, the aim is to uncover the underlying patterns and relationships that drive a video into popularity on YouTube. This analysis will provide valuable insights into the aspects that captivate viewers and influence their engagement with videos.

Moreover, it evaluates the performance of different predictive models in the realm of video popularity prediction. Through rigorous testing and comparison, it assesses the effectiveness of these models in accurately estimating the likelihood of a video becoming popular. This evaluation process will enable identification of the most reliable and accurate predictive models for video popularity.

By integrating exploratory analysis of video characteristics with the evaluation of predictive models, the aim is to make a valuable contribution to the continually expanding domain of social media analytics and predictive modeling. Findings will not only enhance our understanding of YouTube's popularity dynamics but also provide actionable guidance for influencers, marketers, and YouTube administrators in optimizing their video strategies and maximizing audience reach. Ultimately, this research endeavor holds the potential to advance knowledge of online content consumption patterns and shape the future of video analytics and prediction techniques.

1.2 Thesis Problem

The challenge is to accurately predict the popularity of YouTube videos, a critical aspect for influencers, content creators, and marketing companies. The problem encompasses the need to comprehend the diverse factors influencing a video's success.

This entails exploring the effectiveness of machine learning algorithms and predictive approaches in analyzing metadata and user engagement metrics to achieve precise popularity predictions an area that previous studies have not fully addressed[1].

1.3 Thesis Questions

1. Can overall classification accuracy be enhanced in the context of video popularity prediction?
2. What additional features can be incorporated to enhance the richness of the dataset?
3. Can a predictive approach be developed to accurately predict the popularity of YouTube videos?

1.4 Aims of Thesis

The primary aim of this thesis is to develop effective prediction approaches to empower influencers, marketers, and platform administrators in maximizing video popularity across various social platforms like YouTube. To achieve this aim, the following objectives will be covered:

1. Conducting various stages to prepare the dataset for learning algorithms.
2. Extracting visual and textual features that could influence the prediction process. The integration of these features aims to improve prediction rates and provide a more comprehensive understanding of their combined impact on video popularity, thus enhancing the model's performance and accuracy in predicting video trends.
3. Improving the accuracy of prediction, which was limited in previous research[1].

1.5 Related Works

This section reviews various previous research that has addressed the issue of predicting the popularity of specific content on multiple platforms. These studies have explored different features, factors, and algorithms to achieve their predictions.

R. Shreyas et al in 2010. The Random Forest regression model is used in this paper [6]to predict the popularity of articles using the Online News Popularity data set. The Random Forest model's performance is evaluated and compared to that of other models. Standardization, regularization, correlation, strong bias/high variance, and feature selection all have an effect on learning models.

The Random Forest technique predicts popular/unpopular articles with an accuracy of 88.8%, according to the results.

Trzcinski et al in 2015 [7]. proposed a model to predict the popularity of an online video before its content is published, using Support Vector Regression with Gaussian Radial Basis Functions. They showed that predicting popularity patterns with this approach provides more precise and stable results. Furthermore, they demonstrated that combining early distribution patterns with social and visual features improves the accuracy of popularity prediction. In terms of video popularity prediction, social features were found to be a much stronger signal than visual features. The best results were achieved by combining visual features, social features, and early view counts, allowing for the prediction of video popularity on Facebook with a Spearman correlation rank of up to 0.94, just 6 hours after publication.

W. Stokowiec et al in 2017 [8]The researchers used a methodology based on a two-way long-term memory (LSTM) neural network to predict the popularity of content using only textual information from the title, and these approaches showed an improvement in performance compared to traditional shallow approach.

T. Trzcinski et al in 2017[9]The researchers propose a new method based on a Long-term Recurrent Convolutional Network (LRCN) to address the challenge of predicting the popularity of online videos shared on social media. This approach utilizes deep neural network architectures that consider the sequential nature of information in the videos. The popularity prediction problem is formulated as a classification task, with the goal of predicting popularity using only visual cues extracted from the videos. The results of this study demonstrate that their proposed LRCN-based approach outperforms

traditional shallow methods, achieving over a 30% improvement in prediction performance. The experiments are conducted on a dataset comprising more than 37,000 videos published on Facebook.

F. Huang et al in 2018[10]Introduce a thriving application scenario called Social Media Headline Prediction (SMHP), which focuses on predicting the popularity of posts shared on social media. The research proposes a method that utilizes multi-aspect features combined with the random forest (RF) model for popularity predictions. It explains the process of feature extraction by combining metadata of the posts and users' features, as well as strategies for dealing with missing values. The result of this paper indicates that user-related features, such as the number of followers and following, along with the random forest regression model, are the most effective features and model for the current social media headline prediction task. These features and the chosen model have demonstrated strong predictive capabilities in accurately forecasting the popularity of posts on social media platforms.

A. Khan et al in 2018 [11] Researchers successfully predicted the popularity of news articles using news article metadata and concluded that clustering ensemble methods gave the best results in the dataset, with an accuracy of 79.7%.

Y. Li et al in 2019.[1] proposed the use of several machine learning algorithms to predict performance, and backward search is employed to select the most relevant features. As a result, extreme gradient boosting with three features (time gap, category, description) is chosen due to its optimal balance between cost and performance, resulting in an F-score of 0.73.

K. R.Purba et al in 2012. [12] focuses on predicting the popularity of Instagram posts, particularly the Engagement Rate (ER), which is crucial for

marketing purposes. In this study, several regression techniques were compared using a global dataset, and the features for prediction were extracted from hashtags, image analysis, and user history. The results indicate that factors such as image quality, posting time, and the type of image significantly impact ER. The best prediction accuracy, reaching 73.1%, was achieved with Support Vector Regression (SVR), surpassing previous studies on global datasets.

M. U. N. Nisa et al in 2021.[13]proposed a method that predicts the popularity of videos using the XGBoost model. The approach involves features selection, fusion, min-max normalization, and precision parameters such as gamma, eta, and learning rate. The XGBoost model achieved an accuracy of 86% and a precision of 64%, seen table 1.1

Table 1.1 Summary of Previous Related Works

References	Methodology	Content type	Features	Result
[6]	Random Forest model	Articles	Text features	Accuracy =88.8%
[7]	Support Vector Regression with Gaussian Radial Basis Functions	Video	Visual features, social features, early view counts	SCR = 0.94
[8]	Long Short-Term Memory neural network	News articles, News videos	Textual information from the title	Accuracy=0.74 %

[9]	Long-term Recurrent Convolutional Network	Video	Visual cues extracted from videos	Accuracy =0.7
[10]	Random forest model	Posts	Metadata of the posts, Metadata of the users' features	Accuracy =88.8%
[11]	Gradient Boost	News articles	Meta-data of news articles	Accuracy = 79.7%
[1]	Extreme gradient boosting	Video	Time gap, category, description	F-score =0.73
[12]	Several regression techniques	Post	User Features, Post Features, Hashtag Features, Image Assessment (Auto) Features, Image Assessment (Manual) features, User History Features	Accuracy=73.1%
[13]	Xgboost	Video	Video quality, video duration, number of views	Accuracy=86%
The proposed system	Random forest	Video	Text, visual, time interval and statistic features	Accuracy=96.94%

1.6 Thesis Outline

Following Chapter one, which presents a general introduction the rest of the thesis, is structured as follows:

1. Chapter Two (Theoretical Background): This chapter provides an in-depth overview of essential concepts forming the foundation of the thesis. Covered topics include social media, factors affecting video popularity, preprocessing techniques, feature extraction, machine learning algorithms, and performance evaluation metrics.
2. Chapter Three (The Proposed System): In this chapter, it presents the practical aspects of the proposed system, focusing on the algorithms and techniques used in developing the advanced prediction model for video popularity on social networks.
3. Chapter Four (Results and Discussions): This chapter showcases the results obtained from the implementation of the proposed system. The results will be presented using tables, graphs, and visualizations to offer a comprehensive overview of the predictive model's performance. The primary findings and insights derived from the study will be discussed in detail.
4. Chapter Five (Conclusions and Future Works): In this chapter, the essential concepts and findings of the thesis are summarized, accompanied by recommendations for future research directions.

Chapter two

Theoretical Background

CHAPTER TWO

Theoretical Background

2.1 Overview

This chapter provides a comprehensive overview of the fundamental principles underlying social media platforms, datasets, data mining, and natural language processing (NLP). It also covers the techniques involved in preprocessing, text analysis, feature extraction, and prediction algorithms. The primary emphasis of this chapter revolves around the methodologies and strategies employed in this thesis, shedding light on their significance.

2.2 Social Media and Popularity Predication

Social media refers to online platforms that enable users to create, share, and exchange information within virtual communities. It has revolutionized global communication and interaction, offering features for sharing text, photos, videos, and links. Social media's origins can be traced back to the development of computer networks and the internet, with the launch of YouTube in 2005 marking a significant milestone[14]. Since then, influential platforms like Facebook, Twitter, YouTube, LinkedIn, Instagram, and Snapchat have emerged. Social media platforms have not only changed the way people communicate but have also impacted various aspects of society. They have enabled individuals, businesses, organizations, and even governments to engage with audiences, promote products and services, share news and information, and foster communities based on shared interests and values[2]. TikTok, for instance, specializes in video editing and sharing, allowing users to create and share short video clips that can be charming, funny, or even cringe-inducing[15].

The widespread use of smartphones and mobile apps has further fueled its growth[16]. While providing opportunities for self-expression and networking, social media also poses challenges concerning privacy, security, and excessive screen time[15]. Overall, it has transformed the way people connect, communicate, and shape various aspects of personal and professional life.

2.2.1 Predicting Content Popularity

Predicting popularity in social media is a common application of data mining and analytics. With the enormous amount of user-generated content and interactions on platforms such as Facebook, Twitter, Instagram, and YouTube, businesses and researchers are interested in understanding which posts, videos, or content will become popular and gain significant engagement[17]. Data mining techniques can be employed to analyse various factors that contribute to popularity in social media. These factors may include text features like (word embeddings from video descriptions , titles) and visual and metadata features [17]. By examining historical data and patterns, data mining models can be built to predict the potential popularity of web content. Machine learning algorithms can be applied to social media data to uncover patterns and relationships that influence popularity. These algorithms can consider various features, such as textual features ,video descriptions and titles, the source of the content, category, number of views, timing of the post to make predictions about the potential popularity of a post or content[15]. Additionally, sentiment analysis can be used to understand the sentiment or emotional tone expressed in social media posts. Analysing sentiment can provide insights into the factors that contribute to content popularity. For example, positive sentiment in user comments or interactions may indicate a higher likelihood of popularity[18]. By predicting popularity in social media, businesses and influencers can optimize their strategies, identify trends, and improve their chances of reaching a larger

audience. However, it's important to note that predicting popularity in social media is a complex task and can be influenced by various factors, including user preferences, viral trends, and unpredictable events.

2.2.2 Factors Affecting Video Popularity

The factors that affect video popularity can be categorized into two main groups: external factors related to the platform and users, and internal factors related to the video itself [19][4] In figure 2.1, it is displayed. Below is an explanation of each type:

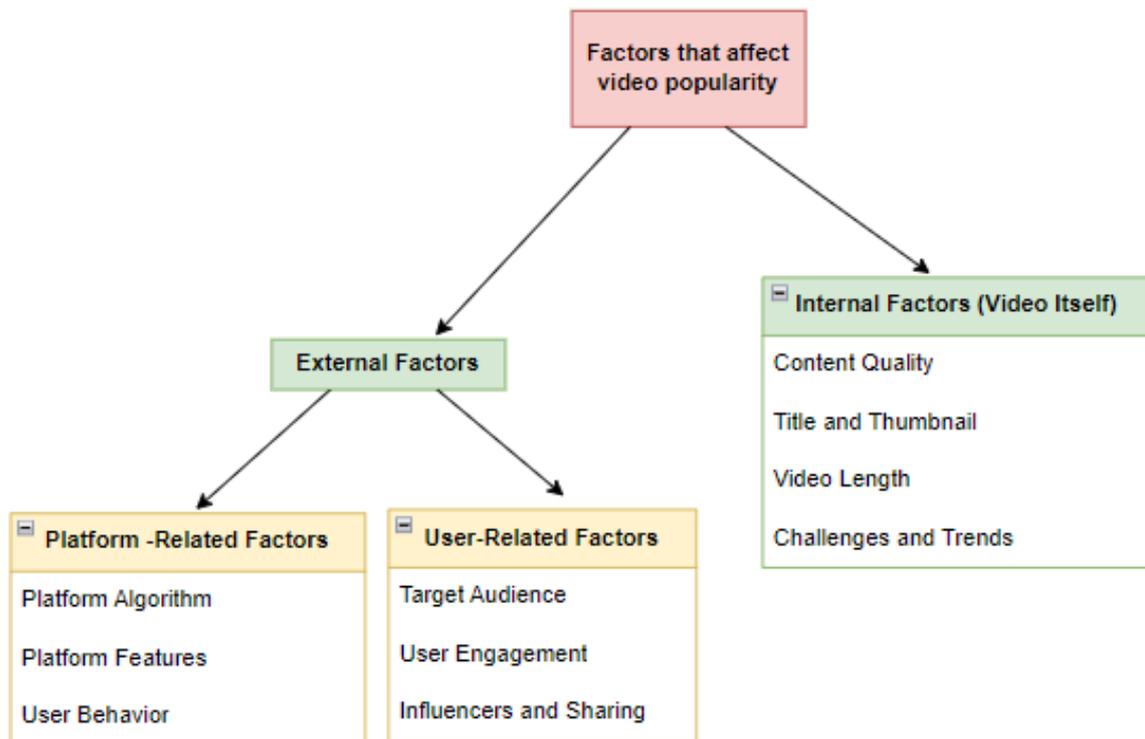


Figure 2.1: Factors That Effect Video Popularity

2.2.2.1 External factors

They encompass elements outside of the video that can impact its popularity, including platform-related and user-related factors.

1. Platform-Related Factors

- a. Platform Algorithm: Each platform has its own algorithm that determines which videos get recommended and shown to a broader audience. Understanding and optimizing for these algorithms can significantly impact video visibility[20].
- b. Platform Features: Utilizing platform-specific features, such as live streaming, stories, or interactive elements, can enhance engagement and popularity[21].
- c. User Behaviour: User behaviour, such as viewing habits, likes, comments, and shares, can influence the visibility and recommendation of a video.

2. User-Related Factors[19]

- a. Target Audience: Understanding the preferences and interests of the target audience is crucial for creating content that resonates with them[19].
- b. User Engagement: Higher engagement metrics like likes, comments, and shares signal to the platform that the video is valuable, potentially leading to increased visibility[2,19].
- c. Influencers and Sharing: When influential users or influencers share or feature a video, it can lead to a significant boost in popularity[2,20, 5].

2.2.2.2 Internal Factors (Video Itself)

Internal factors refer to the characteristics or attributes of the video itself that influence its popularity. These factors include:

1. **Content Quality and Type:** High-quality content, that is informative, entertaining, or valuable to the audience, is more likely to be shared and recommended[3,20].
2. **Challenges and Trends:** Participating in popular challenges or following internet trends can help videos gain visibility and attract a wider audience[22].
3. **Title and Thumbnail:** An attention-grabbing title and thumbnail can encourage users to click on the video, increasing its views[5].
4. **Video Length:** The length of the video is essential. It should be neither too long nor too short, as viewers might lose interest in lengthy videos, while short videos might not provide enough value[5].

It is important to note that the impact of each factor can vary based on the specific platform and content type. Creators should analyse their audience's preferences, experiment with different strategies, and stay adaptable to achieve long-term popularity and growth. In the context of subject, emphasis was placed on relying on internal elements associated with the video itself.

