



**Department of Electrical and Computer Engineering
North South University**

Senior Design Project Report
**Predictive Modeling of Social Movement Lifespan in
Social Media**

ABU BAKAR SIDDIK MINHAZ (2031589642)

SIAM MASUD (2011057642)

MD ARIFUR RAHAMAN (2014198642)

MOHAMMAD SADMAN WASIF (2011832042)

Faculty Advisor:

Dr. Sifat Momen

Associate Professor

ECE Department

Summer, 2024

LETTER OF TRANSMITTAL

December, 2023

To

Dr. Mohammad Abdul Matin
Chairman,
Department of Electrical and Computer Engineering
North South University, Dhaka

Subject: Submission of Capstone Project Report on “Predictive Modeling of Social Movement Lifespan in Social Media”

Dear Sir,

With due respect, we would like to submit our **Capstone Project Report** on “**Predictive Modeling of Social Movement Lifespan in Social Media**” as a part of our BSc program. The report deals with the analysis of Social Media Trends and their lifespans. This project was beneficial to us in gaining experience in the practical field and applying it in real life. We tried to the maximum competence to meet all the dimensions required from this report.

We will be highly obliged if you kindly receive this report and provide your valuable judgment. It would be our immense pleasure if you find this report helpful and informative to have an apparent perspective.

Sincerely Yours,

.....
ABU BAKAR SIDDIK MINHAZ
ECE Department
North South University, Bangladesh

.....
SIAM MASUD
ECE Department
North South University, Bangladesh

.....
MD ARIFUR RAHAMAN
ECE Department
North South University, Bangladesh

.....
MD SADMAN WASIF
ECE Department
North South University, Bangladesh

APPROVAL

ABU BAKAR SIDDIK MINHAZ (2031589642), SIAM MASUD (2011057642), MD ARIFUR RAHAMAN (2014198642) & MD SADMAN WASIF (2011832042) from the Electrical and Computer Engineering Department of North South University have worked on the Senior Design Project titled “Predictive Modeling of Social Movement Lifespan in Social Media” under the supervision of Dr. Sifat Momen partial fulfillment of the requirement for the degree of Bachelors of Science in Engineering and has been accepted as satisfactory.

Supervisor’s Signature

.....

Dr. Sifat Momen

Associate Professor

Department of Electrical and Computer Engineering

North South University

Dhaka, Bangladesh.

Chairman’s Signature

.....

Dr. Mohammad Abdul Matin

Professor

Department of Electrical and Computer Engineering

North South University

Dhaka, Bangladesh.

DECLARATION

It is to declare that this project is our original work. No part of this work has been submitted elsewhere, partially or entirely, for the award of any other degree or diploma. All project-related information will remain confidential and shall not be disclosed without the formal consent of the project supervisor. Relevant previous works presented in this report have been adequately acknowledged and cited. The plagiarism policy, as stated by the supervisor, has been maintained.

Students' names & Signatures

1. ABU BAKAR SIDDIK MINHAZ

2. SIAM MASUD

3. MD ARIFUR RAHAMAN

4. MOHAMMAD SADMAN WASIF

ACKNOWLEDGEMENTS

The authors would like to express their heartfelt gratitude towards their project and research supervisor, Dr. Sifat Momen, Associate Professor, Department of Electrical and Computer Engineering, North South University, Bangladesh, for his invaluable support, precise guidance, and advice on the experiments, research, and theoretical studies carried out during the course of the current project and also in the preparation of the report.

Furthermore, the authors would like to thank the Department of Electrical and Computer Engineering, North South University, Bangladesh, for facilitating the research. We would also like to thank all our friends who helped us with this project. The authors would also like to thank their loved ones for their countless sacrifices and continual support.

ABSTRACT

Predictive Modeling of Social Movement Lifespan in Social Media

Social media has turned out to be a vibrant medium through which social movements are effectuated and through which individuals and groups can make their voices heard on vital issues. It is important to learn what determines the success and sustainability of these movements for activists, policymakers, and researchers. This study will propose a framework that will be able to predict the life of social media-driven movements through sentiment and thematic pattern analyses extracted from public posts. Sentiment analysis then categorized the posts as positive, negative, and neutral, while thematic analysis showed repeating themes such as political and legal contexts.

This dataset was created using a custom web scraper that preprocessed the data into key features including sentiment scores, thematic relevance, and indications of how long the movement lasted. Dimensionality reduction and clustering analyses are conducted to find patterns in the features and provide insight that may be more interpretable. The project used a variety of machine learning and ensemble models: linear regression, support vector regressor, random forest regressor, gradient boosting regressor, XGBoost, LightGBM, and refined deep learning approaches. Simple and optimized weighted ensembles comprised ensemble techniques, where the highest accuracy was achieved by the meta-model. Its predictive accuracy was the best, at minimum MAE and MSE.

Feature interdependencies were analyzed using the method of correlation analysis, while for model interpretability, SHAP-an important technique in machine learning-was used. Accordingly, Legality Nature and Countermovement were identified as strong predictors of movement duration. The current paper explores the potential that might be realized with deeper understanding from such an integration, by focusing on techniques that put together sentiment and thematic analyses with machine learning. While some structural variables-legal frameworks and counter-movement dynamics-played the leading role in determining movement longevity, sentiments played a secondary but complementary role. This project contributes to the greater knowledge of digital activism and gives practical insights into how to have effective and sustained social change.