2.3 Dataset

In this thesis, the "Trending YouTube Video Statistics" dataset obtained from Kaggle[23] was utilized. This dataset consists of information on 24,427 unique trending YouTube videos. The dataset includes various details such as the video id, video title, channel title, publish time, trending date, tags, views, likes, dislikes, description, thumbnail link, and comment count, as explain in table 2.1. Using this dataset, the goal is to explore and analyse the factors that contribute to the popularity and engagement of trending YouTube videos. This may involve applying various data mining and machine learning techniques to uncover

patterns, relationships, and insights within the data (refer to Figure 2.2 for a sample of the dataset).

video_id	trending_date	title	channel_title	category_id	publish_time	views	likes	dislikes	comment_count	thumbnail_link	ratings_disabled	video_error	description
n1WpP7io	17.14.11	Eminem - Eminem	EminemVE	10	2017-11-11	17158579	787425	43420	125882	https://i.ytimg.c	FALSE	FALSE	Eminem's new track Walk on Water ft. Beyonc© is avail
0dBlkQ4M	17.14.11	PLUSH - B!DubbzTV		23	2017-11-11	1014651	127794	1688	13030	https://i.ytimg.c	FALSE	FALSE	STill got a lot of packages. Probably will last for another y
5q9jK5DgC	17.14.11	Racist Sup Rudy Manc		23	2017-11-11	3191434	146035	5339	8181	https://i.ytimg.c	FALSE	FALSE	WATCH MY PREVIOUS VIDEO & SUBSCRIBE &
d380meDc	17.14.11	I Dare You nigahiga		24	2017-11-11	2095828	132239	1989	17518	https://i.ytimg.c	FALSE	FALSE	I know it's been a while since we did this show, but we're
2Vv-BNVoq	17.14.11	Ed Sheera Ed Sheera		10	2017-11-04	33523622	1634130	21082	85067	https://i.ytimg.c	FALSE	FALSE	lectin' §.٥٤
0yIWz1XE	17.14.11	Jake Paul DramaAler		25	2017-11-11	1309699	103755	4613	12143	https://i.ytimg.c	FALSE	FALSE	Follow for News! - https://twitter.com/KEEMSTAR/in
uM5kFkh	17.14.11	Vanoss Sl VanossGa		23	2017-11-11	2987945	187464	9850	26629	https://i.ytimg.c	FALSE	FALSE	Vanoss Merch Shop - https://vanoss.3blackdot.com/in/in
2kyS6SvS	17.14.11	WE WANT CaseyNeis		22	2017-11-11	748374	57534	2967	16959	https://i.ytimg.c	FALSE	FALSE	SHANTELL'S CHANNEL - https://www.youtube.com/sha
JzCsM1vtr	17.14.11	THE LOG/ Logan Pau		24	2017-11-11	4477587	292837	4123	36391	https://i.ytimg.c	FALSE	FALSE	Join the movement. Be a Maverick & - https://ShopLogar
43sm-Qwl	17.14.11	Finally She Sheikh Mu		22	2017-11-11	505181	4135	976	1484	https://i.ytimg.c	FALSE	FALSE	Sheldon is roasting pastor of the church/young Sheldon

Figure 2.2: Sample of Dataset

Table 2.1 Illustration of The Dataset Features

Feature	Description
Video id	Refers to a unique identifier assigned to each individual video on the platform
Video title	The title or name of the video.
Channel title	The title or name of the youtube channel that uploaded the video.
Publish time	The date and time when the video were published on youtube.
Trending date	Refers to the date on which a particular video appeared on the list of trending videos
Tags	Keywords or tags associated with the video, which can provide additional information or context.
Views	The number of times the video has been viewed.
Likes	The number of likes received by the video.
Dis likes	The number of dislikes received by the video.
Description	A text description or summary of the video provided by the uploader.
Comment count	The number of comments posted on the video.
Category ID	A field indicating the category or genre to which the video belongs. This ID can vary depending on the region or classification system used.

2.4 Feature Extraction Technical

Feature extraction also known as feature engineering refers to the process of creating new features or variables from existing data to enhance the performance of a machine learning model. It involves transforming raw data into a format that captures relevant patterns or relationships, making it easier for the model to learn and make accurate predictions[24],[25].

There are several techniques and strategies for feature generation, depending on the nature of the data and the specific problem at hand. Feature extraction is particularly useful for unstructured or high-dimensional data types such as images, audio, text, or sensor readings. Here are some common techniques for feature extraction:

Visual feature extraction

Visual feature extraction is a crucial step in working with image data. Techniques for feature creation involve extracting meaningful visual attributes from images. In the context of image classification and object recognition[26], raw data often comprises pixel values that may not be directly compatible with certain classification algorithms. By extracting higher-level features, such as edges and regions correlated with the presence of human faces, the data can be transformed into a more suitable representation for a wider range of classification techniques[24]. These techniques encompass various approaches, including the utilization of Convolutional Neural Networks (CNN) to extract deep features, or the application of pre-trained models like VGG, ResNet, YOLO[27], Alex Net[28]. These methods enhance the ability to capture meaningful information from images, facilitating accurate classification and object recognition tasks.

YOLO For Object Detection

(YOLO) is a powerful real-time object detection method for deep learning and computer vision. YOLO is a technique for detecting and locating objects within image frames[29]. It treats object detection as a regression issue and calculates the class probabilities for each bounding box[30]. The base model detects images at an astounding 45 frames per second, while the Fast YOLO model detects at 155 frames per second, figure 2.3 shown the YOLO architecture. It is well-known for its precision and speed, and it is frequently employed in a variety of applications[29]. The main advantage of the YOLO model is its speed and accuracy. It can rapidly detect multiple objects in a single pass through the neural network, making it particularly well-suited for real-time object detection applications.

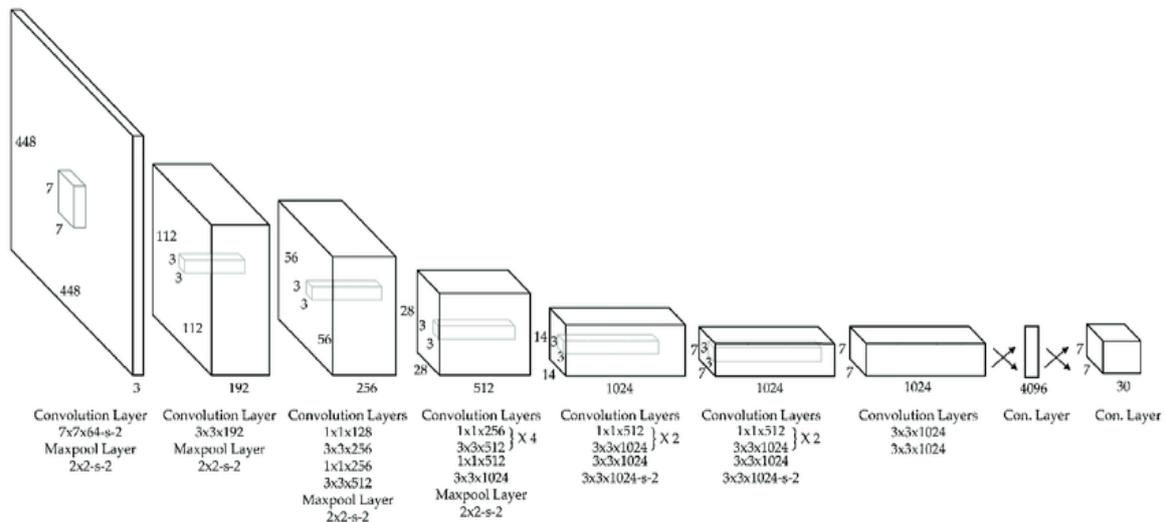


Figure 2.3: The YOLO Architecture[29]

YOLO works dividing the input image into a grid with fixed size ($S \times S$) with S set to 7 as illustrated in Figure 2.4. for each grid cell, it is predicting bounding boxes (rectangles) that might contain objects. These bounding boxes are described by their coordinates (center, width, and height) relative to the grid cell.

YOLO determines the attributes of these bounding boxes, where Y is the final vector representation for each bounding box[31], main steps of YOLO model in(Algorithm 2.1)[32].

$$Y = [P_r, b_x, b_y, b_h, b_w, c_i]$$

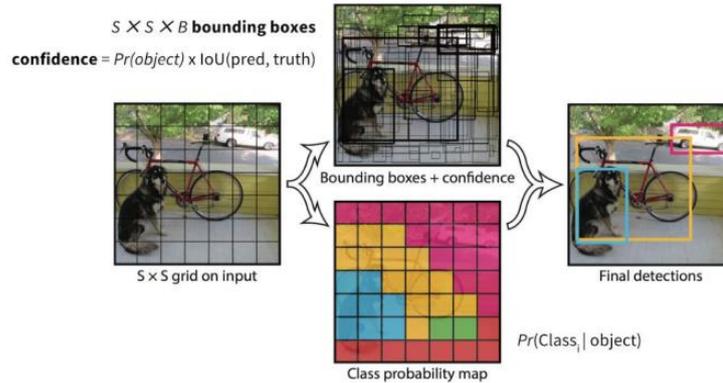


Figure 2.4: YOLO For Object Detection[31]

This is especially important during the training phase of the model.

P_r : corresponds to the probability score of the grid containing an object.

b_x, b_y : are the x and y coordinates of the center of the bounding box with respect to the enveloping grid cell.

b_h, b_w : correspond to the height and the width of the bounding box with respect to the enveloping grid cell.

c_i : correspond to the classes (number of classes).

Algorithm 2.1 Pre-Trained YOLO For Object Detection

Input: An image.
Output: Bounding boxes that represent detected objects along with their class probabilities.
Begin: -YOLO divides an input image into an $S \times S$ grid cell. -For each grid cell:

Calculate confidence scores for each bounding box.

These scores represent the model's confidence that the box contains an object and how accurate it thinks the predicted box is.

-During training, YOLO assign one bounding box predictor to be "responsive" for predicting an object based on which prediction has the highest Intersection over Union (IOU) with the ground truth

-Non-maximum suppression: NMS is a post-processing step that is used to improve the accuracy and efficiency of object detection.

NMS is used to identify and remove redundant or incorrect bounding boxes and to output a single bounding box for each object in the image.

End

2.5 Data Preparation Methods

Preprocessing in data mining refers to the various techniques and procedures used to clean, transform, and prepare raw data before it can be analysed[24]. The quality of the output from data mining algorithms largely depends on the quality of the input data. Therefore, preprocessing is essential to ensure that the data is in a suitable format for the analysis. Overall, the goal of preprocessing in data mining is to ensure that the data is ready for analysis and that the analysis results are accurate, reliable, and meaningful[33][34].

In this thesis, several methods were taken to prepare the dataset for the prediction model. Several important steps include:

2.5.1 Removal of Row Duplications

This method involves identifying and removing any duplicate rows in the dataset. Duplicates may occur due to various reasons such as data entry errors or system issues. Removing duplicates ensures that each data point is unique and avoids any bias or inconsistencies in the analysis[24].

2.5.2 Data Labelling[1]

Data labelling is the process of assigning labels or tags to data points in a dataset. It is typically performed by humans or automated systems to categorize the data according to predefined criteria. The labelled data is then used to train machine learning models, enabling them to learn and make predictions on new, unlabelled data. The classification was based on the equation 2.1 and 2.2 found in the papers.

$$score = (likes - 1.5 * dislikes) * comment_count / views \quad \dots \quad (2.1)$$

$$\text{Classification Rule, } y = \begin{cases} 0, & views < 100,000 \\ 1, & views \geq 100,000, score < 0 \\ 2, & views \geq 100,000, 0 \leq score < 300 \\ 3, & views \geq 100,000, score \geq 300 \end{cases} \quad \dots \quad (2.2)$$

Explanation for each variable:

y (Output Class): The output class assigned based on predefined conditions.

Score: A numerical value calculated using the specified equation.

Views: The number of views for a specific video.

Likes: The number of likes for a specific video.

Dislikes: The number of dislikes for a specific video.

Comment Count: The number of comments for a specific video.

2.5.3 Handling Missing Values

Handling missing data is a crucial aspect of data preprocessing. Choosing the right strategy for addressing missing data depends on factors such as the extent of the data, its nature, analysis goals, and the techniques used. It is important to consider the pros and cons of each strategy to select the best fit for your dataset and analysis[35]. There are several strategies available for handling missing values[24]:

1. Removing Data Objects: If a significant portion of the data objects have missing values, one option is to eliminate those objects entirely.
2. Removing Attributes: Similarly, the choice is to remove attributes with a high number of missing values.
3. Estimating Missing Values: In certain cases, missing values can be estimated or imputed using various techniques. Common methods include mean imputation, regression imputation, or using specialized machine learning algorithms for imputation[36][37].
4. Ignoring Missing Values: it may be possible to ignore missing values during the analysis[36]. Some algorithms can handle missing data by default, or they can be modified to accommodate missing values without imputing them.

2.5.4 Processing and Extracting Text Features

Cleaning textual features involves preprocessing steps to transform raw text data into a cleaner and more suitable format for analysis or modeling. Here are some common techniques used for cleaning textual features[38] 36):

1. Lowercasing Text: Converting all text to lowercase helps to standardize the text and minimize case sensitivity difficulties[40].
2. Removing of Punctuation: Getting rid of punctuation markings like commas, periods, and exclamation points might assist minimize noise in the text and improve future analysis.
3. Handling Special Characters and Numbers: special characters, URLs, or numbers that may not have significant value can be eliminated or substituted with appropriate placeholders[40].

In natural language processing (NLP), text feature creation involves extracting meaningful information from text data using various techniques such as bag-of-

words, term frequency-inverse document frequency (TF-IDF), word embeddings (e.g., Word2Vec or GloVe), sentiment analysis, or document embeddings (e.g., Doc2Vec). These techniques capture the semantic and contextual information present in the text[41].

Natural Language Processing (NLP) and sentiment analysis are closely related fields that both deal with the processing and analysis of human language. NLP is a multidisciplinary field that combines linguistics, computer science, and artificial intelligence to enable computers to understand, interpret, and generate human language[42]. Its primary goal is to bridge the gap between human language and machine language[43].

NLP techniques is a fairly generic term that covers a very wide range of applications[44], Figure 2.5 shows the most popular applications, allowing computers to perform tasks such as language translation, sentiment analysis, text summarization, question answering, part-of-speech tagging, named entity recognition, syntactic parsing, and machine translation[18].

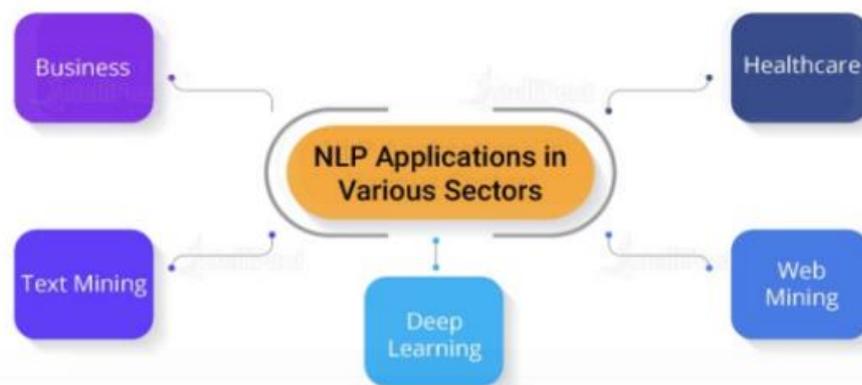


Figure 2.5: Application of NLP[45]

To accomplish these tasks, NLP algorithms rely on a variety of techniques, including statistical modeling, machine learning, deep learning, and rule-

based approaches. Among these applications, sentiment analysis, also known as opinion mining, has gained significant attention in recent years due to the increasing use of social media platforms. Sentiment analysis employs machine learning, data mining, natural language processing, and computational linguistics techniques to identify, extract, and analyze opinions and sentiments present in textual data[18]. Sentiment analysis is a Lexicon-based approach using AFFIN dictionary. its aim of to determine whether a given piece of text expresses a positive, negative, or neutral sentiment, figure 2.6 shows example Process of lexicon-based sentiment analysis[46].

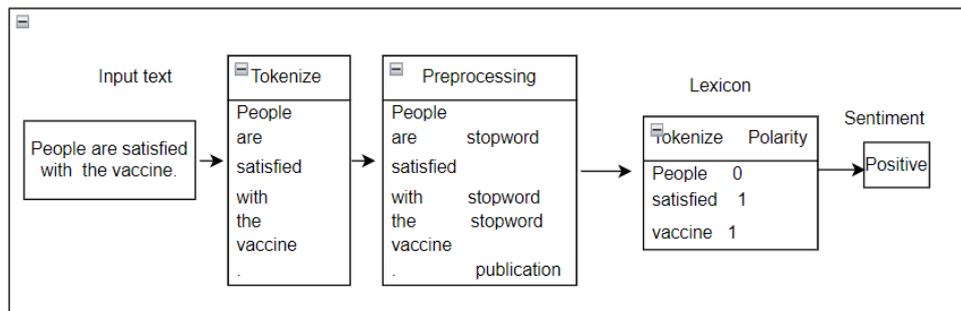


Figure 2.6: Example Process of Lexicon-Based Sentiment Analysis[46]

Instead of training a model on a labelled dataset, this approach relies on predefined sentiment lexicons or dictionaries containing words and their associated sentiment scores (positive, negative, or neutral), Algorithm 2.2 represent main steps in Lexicon-based approach.

Algorithm 2.2 Lexicon-Based Sentiment Analysis

Input: text
Output: class label of text (positive, negative, natural)
Begin: - Input text is usually tokenized into individual words or phrases (tokens).

- For each token, it searches for the corresponding word in the sentiment dictionary.
- The Sentiment Dictionary associates each word or phrase with a range of polarity scores (from -1 to 1), with 0 indicating neutral sentiment.
- Aggregating the individual polarity scores of all tokens in the text to calculate an overall polarity score for the entire text. This aggregation may involve simple averaging or other techniques depending on the specific implementation.
- The final polarity score is often normalized to ensure that it falls within the range of -1 to 1, even if the vocabulary contains scores on a different scale.
- Classify the overall sentiment score into predefined categories such as positive, negative, or neutral.
- Use predefined thresholds or rules for classification, e.g., score surpasses a threshold for positive classification, falls below for negative classification.

End

2.5.5 Encoding Features

It is an important aspect of data preprocessing, especially when working with machine learning algorithms that require numerical inputs[24]. There are several methods to handle non numerical features: Label Encoding, One-hot encoding, Target Encoding, Binary Encoding.

2.6 Data Mining

Data mining means extracting useful information from massive amounts of data[24]. It is a multidisciplinary field that combines elements of statistics, artificial intelligence (AI), and database research. It has become increasingly important due to the exponential growth in the size of datasets. By collecting and analysing large amounts of data, data mining aims to discover patterns and insights that may not be immediately apparent. Data mining has two main goals, shown in figure 2.7, they are prediction and description[47].

1. Prediction: Data mining aims to develop predictive models that can forecast future outcomes or behaviours based on historical data[24]. This is valuable for various applications such as sales forecasting, customer churn prediction, demand forecasting, and stock market prediction[48].
2. Description: Data mining also focuses on providing descriptive insights by summarizing and understanding the characteristics of the data. This involves identifying key attributes, summarizing statistical measures, visualizing data distributions, and generating reports or dashboards[24].

Both prediction and description are crucial in extracting meaningful insights and actionable knowledge from data, enabling informed decision-making and gaining a competitive advantage in different domains[49][8].

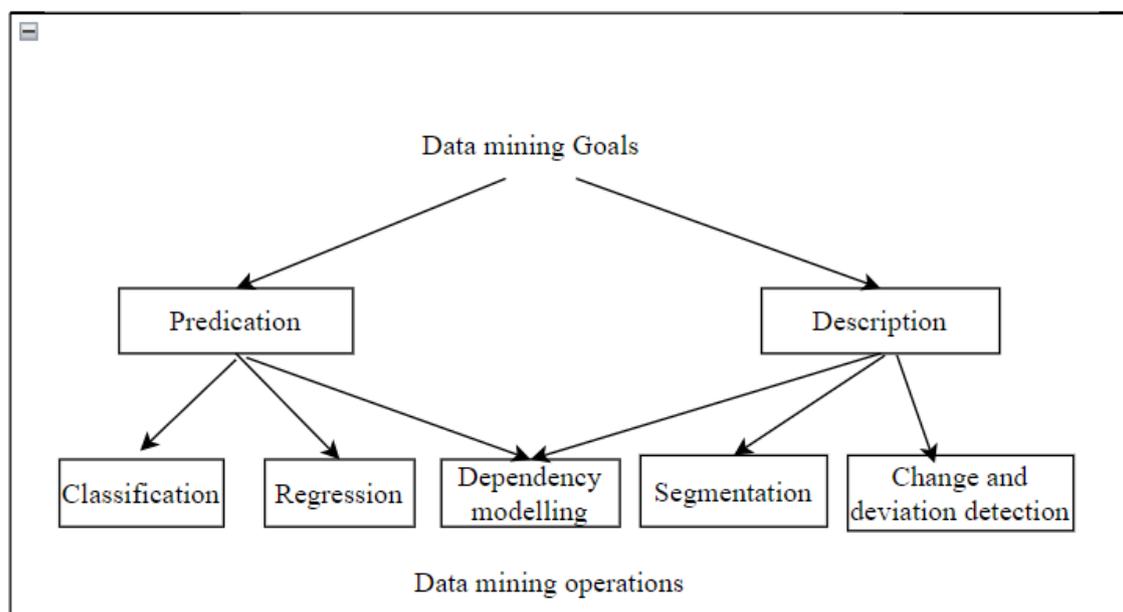


Figure 2.7: Data Mining Goals[47]

Machine Learning Techniques

Machine learning refers to a variety of techniques and approaches that allow computers and systems to learn from data and make predictions[42]. Here are a few examples of common machine learning techniques:

1-Supervised Learning: The set of models trained using a labeled dataset.

2-Unsupervised Learning: The set of models trained using an unlabeled dataset

Classification techniques are methods used in machine learning and statistics to classify data into predefined classes or categories. These techniques are supervised learning algorithms that learn from labelled training data to make predictions or assign class labels to new, unseen data[24]. Here are some used classification techniques in this thesis:

a. Support Vector Machine

Support vector machine (SVM) is a supervised machine learning algorithm that is versatile and can be applied to regression and classification tasks[50]. The basic idea of SVM (Support Vector Machines) is to find an optimal hyperplane that separates data points of different classes in a high-dimensional space, shows in figure 2.8. SVM can handle linear and non-linear problems and is commonly used for classification tasks. Here's types of SVMs and their corresponding equations[24][51]:

In linear SVM, the goal is to find a linear decision boundary. The equation of the hyperplane can be represented as[52]:

$$w^t \cdot x + b = 0 \quad \dots \quad (2.3)$$

w^t represents the weight vector perpendicular to the hyperplane

x represents the input data vector

b is the bias term.

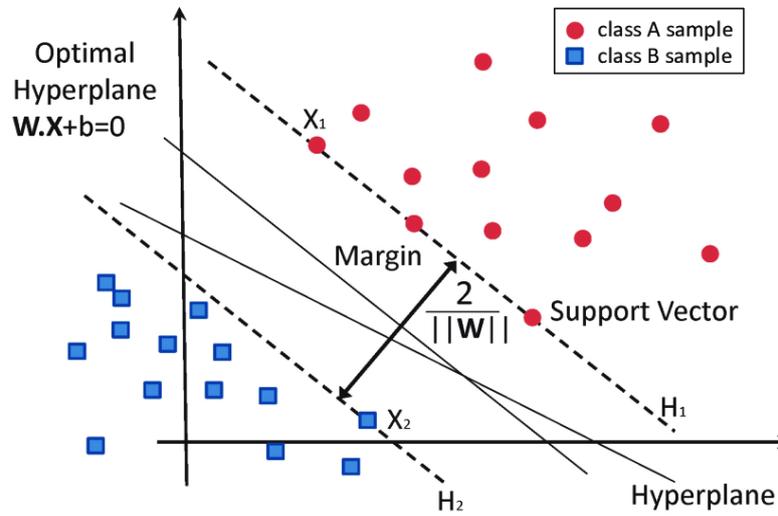


Figure 2.8: Classification of Data by Support Vector Machine

Non-linear SVM, also known as kernel SVM, is used to handle non-linearly separable data by mapping it to a higher-dimensional feature space using a kernel function. The decision boundary equation becomes:

$$\sum(\alpha_i \cdot y_i \cdot k(x_i, x)) + b = 0 \quad \dots \quad (2.4)$$

where α_i and y_i are the Lagrange multipliers and class labels of the support vectors, x_i represents the support vectors, $k(x_i, x)$ is the kernel function that computes the similarity between the support vectors and the input data point x . Commonly used kernel functions include: Gaussian (Radial Basis Function), Polynomial, Sigmoid and Linear Kernel. The choice of kernel function depends on the problem and the underlying data. These equations represent the decision boundaries for different types of SVMs. The goal of SVM is to find the optimal values for the parameters (weights, bias, Lagrange multipliers) that maximize the margin between the classes or minimize classification errors, depending on the problem setup.

b. K-Nearest Neighbor Classifier

The K-Nearest Neighbor (KNN) Classifier is a widely employed technique in machine learning for both classification and regression tasks. It is a non-parametric and instance-based algorithm, meaning it does not assume any particular data distribution and relies on the actual training instances for predictions. The fundamental concept of KNN involves classifying a new data point based on the majority class among its K nearest neighbors [28]. To determine the neighbors, KNN calculates the distance between the new data point and each training data point[24]. Common distance metrics used include Euclidean distance and Manhattan distance. Euclidean distance, for example, calculates the distance between two points (x_1, y_1) and (x_2, y_2) in a 2D space as follows:

$$\text{Distance} = \text{sqrt}((x_2 - x_1)^2 + (y_2 - y_1)^2) \quad (2.5)$$

Once the distances are calculated, the subsequent stage involves determining the K nearest neighbors of the new data point. The value of K , representing the number of neighbors to consider, is specified by the user. The KNN algorithm then employs a voting mechanism to assign a class label to the new data point. In classification scenarios, the predicted class is determined by selecting the majority class among the K neighbors. When K equals 1, the class label of the nearest neighbor is directly assigned to the new data point. For regression tasks, the predicted value is obtained as the average or weighted average of the values associated with the K nearest neighbors. In the k -nearest neighbors (KNN) algorithm, the choice of the right value for the parameter k is crucial. The value of k determines the number of nearest neighbors that will be considered when making predictions for a new or test instance. Choosing the right value for k in K-Nearest Neighbors (KNN) is important. If k is too small, the classifier may

overfit the training data and be sensitive to noise. If k is too large, the classifier may misclassify test instances by including distant neighbors. It's crucial to find the optimal k value to balance between overfitting and misclassification, ensuring accurate predictions (see Figure 2.9).

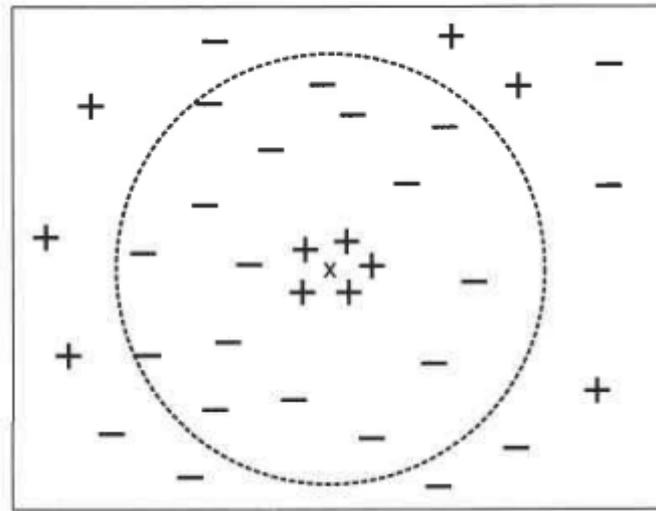


Figure 2.9: KNN Classification with Large K

c. Random Forest Classifier

Random Forest is an ensemble method that combines multiple decision trees to make predictions. It creates a forest of decision trees rather than relying on a single decision tree to improve prediction accuracy. Each decision tree in the Random Forest is trained independently on a random subset of the training data. This random subset is typically selected through bootstrapping, where data points are sampled with replacement[24]. This process introduces diversity among the individual trees. During prediction, the algorithm aggregates the predictions of all the individual decision trees to make the final prediction. The way this aggregation is done depends on the type of task (classification or regression). In classification tasks, the predictions of individual decision trees are combined using majority voting. The class that receives the most votes across

the decision trees becomes the final prediction. In regression tasks, the predictions of individual decision trees are often averaged or aggregated to obtain the final prediction[53]. The main objective of using a Random Forest is to reduce variance and bias, thus improving the generalization ability of the model. Unlike a single decision tree, which can have low bias but high variance (prone to overfitting), the Random Forest's ensemble approach helps to mitigate overfitting by reducing variance, leading to better overall performance on unseen data.

Overall, Random Forest is a powerful and widely used machine learning algorithm known for its ability to handle complex datasets, reduce overfitting, and provide reliable predictions. The main idea of the random forest algorithm is presented in Figure (2.10).

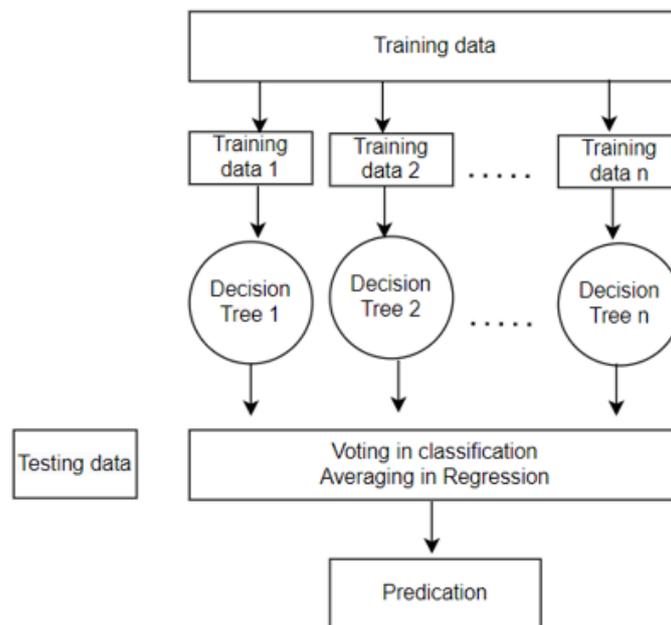


Figure 2.10: The Idea of Random Forest

The impurity measures for regression and classification tasks in the Random Forest algorithm are different. In regression tasks, the Random Forest algorithm

typically uses the mean squared error (MSE) or the mean absolute error (MAE) to measure impurity in equations (2.6),(2.7) In classification tasks, the Random Forest algorithm uses Gini impurity or entropy to measure impurity(2.8),(2.9),whereat the number of features in each tree is identified according to the equations (2.10) [24], Random Forest work steps[54] is presented in Algorithm 2.3:

For regression:

$$\text{variance, MAE} = \frac{1}{N} \sum_{i=1}^N |y_i - \mu| \quad \dots \quad (2.6)$$

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \mu)^2 \quad \dots \quad (2.7)$$

Where y_i is the actual label (target) for the i instance.

N is the number of instances in the current node.

μ = is the mean of the labels (targets) in the current node, given by $\frac{1}{N} \sum_{i=1}^N y_i$

For classification:

$$\text{Gini} = 1 - \sum_i f_i^2 \quad \dots \quad (2.8)$$

$$\text{entropy} = - \sum_i f_i \log_2 \cdot f_i \quad \dots \quad (2.9)$$

Where:

f_i is the frequency (proportion) of instances with class 'i' at the current node.

\sum_i is the sum over all classes.

$$F = \log_2 f + 1 \quad \dots \quad (2.10)$$

f is the number of input features.

Algorithm 2.3 Random Forest[24]

Input: <ul style="list-style-type: none">- Training data- Number of trees in the forest (num_trees)- Stopping criteria for tree growth (e.g., max depth, min-samples-leaf)
Output: Trained Random Forest model
Begin: <ul style="list-style-type: none">-For each tree in the forest (from 1 to number of trees):<ul style="list-style-type: none">Create a bootstrap sample (a random subset with replacement) from the training data.Randomly select a subset of features for this tree.Build a decision tree using the bootstrap sample and selected features:<ul style="list-style-type: none">If stopping criteria are met or only one class remains, make this node a leaf and assign the majority class.Otherwise, choose the best feature and split the data into child nodes.Repeat the splitting process for child nodes until stopping criteria are met.-Store the decision tree in the forest.-To make a prediction for new data (test data):<ul style="list-style-type: none">-For each tree in the forest:<ul style="list-style-type: none">Traverse the tree based on the input features (test data) to reach a leaf node.Record the class (for classification)-Aggregate the predictions use majority voting among tree predictions.Return the final prediction.
End

d. Extreme Gradient Boosting

XGBoost (Extreme Gradient Boosting) is a highly popular and powerful implementation of gradient boosting machines (GBM) used for supervised learning tasks. Gradient boosting is an ensemble learning technique that combines the predictions of multiple weak learners (usually decision trees) to

create a strong predictive model. XGBoost is known for its exceptional performance in various machine learning competitions and real-world applications. It has become a preferred choice for data scientists and machine learning practitioners [55].

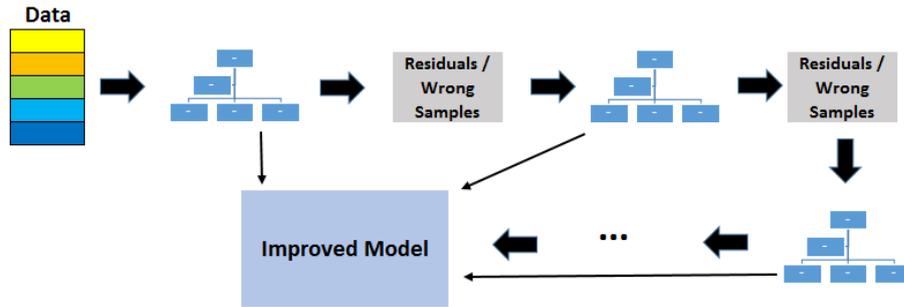


Figure 2.11: Boosting Method[56]

XGBoost combination of high performance, regularization, feature importance analysis, and scalability has made it a popular choice for various machine learning tasks, including classification, regression, ranking, and recommendation systems, among others. Its widespread adoption and continued development in the data science community contribute to its reputation as one of the best performing algorithms for supervised learning. Figure (2.11) shows the boosting method.

The mathematical equations for the XGBoost algorithm involve the objective function, regularization terms, and the optimization process during training[57][52].

The objective function of XGBoost can be written as:

$$\bar{O}_i = \phi(x_j) = \sum_{k=1}^K f_k(x_i), f_k \in y \quad \dots \quad (2.11)$$

$$L(\phi) = \sum_i L(y_i^{\check{v}}, y_i) + \sum_k \Omega(f_k) + \gamma(K) \quad \dots \quad (2.12)$$

To represent the XGBoost algorithm mathematically, we'll use the following notations:

X represents the input features (matrix) of the dataset with n examples and m features.

y_i represents the target output (vector) of the dataset with n examples.

T_k represents the k_{th} tree in the ensemble.

K represents the number of trees in the XGBoost ensemble.

F_k represents the k_{th} tree in the ensemble.

γ is the tree complexity term.

Ω is the regularization term.

L is the loss function.

2.7 Evaluation of the Prediction Models

Evaluating the performance of classification prediction models is essential to assess their effectiveness and make informed decisions. Here are some commonly used evaluation metrics for classification models[24]:

Tp : True positive Tn : True negative

Fn : False negative Fp : False positive

Accuracy: The ratio of the correctly classified instances to the total number of instances [41].

$$Accuracy = \frac{Tp+Tn}{Tp+Tn+Fp+Fn} \quad \dots \quad (2.13)$$

Precision: Precision measures the proportion of true positive predictions (correctly predicted positive instances) out of all positive predictions (true positives + false positives)[41].

$$\textit{Precision} = \frac{Tp}{Tp + Fp} \quad \dots \quad (2.14)$$

Recall (Sensitivity or True Positive Rate): Recall measures the proportion of true positive predictions out of all actual positive instances[41].

$$\textit{Recall}, r = \frac{Tp + Fn}{Tp} \quad \dots \quad (2.15)$$

F1-Score: The F1-score is the harmonic mean of precision and recall and provides a balanced measure that considers both precision and recall. [41].

$$\textit{F1 score} = 2 * \frac{\textit{precision} + \textit{recall}}{\textit{precision} + \textit{recall}} \quad \dots \quad (2.16)$$

Chapter Three

The Proposed System

CHAPTER THREE

The Proposed System

3.1 Overview

This chapter describes the main stages of the proposed system methodology. It begins by presenting the architecture of the proposed system. It covers feature extraction and data preparation, including an explanation of preprocessing techniques. Lastly, the chapter demonstrates the utilization of classifiers and their evaluation to achieve the key aim of this thesis.

3.2 The Proposed System Architecture

The proposed system architecture consists of several main stages, each comprising sub-stages designed to achieve the study objectives. These stages include visual feature extraction, dataset preparation, machine learning models, and their evaluation, as illustrated in Figure 3.1.

To predict popular videos, we utilize a dataset obtained from Kaggle. The first stage of the proposed system focuses on visual feature extraction using pre-trained deep learning techniques such as YOLO v5, specifically designed for object extraction to enrich our dataset.

The second stage involves data preparation, which includes labelling the dataset, processing missing values, and analysing text features using sentiment analysis and interval feature extraction as sub-stages, data preparation encompasses several steps aimed at effectively preparing the input for the system.

In the final stages, a variety of machine learning model are implemented, including K-Nearest Neighbor, Random Forest, Support Vector Machine, Gradient Boosting, and Extreme Gradient Boosting. Subsequently, it evaluated the outcomes of the proposed model using various methods to estimate predictive errors.

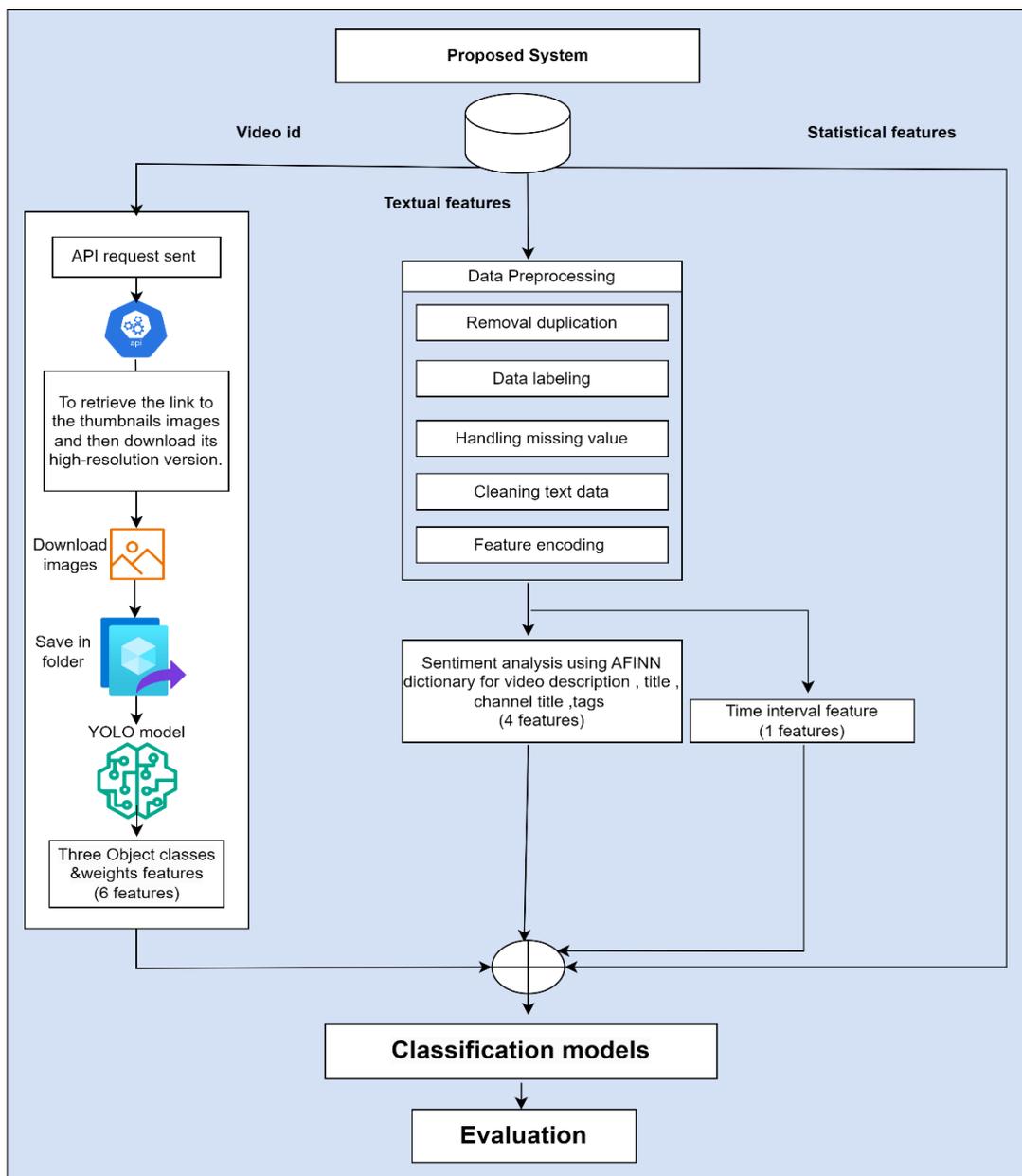


Figure 3.1: The Proposed System Architecture


```

API requests to the YouTube Data API to retrieve video Thumbnail.
  If video Thumbnail is found in the API response:
    It downloads the thumbnails to file.
  Else //no video Thumbnails is found
    It attempts to extract the default thumbnail URL by parsing the video's watch
page HTML using Beautiful Soup
    Downloads the default thumbnail.
  // End if
-Save the image in a folder.
//End for
End

```

The image processing begins by loading the pre-trained deep learning model (YOLO V5) and its weight file, followed by setting initial values for specific parameters and resizing the image as a preprocessing step. Performing object detection involves passing the image through the YOLO model to obtain predictions, Algorithm 3.2 main steps to extract features from images.

Algorithm 3.2 Extract Features from Images

Input: Thumbnail images
Output: set of features (Three largest objects and each object weight) in csv file.
<p>Begin:</p> <ul style="list-style-type: none"> -Loading a pre-trained YOLOv5 model from specified weights and setting NMS parameters. -Iterating over images in a folder. -Using YOLOv5 to detect objects in the pre-processed image, obtaining attributes like bounding boxes and categories. -Filtering the detected objects and selecting the three largest based on object size, assigning weights based on object size. <p>Calculate the weights using the following formula:</p> $\text{Weights} = \text{object size (bounding box)} / \text{size of image}$

-Storing selected objects' categories, weights, and image IDs in CSV file for further analysis or use.

End

Subsequently, Non-Maximum Suppression is applied to remove redundant bounding box predictions. Following this, process the output predictions by extracting class labels. Then filter the predictions based on a confidence threshold or other criteria.

In the next step, select the three largest objects based on their size. Weights are assigned to these objects according to their size, and these values are stored in a CSV file for further analysis or use (see Algorithm 3.2). The weight formula implies that the weight assigned to an object is determined by its size relative to the overall size of the image. The result represents a measure of the significance or contribution of the object within the context of the entire image. Addressing missing values in the new features involves a proactive approach. When encountering instances in dataset with incomplete information, objective is to avoid leaving these gaps unattended. Instead, systematically examine the entire dataset to identify other instances resembling those with missing data.

This is pretty important because information from those similar cases is allowed to be used. It's kind of like help being sought from friends when some details are missing; the answers might be with them. When those similar instances are picked to help us out, the same ones are not always gone for. They are mixed up, and they are chosen randomly. Why? Due It ensures that our new information comes from all over the dataset, not just one specific place. This keeps our data diverse and prevents relying too much on just one

source; it helps reduce any biases that might sneak in if the same source is always picked for filling in the gaps. It makes our data tougher and more reliable for whatever is wanted to be done with it. In the final step, these features are combined with the original data using the video ID as additional features.

3.2.2 Data Preprocessing

This phase encompasses various steps aimed at preparing the data for optimal compatibility with machine learning algorithms. Typically, these steps involve:

3.2.2.1 Removal Duplication

Duplicate in the dataset arises from the collection process being conducted at different times. Eliminating duplicate entries becomes a critical step to ensure that each video is represented only once. This process aids in averting redundancy and potential bias in subsequent analyses. By deleting duplicate records based on unique video IDs, the dataset is refined, resulting in a more streamlined and accurate representation. This refinement enhances the dataset's appropriateness for future analyses and modelling.

3.2.2.2 Dataset Labelling

Data labeling, as discussed in Chapter 2, Section 2.5.2, involves assigning labels to instances in a dataset for classification purposes. Algorithm 3.3 provides a comprehensive explanation of this labeling process.

Algorithm 3.3 Labelling Video

Input:

Number of likes (likes)

Number of dislikes (dislikes)

Number of comments (comment_count)
Number of views (views)
Output: label (0,1,2,3)
<p>Begin</p> <p>-Calculate the score using the following formula:</p> $score = (likes - 1.5 * dislikes) * comment_count / views$ <p>-Determine Label:</p> <p>If views are less than 100,000:</p> <p style="padding-left: 40px;">Set label to 0.</p> <p>Else</p> <p style="padding-left: 40px;">If score is less than 0:</p> <p style="padding-left: 80px;">Set label to 1.</p> <p style="padding-left: 40px;">Else</p> <p style="padding-left: 80px;">if $0 < score < 300$:</p> <p style="padding-left: 120px;">Set label to 2.</p> <p style="padding-left: 40px;">Else</p> <p style="padding-left: 80px;">Set label to 3.</p> <p>-Return label</p> <p>End</p>

3.2.2.3 Handling Missing Value of Video Description

There are very few missing values in one of the features used, specifically the video description. The method employed to handle these missing values is imputation. The imputation methods used include replacing missing values with a constant value, in this case, 'N/A.' This allows the dataset to indicate that description information for these videos is not available or has not been provided, and the specific steps are outlined in Algorithm 3.4.

Overall, these imputation methods help maintain the integrity of the dataset by properly handling missing information and enabling subsequent analyses and modelling to be conducted accurately.

Remove words containing digits and individual digits: This step uses regular expressions to remove words that contain digits or individual digits from the text.

Split the text using the provided delimiters: This step splits the text into a list of words using specified delimiters: ',', ':', '.', '>', '<', and whitespace. It converts the text into individual words, which will be further processed.

Remove empty elements from the list: This step removes any empty elements from the list of words.

Merage the cleaned words back into a single string: This step joins the cleaned words from the list back into a single string using a space (' ') as the separator.

3.2.2.5 Features Encoding

After cleaning text data and handling missing values, the next step transforming non numerical features into numerical features for ease of use in data analysis and machine learning algorithms. In the used dataset, the 'comment disable' feature can be transformed into a binary representation where 'True' indicates that comments are disabled, and it is replaced with '1,' while 'False' indicates that comments are enabled and is replaced with '0.'

The steps to perform this coding are outlined in Algorithm 3.5, which provides a systematic approach to convert the non-numerical feature into the desired numerical representation.

Algorithm 3.5 Convert Non numerical Feature into numerical Feature

Input: Non numerical feature (True, False)
Output: Numerical feature (1,0)
Begin:

- Identify the non-numerical feature in the dataset that needs to be converted into a numerical feature.
- For each value in non-numerical: // Iterate over the instances
- If value is True, replace it with 1.
- Else , replace it with 0.

End

3.2.2.6 Processing Text Data Using Sentiment Analysis

When choosing the text features to work with, such as the video description, tags, video title, and channel title, these texts undergo a bit of cleanup and preparation.

Once the text is cleaned, it is broken down into individual words. A sentiment dictionary, like the AFINN Dictionary, is used to rate words in English as positive, negative, or neutral. These ratings range from super negative to very positive, with zero meaning totally neutral.

After assigning sentiment scores to each word in the text based on the AFINN sentiment list, calculate an overall sentiment score for the entire text. This can be done by adding up all the individual scores or by determining a weighted average.

Finally, categorize the overall sentiment score into different sentiment categories, such as positive, negative, or neutral. This classification depends on specific thresholds or rules that have been set. For example, if the score goes above a certain threshold, call it positive; if it drops below, consider it negative.

3.2.2.7 Time Interval Features Extraction

In the context of video analysis, the feature " time interval " captures the duration it took for a video to attract significant attention or engagement from users. The time interval is calculated by subtracting the publication date from the trending date, as denoted in algorithm 3.6. This calculation yields a time duration, typically measured in hours, days, or weeks, which represents how long it took for the video to gather significant attention or become trending.

This calculated time interval is then incorporated as feature in the dataset, enabling its use alongside other features to analyse patterns, correlations, or predictive relationships. By utilizing the " time interval " analysts can gain insights into the time dynamics and popularity growth of videos. It becomes a valuable tool for understanding the factors that contribute to a video's success and for predicting future trends or user engagement patterns.

Algorithm 3.6 Time Interval Feature Extraction

Input: published date, trending date.
Output: time interval value in days
<p>Begin:</p> <ul style="list-style-type: none"> -Convert publish date to suitable format. - Convert trending date to suitable format. -For each instance: <ul style="list-style-type: none"> Subtract the trending date value from the corresponding publish date value. Store the result in the time interval feature - Return time interval features and add them to the dataset <p>End</p>

3.2.3 Machine Learning Models

Machine learning models were used to predict the popularity of YouTube videos by classifying them into four categories. To do this, firstly a set of features are extracted from the thumbnail video images using the pre-trained Yolo model.

Additionally, textual features were extracted and processed, along with the time interval feature, and the statistical features that are originally present in the data set are also used, such as the number of likes, comments, dislikes, video disabling, and video category. The total number of features used is 16. After processing the data, an appropriate classifier was applied, and predictions were made in multiple stages, depending on the features used.

The total dataset was divided into a training set (70%) and a testing set (30%). Four classifiers were applied to find the best, these are:

- a. Support Vector Machine
- b. K-Nearest Neighbor Classifier
- c. Random Forest Classifier
- d. Extreme Gradient Boosting Classifier

After applying the four classifiers, the Random Forest and XGBoost algorithms demonstrated superior performance, achieving the highest accuracy and outperforming other classifiers.

3.2.4 Evaluating the Performance of the Proposed Model

Different separate measures were employed to assess the effectiveness of a specific categorization system. Accuracy, f1-score, precision, and recall are all included. These measurements are calculated using a confusion matrix, which is a matrix that represents the number of cases that are correctly or incorrectly

predicted by a classification model. In this thesis, the predicted classes are four, so the form of the confusion matrix would be as follows (see Figure 3.3)

	Predicted Class				
		0	1	2	3
Actual Class	0	TP	FP	FP	FP
	1	FN	TP	FP	FP
	2	FN	FN	TP	FP
	3	FN	FN	FN	TP

Figure 3.3: The Confusion Matrix of a Four Class

Chapter Four

Result and Discussion

CHAPTER FOUR

Results and Discussion

4.1 Introduction

This chapter presents the outcomes of the various stages of the proposed system, as discussed in Chapter three. It provides a comprehensive overview of the results obtained from each stage, including feature extraction from images, text cleaning and processing, and the performance of the machine learning models employed.

4.2 Research Environment

Hardware: Processor Intel i7, RAM 16GB.

Operating System: Windows10 (64) bit.

The system was implemented using Python 3.11 and the PyCharm.

4.3 Proposed System Results

In this section, after applying the selected classifiers to the dataset, the proposed system exhibited promising outcomes at each stage, showcasing its effectiveness in predicting YouTube video popularity.

4.3.1 Extract Visual Features from Images

After using pre-trained YOLOv5 to process images and extract objects, as in figure 4.1.



Figure 4.1: The Result of Using Pre-Trained Yolo V5 Model

The result of this step is a group of objects and related information. The three largest objects are selected based on their size in the image (boundary box), and a weight is assigned to each object depending on its size in the image. Each object is represented by a number belonging to 80 classes with a weight for each object, as shown in the Table 4.3.

Table 4.1: Feature Set Extracted from Images (Before Processing Step)

video_id	ObjectCategory 1	ObjectCategory 2	ObjectCategory 3	ObjectWeight 1	ObjectWeight 2	ObjectWeight 3
0dBIKQ4Mz1M	0	0	nan	0.2400	0.3793	nan
d380meDoWoM	0	0	0	0.517	0.296	0.186
2Vv-BfVog4g	16	74	nan	0.945	0.054	nan
0yIWz1XEeyc	0	0	0	0.5519	0.2490	0.1989
6ylqz1qWeyu	16	32	4	0.329	0.094	0.502

There are missing values in the extracted features, likely due to the absence of three objects in some images, as previously discussed in Section 3.2.1 Which were processed, (Table 4.1: illustrates the data before processing step).

When there are missing values in object category 2 or object category 3, the dataset will be searched for videos that are similar to object category 1. From the found videos, one will be randomly selected to fill the missing values in object category 2 or object category 3, along with their respective weights. and the (Table 4.2: shows the results after processing the empty values).

Table 4.2: Feature Set Extracted from Images (After Processing Step)

video_id	ObjectCategory 1	ObjectCategory 2	ObjectCategory 3	ObjectWeight 1	ObjectWeight 2	ObjectWeight 3
0dBIKQ4Mz1M	0	0	0	0.2400	0.3793	0.1989
d380meDoVWoM	0	0	0	0.517	0.296	0.186
2Vv-BfVoq4g	16	74	4	0.945	0.054	0.502
0yIWz1XEeyc	0	0	0	0.5519	0.2490	0.1989
6ylqz1qWeyu	16	32	0	0.329	0.094	0.502

Subsequently, the extracted features are integrated with the original features using the Image ID, which corresponds to the Video ID in the original dataset. This integration is a crucial step as it allows for the combination of both sets of features for further analysis and modeling.

It is important to note that each value in the first three columns represent a category(classes), with the number 0 denoting a person and the number 27 representing a tie (To find out more, see Table 4.3), the other three columns are the weights for each object, its importance according to its size in the image, as indicated in Table 4.1.

Table 4.3 Shows the Number and What It Represents (Category)

Number	Class2	Number	Class
1	person	41	wine glass
2	bicycle	42	cup
3	car	43	fork
4	motorbike	44	knife
5	aeroplane	45	spoon
6	bus	46	bowl
7	train	47	banana
8	truck	48	apple
9	boat	49	sandwich
10	traffic light	50	orange
11	fire hydrant	51	broccoli
12	stop sign	52	carrot
13	parking meter	53	hot dog
14	bench	54	pizza
15	bird	55	donut
16	cat	56	cake
17	dog	57	chair
18	horse	58	sofa
19	sheep	59	pottedplant
20	cow	60	bed
21	elephant	61	diningtable
22	bear	62	toilet
23	zebra	63	tvmonitor
24	giraffe	64	laptop
25	backpack	65	mouse
26	umbrella	66	remote
27	handbag	67	keyboard
28	tie	68	cell phone
29	suitcase	69	microwave
30	frisbee	70	oven
31	skis	71	toaster
32	snowboard	72	sink
33	sports ball	73	refrigerator
34	kite	74	book
35	baseball bat	75	clock
36	baseball glove	76	vase
37	skateboard	77	scissors
38	surfboard	78	teddy bear
39	tennis racket	79	hair drier
40	bottle	80	toothbrush

4.3.2 Data Preparation Results

The results of the data preparation phase are pivotal to the success of the proposed system, being influenced by the quality and reliability of subsequent predictions. In this section, the results of the data preparation process are delved into, with key observations and decisions that shaped the dataset for optimal model training and evaluation being highlighted.

4.3.2.1 Dataset Labelling

The result of labeling the data into four classes, as mentioned in the proposed system, explain in table 4.4.

Table 4.4 The Result of Labelling the Data into Four Classes

<u>video_id</u>	<u>trending_date</u>	<u>category_id</u>	<u>.....</u>	<u>likes</u>	<u>Label</u>
n1WpP7iowLc	17.14.11	23	127794	1
0dBikQ4Mz1M	17.14.11	23	146035	1
5qpjK5DgCt4	17.14.11	24	132239	1
d380meD0W0M	17.14.11	10	1634130	1

At this stage of pre-processing, duplicate records in the video dataset have been successfully removed based on duplicate video IDs. The dataset now contains only unique records (see table 4.5,4,6), ensuring that each video is represented by a single entry without any redundant observations.

Table 4.5: The Results Before of Removing Duplicate Records

<u>video_id</u>	<u>trending_date</u>	<u>category_id</u>	<u>.....</u>	<u>likes</u>
n1WpP7iowLc	17.14.11	23	127794
0dBikQ4Mz1M	17.14.11	23	146035
5qpjK5DgCt4	17.14.11	24	132239
0dBikQ4Mz1M	17.14.11	23	146035
d380meD0W0M	17.14.11	10	1634130
n1WpP7iowLc	17.14.11	23	127794

Table 4.6: The Result After of Removing Duplicate Records

<u>video_id</u>	<u>trending_date</u>	<u>category_id</u>	<u>.....</u>	<u>likes</u>
n1WpP7iowLc	17.14.11	23	127794
0dBikQ4Mz1M	17.14.11	23	146035
5qpjK5DgCt4	17.14.11	24	132239
d380meD0W0M	17.14.11	10	1634130

4.3.2.2 Missing Value in Video Description

For some videos, the description field contains missing values, which are processed by assigning the value 'Not available' This allows the dataset to indicate that the description information for those videos is not available or not provided. The table (4.7,4.8) shows the results of this step.

Table 4.7: Description Feature Before Handling the Missing Values

<u>video_id</u>	<u>description</u>
n1WpP7iowL	Eminem's new track Walk on Water ft. Beyoncé i...
0dBikQ4Mz1	nan
5qpjK5DgCt4	I know it's been a while since we did this sho...
8HNuRNi8t70	nan

Table 4.8: Description Feature After Handling the Missing Values

<u>video_id</u>	<u>description</u>
n1WpP7iowL	Eminem's new track Walk on Water ft. Beyoncé i...
0dBikQ4Mz1	Not available
5qpjK5DgCt4	I know it's been a while since we did this sho...
8HNuRNi8t70	Not available

4.3.2.3 Cleaning text features

The results of the pre-processing stage are shown, where textual features such as video descriptions, video titles, and video channel titles are cleaned by removing irrelevant information. Table 4.9 shows the results of the pre-processing steps for the textual features within its dataset.

Table 4.9: The Results of Each Step on the Sample Video Description

Original text	<p>Ending Explained for the latest from master Guillermo Del Toro, the moving romantic monster movie2 THE SHAPE OF WATER starring Sally Hawkins and Doug Jones. Plus, analyzing the films bigger meaning and themes.</p> <p>Subscribe! ▶▶▶ http://bit.ly/2jrstgM</p> <p>Support FoundFlix on Patreon! ▶▶▶</p> <p>http://www.patreon.com/foundflix</p> <p>=== Connect with us on Social Media! ===</p> <p>FACEBOOK ▶▶▶ www.facebook.com/foundflix</p> <p>TWITTER ▶▶▶ www.twitter.com/foundflix</p> <p>INSTAGRAM ▶▶▶ www.instagram.com/foundflix</p>
Step 1: After Convert to lowercase	<p>ending explained for the latest from master guillermo del toro, the moving romantic monster movie2 the shape of water starring sally hawkins and doug jones. plus, analyzing the films bigger meaning and themes.</p> <p>subscribe! ▶▶▶ http://bit.ly/2jrstgm</p> <p>support foundflix on patreon! ▶▶▶</p> <p>http://www.patreon.com/foundflix</p> <p>=== connect with us on social media! ===</p>

	<p>facebook ►► www.facebook.com/foundflix</p> <p>twitter ►► www.twitter.com/foundflix</p> <p>instagram ►► www.instagram.com/foundflix</p>
Step 2: After Remove URLs from the text	<p>ending explained for the latest from master guillermo del toro, the moving romantic monster movie2 the shape of water starring sally hawkins and doug jones. plus, analyzing the films bigger meaning and themes. subscribe! ►►</p> <p>support foundflix on patreon! ►►</p> <p>=== connect with us on social media! ===</p> <p>facebook ►►</p> <p>twitter ►►</p> <p>instagram ►►</p>
Step 3: After Remove punctuation marks	<p>ending explained for the latest from master guillermo del toro the moving romantic monster movie2 the shape of water starring sally hawkins and doug jones plus analyzing the films bigger meaning and themes subscribe ►►</p> <p>support foundflix on patreon ►►</p> <p>connect with us on social media</p> <p>facebook ►►</p> <p>twitter ►►</p> <p>instagram ►►</p>
Step 4: After Remove special characters,	<p>ending explained for the latest from master guillermo del toro the moving romantic monster movie2 the shape of water starring sally hawkins and doug jones plus analyzing the films bigger meaning and themes</p>

<p>numbers, and extra spaces</p>	<p>subscribe support foundflix on patreon connect with us on social media facebook twitter instagram</p>
<p>Step 5: After remove words containing digits and individual digits</p>	<p>ending explained for the latest from master guillermo del toro the moving romantic monster movie the shape of water starring sally hawkins and doug jones plus analyzing the films bigger meaning and themes subscribe support foundflix on patreon connect with us on social media facebook twitter instagram</p>
<p>Step 6: After split the text using the provided delimiters</p>	<p>[', 'ending', 'explained', 'for', 'the', 'latest', 'from', 'master', 'guillermo', 'del', 'toro', ", 'the', 'moving', 'romantic', 'monster', 'movie', 'the', 'shape', 'of', 'water', 'starring', 'sally', 'hawkins', 'and', 'doug', 'jones', ", 'plus', ", 'analyzing', 'the', 'films', 'bigger', 'meaning', 'and', 'themes', ", ", 'subscribe', ", ", ", 'support', 'foundflix', 'on', 'patreon', ", ", ", ", ", ", ", 'connect', 'with', 'us', 'on', 'social', 'media', ", ", ", ", ", ", 'facebook', ", ", 'twitter', ", ", 'instagram', ", ", "]</p>

Step 7: After remove empty element in list	['ending', 'explained', 'for', 'the', 'latest', 'from', 'master', 'guillermo', 'del', 'toro', 'the', 'moving', 'romantic', 'monster', 'movie', 'the', 'shape', 'of', 'water', 'starring', 'sally', 'hawkins', 'and', 'doug', 'jones', 'plus', 'analyzing', 'the', 'films', 'bigger', 'meaning', 'and', 'themes', 'subscribe', 'support', 'foundflix', 'on', 'patreon', 'connect', 'with', 'us', 'on', 'social', 'media', 'facebook', 'twitter', 'instagram']
Step 8: After Join the cleaned words back into a single string	ending explained for the latest from master guillermo del toro the moving romantic monster movie the shape of water starring sally hawkins and doug jones plus analyzing the films bigger meaning and themes subscribe support foundflix on patreon connect with us on social media facebook twitter instagram

4.3.2.4 Encoding Features

In this step, the comment disable feature is converted from True/False representation to numerical values (0 or 1). This transformation prepares the data for machine learning algorithms that require numerical inputs, facilitating analysis, training, and predictions. Table 4.10, 4.11 illustrates the outcome of this conversion process.

Table 4.10: The Result Before Handling Non-Numerical Features

<u>Video_id</u>	<u>Comment_disable</u>
n1WpP7iowLc	FALSE
0dBIkQ4Mz1M	FALSE
5qpjK5DgCt4	TRUE
8HNuRNi8t70	FALSE

Table 4.11: The Result After Handling Non-Numerical Features

<u>Video_id</u>	<u>Comment_disable</u>
n1WpP7iowLc	0
0dBikQ4Mz1M	0
5qjK5DgCt4	1
8HNuRNi8t70	0

4.3.2.5 Processing Text Data Using Sentiment Analysis

The table 4.12 shows the proportion of sentiment conveyed in video descriptions, titles tags, and channel titles.

Table 4.12: The Textual Features Extraction

<u>Video_id</u>	<u>Video Description</u>	<u>Video Title</u>	<u>Channal Title</u>	<u>Tags</u>
0dBikQ4Mz1M	0.250000	-0.700000	0.0	0.150000
d380meD0W0M	0.459091	0.000000	0.0	-0.350000
2Vv-BfVoq4g	0.200000	1.000.000	0.0	0.136364
0yIWz1XEeyc	-0.131818	0.000000	0.0	0.000000
_uM5kFfkhB8	0.250000	0.136364	0.0	0.056667

Allowing the extraction and analysis of sentiment from textual data provides valuable insights into the emotions conveyed in the videos' descriptions, titles, tags and channel titles. The objective was to categorize the overall sentiment as positive, negative, or neutral, providing a deeper understanding of how viewers may perceive the videos based on their titles and descriptions. Understanding the sentiment behind the textual content helps reveal the emotional impact and tone of the videos. This information is valuable for influencers, marketers, and researchers, as it can be used to optimize video titles, descriptions, tags, and channel titles to enhance engagement and audience response.

To illustrate the distribution of positive, negative, and neutral sentiments across the textual features, Figures (4.2), (4.3), and (4.4) are presented.

Additionally, Figures (4.5), (4.6), and (4.7) show the ratio of positive, negative, and neutral (zero) sentiments in the text features.

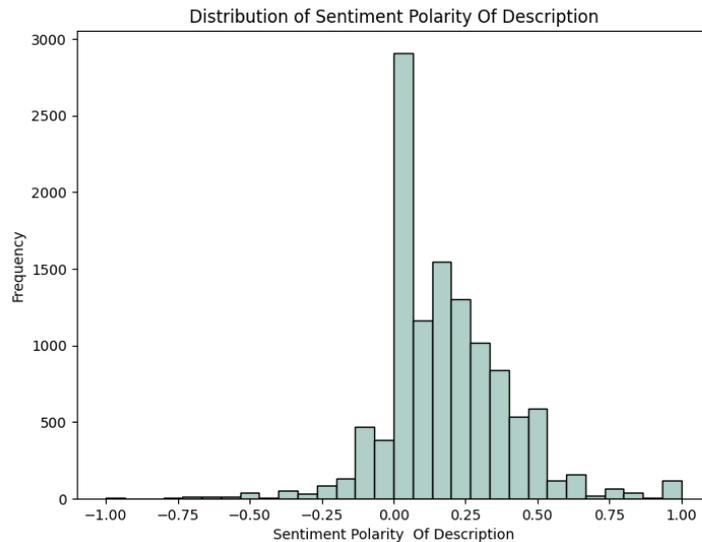


Figure 4.2: Distribution Sentiments Across the Description Feature

Overall, the figure 4.12 illustrates the diversity of sentiments in video descriptions, including neutral, moderately positive, very positive, and very negative sentiments.

The sentiment analysis results provide a nuanced view of the emotional tone conveyed in the video descriptions, which could be valuable for understanding the content's overall tone and potentially predicting user engagement or response based on sentiment

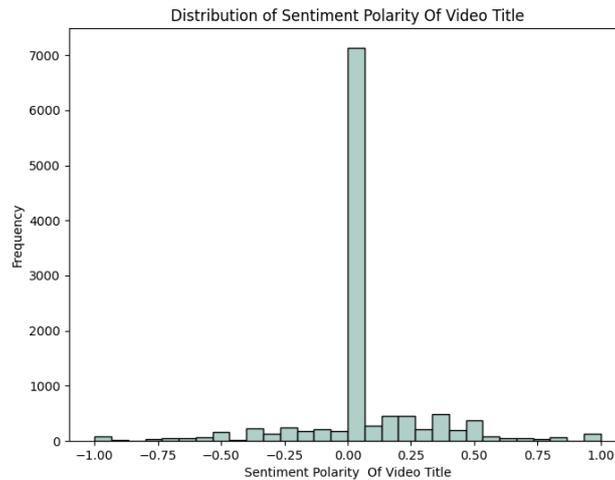


Figure 4.3: Distribution Sentiments Across the Video Title Feature

The figure 4.3 indicates that a substantial proportion of sentiments in channel titles is neutral, with a smaller number being highly positive or highly negative. This suggests the existence of extreme sentiments in both the positive and negative directions. The presence of such strong sentiments implies a diversity of emotional tones across various channel titles.

The inclusion of highly positive and highly negative sentiments offers insights into the emotional characteristics conveyed through channel titles. Understanding these emotional nuances is valuable for comprehending how the language employed in titles may influence viewer perception and engagement.

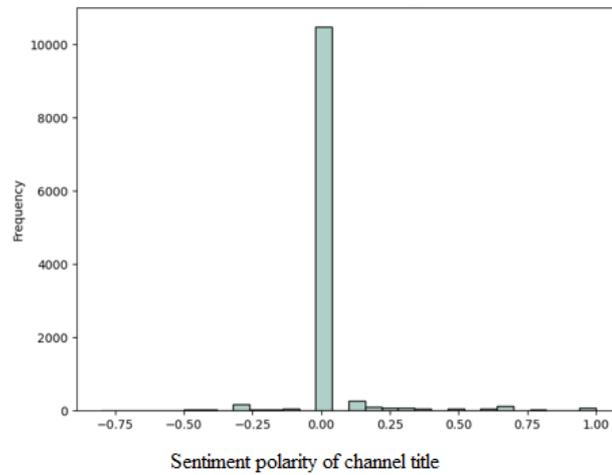


Figure 4.4: Distribution Sentiments Across Channel Title Feature

Similarly, concerning the channel title, the figure 4.14 illustrates a significant proportion of neutral sentiments, accompanied by a few instances of highly positive and highly negative expressions.

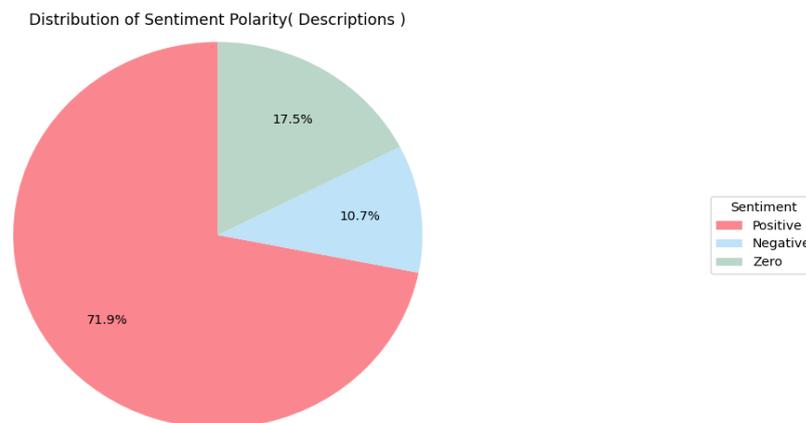


Figure 4.5: The Ratio of Positive (+), Negative (-), and Neutral (Zero) Sentiments in the Video Title Feature Description

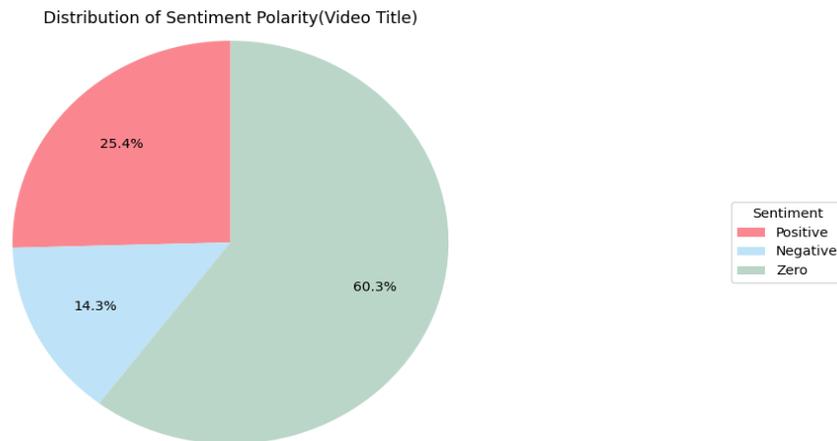


Figure 4.6: The Ratio of Positive (+), Negative (-), And Neutral (Zero) Sentiments in The Video Title Feature

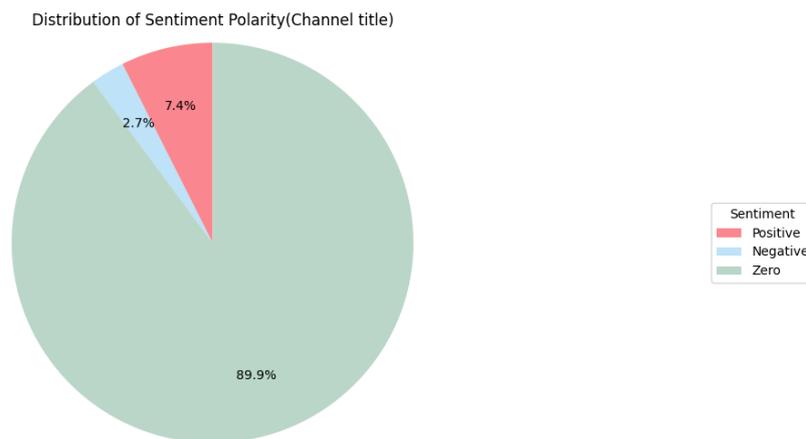


Figure 4.7: The Ratio of Positive (+), Negative (-), And Neutral (Zero) Sentiments in the Channel Title Feature

The figure 4.5,4.6 and 4.7 allows a comprehensive comparison of sentiment distribution across the mentioned features, providing insights into the emotional tones associated with each.

The percentages reveal the relative prevalence of positive, negative, and neutral sentiments within each feature.

Understanding the sentiment distribution in these features is crucial for gauging the overall emotional appeal and tone of video content. For example, positive sentiments may indicate engaging and appealing content, while negative sentiments might suggest controversial.

4.3.2.6 Time Interval Features Extraction

Time interval feature provides valuable insights into how quickly a video gained popularity after being published. A smaller Time interval indicates that a video became popular shortly after its release, while a larger Time interval suggests that the video took more time to gain traction and trend among viewers. Table (4.13) displays the calculated Time interval values for each video in the dataset.

Table 4.13: Time interval Feature for Each Video in the Dataset

video_id	publish_time	trending_date	Time interval(in day)
0dBikQ4Mz1M	13/11/2017 17:00	14/11/2017	0.291667
d380meD0W0M	12/11/2017 18:01	14/11/2017	1.248.831
2Vv-BfVoq4g	09/11/2017 11:04	14/11/2017	4.538.727
0yIWz1XEeyc	13/11/2017 07:37	14/11/2017	0.682049
_uM5kFfkhB8	12/11/2017 23:52	14/11/2017	1.005.405

4.3.3 Results of the Classification Methods

In this thesis, a detailed comparison of video classification techniques was conducted, employing various machine learning algorithms and combinations of features. To evaluate the models, a dataset comprising 24,427 categorized videos was employed, which was divided into 70% for training and 30% for testing, as illustrated in Figure 4.8.



Figure 4.8: Separation of Data into Training and Testing Data

Types of features used:

- Metadata (Statistic ,5 features): such Like, Dislike, Comment, Comment Disable, Category ID which exist in dataset.
- Extracted Features:
 - Extracted Text Features using Sentiment Analysis (4 features): Video Description, tags, Video Title, Video Channel Title.
 - Time Interval (1 features): The time between when a video was published and when it became popular.
 - Visual or Extracted Features from Images (6 features): Object Category, Object Weights.

Now, a brief overview of the key findings from my evaluations will be provided:

4.3.3.1 Evaluation Using All Features

Following the preprocessing phases, the stages of feature extraction, classification, and performance evaluation are implemented. The prediction accuracy is obtained in the first phase based on all extracted and existing features in the dataset.

Table 4.14 summarizes the findings. The results indicates that Random Forest and XGBoost demonstrated superior performance compared to other algorithms, suggesting their efficacy in managing the dataset's intricacies and patterns. This effectiveness may be attributed to their capacity to capture nonlinear relationships, address feature interactions, and mitigate overfitting. Furthermore, the selection of attributes for analysis appears to align well with these algorithms. Notably, these algorithms are recognized for their capability to evaluate feature importance, offering valuable insights into the attributes that substantially influence the predictive capability of the model. In terms of performance, K-Nearest Neighbor sit in between RF/XGB and SVC. They outperform SVC but fall short of RF and XGB, seen figure 4.9.

- a. Random Forest accuracy reaching to 96.94.
- b. Extreme Gradient Boosting accuracy reaching to 96.33%.
- c. K-Nearest Neighbor accuracy to reaching 91.26%.
- d. Support Vector Machine accuracy reaching to 87.55%.

Table 4.14: Accuracy Comparison of Classifiers Using all features

Algorithms	Accuracy	Precision	Recall	f-score
KNN	91.26	91.61	91.02	90.92
RF	96.94	97.00	96.86	96.89
SVC	87.55	87.77	87.37	87.40
XGB	96.33	96.38	96.24	96.27

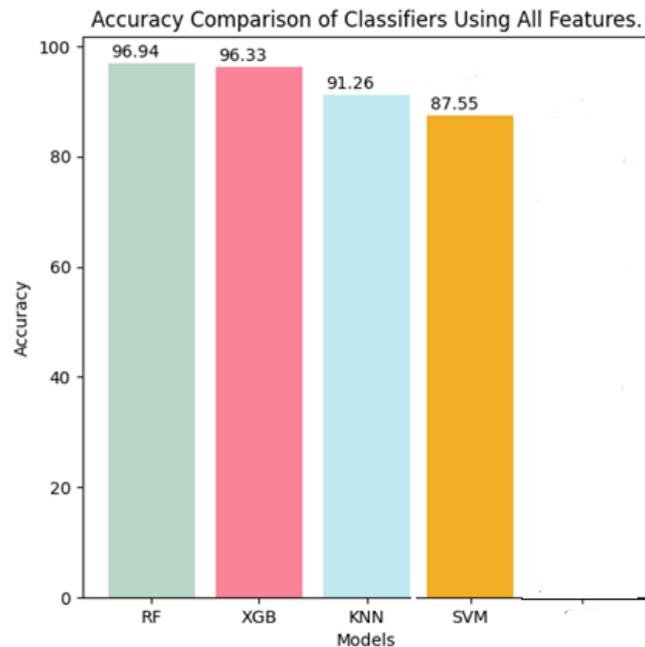


Figure 4.9 Accuracy Comparison of Classifiers Using all features

4.3.3.2 Evaluation Using Extracted Features Only

The initial assessment revealed that the Random Forest (RF) algorithm achieved a baseline accuracy of 90.95%, surpassing the slightly lower accuracy of 79.98% attained by XGBoost (XGB). This initial performance sets the stage as a reference point for subsequent enhancements.

The next phase involved integrating visual features extracted from images into the models. As a result, the Random Forest accuracy saw a notable improvement to 93.00%, and XGBoost also benefited from the inclusion of visual features, reaching an accuracy of 84.48%. This indicates the effectiveness of these algorithms for video classification tasks, as shown in Table 4.15.

Table 4.15: Shows That Random Forest (RF) And XGBoost (XGB) Outperformed the Previous Research on the Same Dataset Using Extracted Features Only.

References	Algorithms	F1 score using Extracted features (Gap, text features)	
[1]	RF	70%	
	XGB	74%	
This study	Algorithms	F1 score using text and time interval features	F1 score using text features, time interval and visual features
	RF	90.95%	93.00 %
	XGB	79.98%	84.48 %

These outcomes strongly suggest that incorporating visual features played a pivotal role in elevating the performance of both Random Forest and XGBoost. The consistent trend of accuracy improvement with the inclusion of visual features emphasizes the significance of a comprehensive feature set in achieving enhanced model performance.

In summary, the initial findings underscore the effectiveness of text and time-lapse features, with a substantial accuracy boost observed upon integrating visual features. This discussion sheds light on the iterative nature of feature engineering and model enhancement, offering valuable insights for informed decision-making in future analyses or applications.

Random Forest consistently achieved well across all evaluations, indicating its reliability and making it an appropriate option across different feature sets. In contrast, Extreme Gradient Boosting performed well in the first two evaluations but fell short in the third evaluation when only extracted features

were used. Carefully chosen features extracted from video thumbnails, descriptions, tags, and titles can benefit marketers and influencers by providing insights into the nature of the video, ultimately helping them make their videos more popular, figure 4.10 Explain Confusion Matrix of Random Forest Classifier.

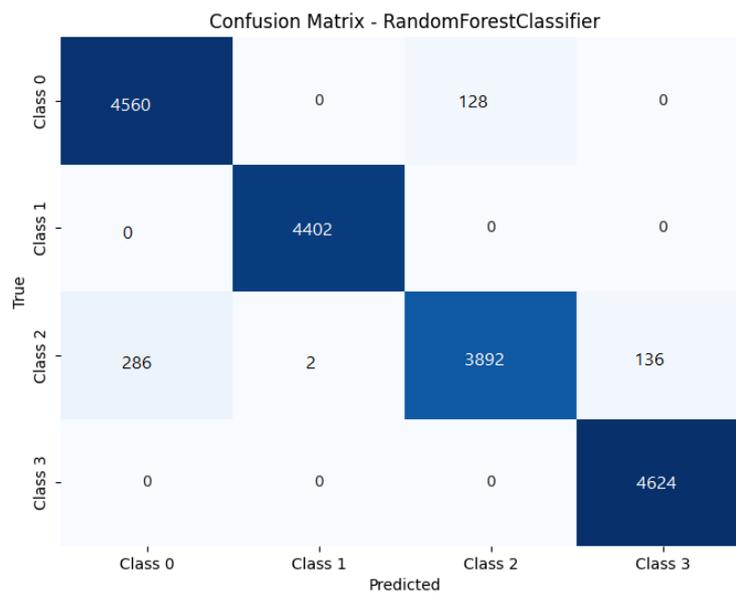


Figure 4.10 Confusion Matrix of Random Forest Classifier

Chapter Five

Conclusions and Future Works

CHAPTER FIVE

Conclusions and Future Works

5.1 Conclusions

In summary, valuable insights into the effective strategies for predicting video popularity on social media platforms are provided by the findings of the study. The following conclusions highlight key approaches and outcomes that contribute to understanding influential factors and model performance.

1. A highly effective approach is employed when the prediction model combines various features.
2. An effective strategy has been demonstrated through the extraction of additional features, including those related to video thumbnails.
3. Enhanced insights are achieved through the utilization of sentiment analysis, assisting in the better understanding of emotions and feelings in textual features.
4. Remarkable accuracy levels of nearly 96.94% and 96.33%, respectively, were achieved through the utilization of robust classification algorithms like Random Forest and Extreme Gradient Boosting.
5. Lower accuracy levels are observed for SVM compared to Random Forest and XGBoost, indicating that these models may not be the best choices for this specific dataset and feature set.
6. The real value of analysis is found in gaining insights into the factors influencing video popularity like Thumbnail-related Features, video description.

Finally, the opportunity to customize video content for increased appeal to the target audience is offered by understanding the impact of titles, tags, descriptions, and thumbnails on viewer engagement. This, in turn, increases the likelihood of video success on social media platforms.

5.2 Future Work

Building upon the findings and methodologies presented in this thesis offers numerous opportunities for future work aimed at enhancing the prediction of video popularity on social media platforms. Potential areas for exploration include:

1. Explore and scrutinize the intricate relationship between video thumbnails and the corresponding video content in more detail. Investigate how specific elements in thumbnails contribute to viewer engagement and popularity.
2. Examine images using models trained on a broader array of classes to extract supplementary features that demonstrate enhanced effectiveness and accuracy. Evaluate the impact of diverse training data on the model's ability to capture nuanced visual information.
3. Refine and optimize the prediction model by exploring additional combinations of features using API (Application Programming Interface). Investigate the incorporation of real-time data and evolving trends to improve the model's predictive capabilities.
4. Extend the research beyond YouTube and include other social media platforms, such as TikTok. Analyse the platform-specific factors influencing video popularity, providing a comprehensive understanding of the diverse dynamics across platforms.

5. Compare and contrast the factors influencing video popularity on different platforms. This comparative analysis will enable influencers and marketers to tailor their strategies based on the unique characteristics and preferences of each platform's audience.

By addressing these areas of future work, researchers can contribute to the advancement of the video popularity prediction field. This, in turn, will facilitate the development of more effective content marketing and user engagement strategies on a variety of social media platforms.

References

References

- [1] Y. Li, K. Eng, and L. Zhang, “YouTube Videos Prediction: Will this video be popular?,” 2019, [Online]. Available: http://cs229.stanford.edu/proj2019aut/data/assignment_308832_raw/26647615.pdf (accessed 9 May 2020).
- [2] M. L. Khan, “Social media engagement: What motivates user participation and consumption on YouTube?,” *Comput. Human Behav.*, vol. 66, pp. 236–247, Jan. 2017, doi: 10.1016/J.CHB.2016.09.024.
- [3] F. Figueiredo, F. Benevenuto, and J. M. Almeida, “The tube over time: Characterizing popularity growth of YouTube videos,” *Proc. 4th ACM Int. Conf. Web Search Data Mining, WSDM 2011*, pp. 745–754, 2011, doi: 10.1145/1935826.1935925.
- [4] Q. I. A. O. Sibó, P. A. N. G. Shanchen, W. A. N. G. Min, Z. H. A. I. Xue, and D. A. I. Feng, “Online Video Popularity Regression Prediction Model with Multichannel Dynamic Scheduling Based on User Behavior,” *Chinese J. Electron.*, vol. 30, no. 5, pp. 876–884, Sep. 2021, doi: 10.1049/cje.2021.06.010.
- [5] L. Jiang, Y. Miao, Y. Yang, Z. Lan, and A. G. Hauptmann, “Viral video style: A closer look at viral videos on YouTube,” *ICMR 2014 - Proc. ACM Int. Conf. Multimed. Retr. 2014*, pp. 193–200, 2014, doi: 10.1145/2578726.2578754.
- [6] R. Shreyas, D. M. Akshata, B. S. Mahanand, B. Shagun, and C. M. Abhishek, “Predicting popularity of online articles using Random Forest regression,” *Proc. - 2016 2nd Int. Conf. Cogn. Comput. Inf. Process. CCIP 2016*, 2016, doi: 10.1109/CCIP.2016.7802890.

- [7] T. Trzcinski and P. Rokita, “Predicting popularity of online videos using Support Vector Regression,” Oct. 2015, doi: 10.1109/TMM.2017.2695439.
- [8] W. Stokowiec, T. Trzciński, K. Wołk, K. Marasek, and P. Rokita, “Shallow reading with deep learning: Predicting popularity of online content using only its title,” *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 10352 LNAI, pp. 136–145, 2017, doi: 10.1007/978-3-319-60438-1_14.
- [9] T. Trzcinski, P. Andruszkiewicz, T. Bochenski, and P. Rokita, “Recurrent Neural Networks for Online Video Popularity Prediction,” Jul. 2017, doi: 10.1007/978-3-319-60438-1_15.
- [10] F. Huang, J. Chen, Z. Lin, P. Kang, and Z. Yang, “Random forest exploiting post-related and user-related features for social media popularity prediction,” *MM 2018 - Proc. 2018 ACM Multimed. Conf.*, pp. 2013–2017, 2018, doi: 10.1145/3240508.3266439.
- [11] A. Khan, G. Worah, M. Kothari, Y. H. Jadhav, and A. V. Nimkar, “News Popularity Prediction with Ensemble Methods of Classification,” *2018 9th Int. Conf. Comput. Commun. Netw. Technol. ICCCNT 2018*, pp. 1–6, 2018, doi: 10.1109/ICCCNT.2018.8494095.
- [12] K. R. Purba, D. Asirvatham, and R. K. Murugesan, “Instagram post popularity trend analysis and prediction using hashtag, image assessment, and user history features,” *Int. Arab J. Inf. Technol.*, vol. 18, no. 1, pp. 85–94, 2021, doi: 10.34028/iajit/18/1/10.
- [13] M. U. N. Nisa, D. Mahmood, G. Ahmed, S. Khan, M. A. Mohammed, and R. Damaševičius, “Optimizing prediction of youtube video

- popularity using xgboost,” *Electron.*, vol. 10, no. 23, 2021, doi: 10.3390/electronics10232962.
- [14] M. L. Khan, “Social media engagement: What motivates user participation and consumption on YouTube?,” *Comput. Human Behav.*, vol. 66, pp. 236–247, 2017, doi: 10.1016/j.chb.2016.09.024.
- [15] A. Neyaz, A. Kumar, S. Krishnan, J. Placker, and Q. Liu, “Security, Privacy and Steganographic Analysis of FaceApp and TikTok,” *Int. J. Comput. Sci. Secur.*, vol. 14, no. 2, pp. 38–59, 2020, [Online]. Available: <https://www.researchgate.net/publication/341782197>
- [16] Y. Fan, B. Yang, D. Hu, X. Yuan, and X. Xu, “Social- And Content-Aware Prediction for Video Content Delivery,” *IEEE Access*, vol. 8, pp. 29219–29227, 2020, doi: 10.1109/ACCESS.2020.2972920.
- [17] S. L. de Sá, A. A. d. A. Rocha, and A. Paes, “Predicting popularity of video streaming services with representation learning: A survey and a real-world case study,” *Sensors*, vol. 21, no. 21, 2021, doi: 10.3390/s21217328.
- [18] L. Yue, W. Chen, X. Li, W. Zuo, and M. Yin, “A survey of sentiment analysis in social media,” *Knowl. Inf. Syst.*, vol. 60, no. 2, pp. 617–663, 2019, doi: 10.1007/s10115-018-1236-4.
- [19] M. Ham and S. W. Lee, “Factors affecting the popularity of video content on live-streaming services: Focusing on V live, the South Korean live-streaming service,” *Sustain.*, vol. 12, no. 5, pp. 1–17, 2020, doi: 10.3390/su12051784.
- [20] J. Davidson, B. Liebald, J. Liu, P. Nandy, and T. Van Vleet, “The YouTube video recommendation system,” *RecSys’10 - Proc. 4th ACM Conf. Recomm. Syst.*, no. August 2014, pp. 293–296, 2010, doi:

- 10.1145/1864708.1864770.
- [21] “Key Features of Popular Social Media Platforms.” <https://techeconomy.ng/key-features-of-popular-social-media-platforms/> (accessed Aug. 01, 2023).
- [22] L. H. X. Ng, J. Y. H. Tan, D. J. H. Tan, and R. K. W. Lee, “Will you dance to the challenge?: Predicting user participation of TikTok challenges,” in *Proceedings of the 2021 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, ASONAM 2021*, Nov. 2021, pp. 356–360. doi: 10.1145/3487351.3488276.
- [23] Mitchell J., “Trending YouTube Video Statistics | Kaggle,” *Kaggle*, 2019. <https://www.kaggle.com/datasnaek/youtube-new> (accessed Mar. 11, 2023).
- [24] J. M. Luna, *Introduction to Data Mining*. 2021. doi: 10.1007/978-981-16-3964-7_1.
- [25] S. Khalid, T. Khalil, and S. Nasreen, “A survey of feature selection and feature extraction techniques in machine learning,” *Proc. 2014 Sci. Inf. Conf. SAI 2014*, no. July, pp. 372–378, 2014, doi: 10.1109/SAI.2014.6918213.
- [26] O. I. Obaid, M. A. Mohammed, A. O. Salman, S. A. Mostafa, and A. A. Elngar, “Comparing the Performance of Pre-trained Deep Learning Models in Object Detection and Recognition,” *J. Inf. Technol. Manag.*, vol. 14, no. 4, pp. 40–56, 2022, doi: 10.22059/JITM.2022.88134.
- [27] M. Jogin, “Feature Extraction using Convolution Neural Networks (CNN) and Deep Learning,” *2018 3rd IEEE Int. Conf. Recent Trends*

- Electron. Inf. Commun. Technol.*, no. November, pp. 2319–2323, 2020, doi: 10.1109/RTEICT42901.2018.9012507.
- [28] M. Zabir, N. Fazira, Z. Ibrahim, and N. Sabri, “Evaluation of pretrained Convolutional Neural Network models for object recognition,” *Int. J. Eng. Technol.*, vol. 7, no. 3, pp. 95–98, 2018, doi: 10.14419/ijet.v7i3.15.17509.
- [29] D. G. K. Mahshiya VM, Manju VM, “A SURVEY ON OBJECT DETECTION TECHNIQUES USING TENSORFLOW , KERAS AND YOLO,” no. 4, pp. 495–500, 2022.
- [30] B. S. Rekha, A. Marium, G. N. Srinivasan, and S. A. Shetty, “Literature Survey on Object Detection using YOLO,” *Int. Res. J. Eng. Technol.*, vol. 07, no. 06, pp. 3082–3088, 2020, [Online]. Available: <https://www.irjet.net/archives/V7/i6/IRJET-V7I6576.pdf>
- [31] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, vol. 2016-Decem, pp. 779–788, 2016, doi: 10.1109/CVPR.2016.91.
- [32] A. P. Jana, A. Biswas, and Mohana, “YOLO based detection and classification of objects in video records,” *2018 3rd IEEE Int. Conf. Recent Trends Electron. Inf. Commun. Technol. RTEICT 2018 - Proc.*, no. April, pp. 2448–2452, 2018, doi: 10.1109/RTEICT42901.2018.9012375.
- [33] S. A. Alasadi and W. S. Bhaya, “Review of data preprocessing techniques in data mining,” *J. Eng. Appl. Sci.*, vol. 12, no. 16, pp. 4102–4107, 2017, doi: 10.3923/jeasci.2017.4102.4107.
- [34] M. H. Zhu, “Research on data preprocessing in exam analysis

- system,” *Lect. Notes Electr. Eng.*, vol. 100 LNEE, no. VOL. 4, pp. 333–338, 2011, doi: 10.1007/978-3-642-21762-3_43.
- [35] S. García, S. Ramírez-Gallego, J. Luengo, J. M. Benítez, and F. Herrera, “Big data preprocessing: methods and prospects,” *Big Data Anal.*, vol. 1, no. 1, pp. 1–22, 2016, doi: 10.1186/s41044-016-0014-0.
- [36] C. Fan, M. Chen, X. Wang, J. Wang, and B. Huang, “A Review on Data Preprocessing Techniques Toward Efficient and Reliable Knowledge Discovery From Building Operational Data,” *Front. Energy Res.*, vol. 9, no. March, pp. 1–17, 2021, doi: 10.3389/fenrg.2021.652801.
- [37] E. Cho, T. W. Chang, and G. Hwang, “Data Preprocessing Combination to Improve the Performance of Quality Classification in the Manufacturing Process,” *Electron.*, vol. 11, no. 3, pp. 1–15, 2022, doi: 10.3390/electronics11030477.
- [38] V. Gupta and G. S. Lehal, “A survey of text mining techniques and applications,” *J. Emerg. Technol. Web Intell.*, vol. 1, no. 1, pp. 60–76, 2009, doi: 10.4304/jetwi.1.1.60-76.
- [39] A. Negi, “A Brief Survey On Text Mining, Its Techniques, And Applications,” *Int. J. Mob. Comput. Appl.*, vol. 8, no. 1, pp. 1–6, 2021, doi: 10.14445/23939141/ijmca-v8i1p101.
- [40] A. Shiri, “Introduction to Modern Information Retrieval (2nd edition),” *Libr. Rev.*, vol. 53, no. 9, pp. 462–463, 2004, doi: 10.1108/00242530410565256.
- [41] M. Mujahid *et al.*, “Sentiment analysis and topic modeling on tweets about online education during covid-19,” *Appl. Sci.*, vol. 11, no. 18,

- 2021, doi: 10.3390/app11188438.
- [42] I. H. Sodhar, “Computer Science,” no. December 2020, 2022, doi: 10.22271/ed.book.784-CITATIONS.
- [43] A. G. Reece and C. M. Danforth, “Instagram photos reveal predictive markers of depression,” *EPJ Data Sci.*, vol. 6, no. 1, 2017, doi: 10.1140/epjds/s13688-017-0110-z.
- [44] D. Khurana, A. Koli, K. Khatter, and S. Singh, “Natural language processing: state of the art, current trends and challenges,” *Multimed. Tools Appl.*, vol. 82, no. 3, pp. 3713–3744, 2023, doi: 10.1007/s11042-022-13428-4.
- [45] B. Priya, J. M. Nandhini, and T. Gnanasekaran, “An analysis of the applications of natural language processing in various sectors,” *Adv. Parallel Comput.*, vol. 38, pp. 598–602, 2021, doi: 10.3233/APC210109.
- [46] B. S. Ainapure *et al.*, “Sentiment Analysis of COVID-19 Tweets Using Deep Learning and Lexicon-Based Approaches,” *Sustain.*, vol. 15, no. 3, pp. 1–22, 2023, doi: 10.3390/su15032573.
- [47] S. Velickov and D. Solomatine, “Predictive Data Mining : Practical Examples Abstract : 2 . Data Mining – theoretical and practical aspects,” *Neural Networks*, no. March, pp. 1–17, 2000.
- [48] U. M. D. E. C. D. E. Los, *Data Mining Concepts, Models and Techniques*, vol. 12. 2011. doi: 10.1007/978-3-642-19721-5_1.
- [49] P. Jaganathan, S. Vinothini, and P. Backialakshmi, “A Study of Data Mining Techniques to Agriculture,” *Int. J. Res. Inf. Technol.*, vol. 2, no. 4, pp. 306–313, 2014.
- [50] S. Mahdevari, K. Shahriar, S. Yagiz, and M. Akbarpour Shirazi, “A

- support vector regression model for predicting tunnel boring machine penetration rates,” *Int. J. Rock Mech. Min. Sci.*, vol. 72, pp. 214–229, 2014, doi: 10.1016/j.ijrmms.2014.09.012.
- [51] X. Wu *et al.*, *Top 10 algorithms in data mining*, vol. 14, no. 1. 2008. doi: 10.1007/s10115-007-0114-2.
- [52] A. Ibrahim Ahmed Osman, A. Najah Ahmed, M. F. Chow, Y. Feng Huang, and A. El-Shafie, “Extreme gradient boosting (Xgboost) model to predict the groundwater levels in Selangor Malaysia,” *Ain Shams Eng. J.*, vol. 12, no. 2, pp. 1545–1556, 2021, doi: 10.1016/j.asej.2020.11.011.
- [53] M. A. Ganaie, M. Hu, A. K. Malik, M. Tanveer, and P. N. Suganthan, “Ensemble deep learning: A review,” *Eng. Appl. Artif. Intell.*, vol. 115, 2022, doi: 10.1016/j.engappai.2022.105151.
- [54] R. Forests, “Chapter 15.1 Random Forest for Regression or Classification,” pp. 587–604, 2001, [Online]. Available: <http://www.math.usu.edu/>
- [55] D. A. Otchere, T. O. A. Ganat, J. O. Ojero, B. N. Tackie-Otoo, and M. Y. Taki, “Application of gradient boosting regression model for the evaluation of feature selection techniques in improving reservoir characterisation predictions,” *J. Pet. Sci. Eng.*, vol. 208, no. May, p. 109244, 2022, doi: 10.1016/j.petrol.2021.109244.
- [56] S. Wassan *et al.*, “Gradient Boosting for Health IoT Federated Learning,” *Sustain.*, vol. 14, no. 24, 2022, doi: 10.3390/su142416842.
- [57] T. Chen and C. Guestrin, “XGBoost: A scalable tree boosting system,” *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*,

vol. 13-17-Augu, pp. 785–794, 2016, doi:
10.1145/2939672.2939785.

الخلاصة

لقد أحدث النمو المتزايد لمنصات الوسائط الاجتماعية ومحتوى الفيديو عبر الإنترنت ثورة في الطريقة التي نتواصل بها ونتفاعل ونستهلك المعلومات. مع مشاركة الملايين من مقاطع الفيديو يوميًا على منصات مثل يوتيوب، أصبح التنبؤ بشعبية مقاطع الفيديو جانبًا حاسمًا بالنسبة للمؤثرين ومسؤولي المنصات والمسوقين. كذلك إن فهم العوامل التي تساهم في التنبؤ بشعبية مقطع الفيديو يمكن أن يؤثر بشكل كبير على تسويق المحتوى واستراتيجيات منصات يوتيوب.

تتناول هذه الأطروحة التحدي المتمثل في التنبؤ بدقة بشعبية الفيديو على منصات التواصل الاجتماعي، مع التركيز بشكل خاص على منصات يوتيوب. أقترح نهجًا مشتركًا جديدًا يستخدم تحليل البيانات الوصفية (بما في ذلك اختيار الميزات المؤثرة) وتحليل الصور المصغرة لإنشاء نموذج تنبؤ أكثر فعالية. ومن خلال استخراج الميزات ذات الصلة من هذه المصادر، أهدف إلى تحسين دقة التنبؤ بشعبية الفيديو. ولتحقيق هذا الهدف، تم استخدام مجموعة من خوارزميات التصنيف القوية، بما في ذلك XGB,GB,RF,KNN.

تم إجراء بعض عمليات المعالجة المسبقة على البيانات لجعلها مناسبة لخوارزميات التعلم الآلي. وبعد ذلك تم استخراج ثلاث أنواع من الميزات: الميزات النصية، والميزات المرئية، والميزات المرتبطة بالزمن. وقد أدت هذه الميزات المضافة حديثًا، جنبًا إلى جنب مع الميزات الموجودة، إلى إثراء مجموعة البيانات بميزات مؤثرة إضافية قادرة على تعزيز دقة النموذج. من خلال تدريب النموذج باستخدام هذه الخوارزميات وتنفيذ تقنيات استخراج الميزات، تم تحقيق دقة كبيرة، خاصة مع خوارزميات الغابة العشوائية وخوارزميات تعزيز التدرج الشديد أو الإضافي. تمت عملية التنبؤ على مرحلتين: التنبؤ باستخدام الميزات الأصلية والمستخرجة، والتنبؤ بالميزات المستخرجة فقط.

تم تقييم أداء النموذج باستخدام مقاييس مختلفة، بما في ذلك الدقة، والتي تصل إلى معدلات 96,33% و 96,94% لخوارزميات XGB وخوارزميات RF على التوالي باستخدام جميع الميزات و 84 و 93 باستخدام الميزات المستخرجة فقط. تكمن قيمة الجمع بين البيانات الوصفية وتحليل الصور المصغرة في اكتساب رؤية حول العوامل التي تؤثر على شعبية الفيديو. وهذا يمكن المؤثرين

والمسوقين من تخصيص المحتوى الخاص بهم بشكل أفضل، مما يلقي صدى لدى جمهورهم المستهدف ويزيد من احتمالية النجاح على منصات التواصل الاجتماعي.

Thesis Related Publication

Some of the works presented in this thesis have been accepted as below.

First Paper

Name of Conference: 2nd Al Hikma International Conference on Natural and Applied Sciences (HICNAS2023).

Name of Journal: Lecture Notes in Networks and Systems

Paper Title: Predicting Content Popularity on social media: An Analytical Approach using Regression Modeling

Authors: 1) Heba Hussien 2) Wadhah R. Baiee. Software Department, College of Information Technology, Babylon University, Iraq.

Second Paper

Name of Conference: 2nd Al Hikma International Conference on Natural and Applied Sciences (HICNAS2023).

Name of Journal: AIP Conference Proceedings.

Paper Title: Maximizing Video Popularity Prediction: A Holistic Approach Utilizing Metadata and Thumbnail Analysis.

Authors: 1) Heba Hussien 2) Wadhah R. Baiee. Software Department College of Information Technology, Babylon University, Iraq.

Email: hibaha.sw.msc@student.uobabylon.edu.iq

Email: wadhah.baiee@uobabylon.edu.iq



LETTER OF ACCEPTANCE

Date: August 8, 2023

Paper ID: HICNAS_377

Dear Authors,

On behalf of the HICNAS-23 Scientific Committee, and based on our (Reviewers' Evaluation, Scientific Committee Decision, and Guest Editors' Approval), we are pleased to inform you that your paper entitled:

[Maximizing Video Popularity Prediction: A Holistic Approach Utilizing Metadata and Thumbnail Analysis](#)

Written By

[Heba Al-Mamouri and Wadhah R. Boiee](#)

Has been accepted in HICNAS Conference and selected to be processed for possible publication in **AIP Conference Proceedings** (ISSN: 0094-243X, 1551-7616). We have accepted your paper based on our initial technical checking, the positive double-blind reviews, and your Paper's Scope and Subject. Remember that your paper in Camera Ready Version should accurately follow our reviewer comments and the technical notes. The publication Schedule of the accepted paper will be provided after passing the Internal Check of AIP Editors. The paper should not contain plagiarism till that date (more than 15%), and the Camera Ready version should follow our conference Rules and AIP Ethics and Guidelines. Publication time depends on Publication Queue on AIP, and we will provide you with the required information after we finalize all selected Papers for our Conference and present the accepted papers at the conference sessions.

Thank you for considering this Conference as a venue for your work

sincerely yours

Prof. Dr. Shubham Sharma
Conference Guest Editor

Prof. Dr. Mohammed I. Mohammed
Chair of Scientific Committee

CAUTION: This Acceptance Letter is Made by ICAS (Special Sessions) Conference Guest Editors, All Approval and other inquiries Should be addressed to Conference Editorial Board and Patrons Committee, as all Accepted Papers will be Process for Possible Publication in the AIP Conference Proceedings, and the Final Decision upon Publish of your Paper will be Made by AIP Publication Editors Only after Review the Paper Contents and Writing Quality

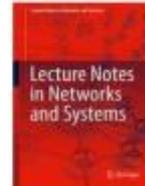


+964 7819166679
+964 7821601603
+964 7716450783
+964 7713997116



info.hicnas23@hiuc.edu.iq

LETTER OF ACCEPTANCE



Paper ID: HICNAS_477
Submission ID: 202305
Date: 13 August 2023

Paper Title: Predicting Content Popularity on Social Media: An Analytical Approach using Regression Modeling
Author (s): Heba Al-Mamouri and Wadhah R. Baiee

Congratulations!

Based on the recommendations of the Technical Program Committee of (ICCCNet-2023), we are pleased to inform you that your manuscript has been **Accepted as a REGULAR paper** in our Conference Partner ICCCN₂₀₂₃ at MANCHESTER METROPOLITAN UNIVERSITY, UNITED KINGDOM and will be processed for Publication in the Springer Series **"Lecture Notes in Networks and Systems"** [ISSN: 2367-3389; 2367-3370] (**Scopus Indexed**). The paper shall appear in ICCCN-2023 in Lecture Notes in Networks and Systems. Remember that your paper in Camera Ready Version should accurately follow our reviewer comments and technical notes. The publication Schedule of the accepted paper will be provided after passing the Internal Check of Springer Editors. The paper should not contain plagiarism till that date (**more than 15%**), and the Camera Ready version should follow our conference Rules and Springer Guidelines. Publication time will provide you with the required information after we finalize all selected Papers; *please do not hesitate to call us for any assistance.*

We will encourage more quality submissions from you and your colleagues in the future.

sincerely yours

Prof. Dr. Mohammed I. Mohammed
 Chair of Scientific Committee

TPC Cahir
 Conference Chair
 SPRINGER ICCCN-2023

CAUTION: This Acceptance Letter is Made by HICNAS Conference Guest Editors, and its Subjected to All Approval of any Inquiries Should be addressed to Conference Editorial Board and Patrons Committee, as all Accepted Papers will be Process for Possible Publication in the LNNS Springer Series in our Partner conferences ICCCN at MANCHESTER METROPOLITAN UNIVERSITY, at UNITED KINGDOM.



+964 7819166679
 +964 7821601603
 +964 7716450783
 +964 7713997116



 info.hicnas23@hiuc.edu.iq



جمهورية العراق
وزارة التعليم العالي والبحث العلمي
جامعة بابل
كلية تكنولوجيا المعلومات
قسم البرمجيات

التنبؤ بشعبية الفيديو استنادًا إلى بيانات اليوتيوب الوصفية والصور المصغرة

رسالة مقدمة إلى

مجلس كلية تكنولوجيا المعلومات - جامعة بابل كجزء من متطلبات
نيل درجة الماجستير في تكنولوجيا المعلومات / البرمجيات

من قبل

هبة حسين عبد العباس نايف

بإشراف

م.د وضاح رزوقي عبود حسن بيبي

م ٢٠٢٣

١٤٤٥ هـ