TABLE OF CONTENTS

Chapter 1 Introduction	1
1.1 Background and Motivation	1
1.2 Purpose and Goal of the Project	2
1.3 Organization of the Report	2
Chapter 2 Research Literature Review	4
2.1 Introduction	4
2.2 Sentiment Analysis in Social Media	4
2.2.1 Approaches to Sentiment Analysis	4
2.3 Challenges in Sentiment Analysis	5
2.4 Explainable AI (XAI) and Its Role in Sentiment Analysis	6
2.4.1 Applications of SHAP in Social Movement Analysis	6
2.5 Predictive Modeling of Social Movement Lifespan	6
2.5.1 Machine Learning and Deep Learning Models	7
2.6 Multimodal Sentiment Analysis and Future Directions	7
2.6.1 Addressing Cultural and Linguistic Biases	7
2.7 Summary	8
Chapter 3 Methodology	9
3.1 System Design	9
3.2 Hardware and Software Components	10
3.3 Predictive Models	11
3.3.1 Meta-Model	14
3.4 Dimensionality Reduction and Clustering	15
3.5 Hardware and Software Implementation	15
Chapter 4 Investigation/Experiment, Result, Analysis and Discussion	16
4.1 Investigation and Experiment	16
4.2 Results and Analysis	17
4.2.1 Correlation Analysis	17
4.2.2 Dimensionality Reduction	21
4.2.3 Clustering Analysis	25
4.2.4 Feature Importance	31
4.2.5 Model Evaluation	33
4.2.6 Explainability with SHAP	35
4.3 Discussion	37
Chapter 5 Impact of the project	39
5.1 Impact of This Project on Societal, Health, Safety, Legal, and Cultural Issues	39
5.2 Impact of This Project on Environment and Sustainability	39
5.3 Technological Impact	40

5.4 Educational Impact	41
5.5 Economic Impact	41
5.6 Summary of Impacts	42
Chapter 6 Project Planning and Budget	43
6.1 Project Timeline and Phases	43
6.2 Resource Requirements	44
6.3 Human Resources	45
6.4 Budget Breakdown	46
6.5 Risk Management	47
6.6 Summary	47
Chapter 7 : Complex Engineering Problems and Activities	48
7.1 Complex Engineering Problems (CEP)	48
7.2 Complex Engineering Activities (CEA)	49
Chapter 8 Conclusions	51
8.1 Summary	51
8.2 Limitations	52
8.3 Future Improvements	52
References	54

LIST OF FIGURES

Figure 4.1: Correlation Matrix	20
Figure 4.2: PCA Scatter Plot	22
Figure 4.3: t-SNE Visualization	24
Figure 4.4: K-Means Clustering Results	27
Figure 4.5: DBSCAN Clustering Results	28
Figure 4.6: Gaussian Mixture Clustering Results	29
Figure 4.7: Feature Importance (Random Forest)	31
Figure 4.8: SHAP Summary Plot	35
Figure 4.9: SHAP Dependence Plot	36
Figure 6.1: Gantt Chart	44

LIST OF TABLES

Table 3.1: Tools and Technologies	10
Table 4.1: Correlation Analysis	20
Table 4.2: Importance of Features	31
Table 4.3: Model Evaluation	33
Table 6.1 Budget Summary	46
Table 7.1: Complex Engineering Problem Attributes	48
Table 7.2: Complex Engineering Problem Activities	49

Chapter 1 Introduction

1.1 Background and Motivation

In the era of digital transformation, people are feeling more comfortable to express opinions on social medias. Platforms such as Facebook, Instagram, Twitter are playing significant roles in social movements. Because of this reason, geographical barrier is not a problem now. Social media platforms are enabling global participation. Powerful Movements like #FreePalestine, #MeToo and #BlackLivesMatter are showing the power of social media in amplifying marginalized voices, fostering solidarity, and mobilizing public sentiment [1][2].

The collective power of public opinion on platforms such as Twitter, Facebook can successfully influence any social movement. For example, the #MeToo movement is inspiring millions of people in this planet to share their experiences of sexual harassment without hesitation. In the same way, #BlackLivesMatter has become a global rallying cry against police brutality demonstrating the capacity of social media to sustain public discourse[1][3].

The rise of social media-driven movements has introduced new challenges in understanding the dynamics in online activism. Social movements are becoming popular in online. It has become more challenging to understand the influence of social media platforms in these movements. Predicting the lifespan of these movements has become an important area of research which is providing valuable insights for activists, policymakers, and organizations to optimize their strategies and ensure long-term impact [4][5]. Popular social movements helped us to understand the importance of social media platforms. Social medias have capabilities to bring positive changes in society. Because of this reason, it has become essential to do more research on it.

Sentiment analysis is a subfield of Natural Language Processing (NLP). It helps us to analyze public opinion by categorizing text into positive, negative, or neutral sentiments. When coupled with thematic analysis, which identifies dominant themes within discussions, sentiment analysis provides a deeper understanding of the dynamics that shape public engagement and movement longevity. This research leverages these analytical techniques to address the complexity of

understanding the persistence of social media movements, laying the groundwork for actionable insights into digital activism.

1.2 Purpose and Goal of the Project

In this project we have tried to analyze and predict the lifespan of social media-driven movements. Sentiment and thematic analysis are done because of this purpose. The project wants to identify important factors such as sentiment trends and thematic relevance. Sentiment trends and thematic relevance influence the duration and effectiveness of these movements. Because of this reason, researchers develop predictive models.

It is an innovative project as it is successfully incorporating sentiment and thematic analysis with machine learning techniques in order to predict movement lifespans. Further novelty is brought in by using thematic scores, which quantify the prominence of certain themes, like political or legal relevance, in the texts. For example, the novelty in this research resides in the fact that most of the previous approaches do not consider a holistic approach, like taking a qualitative-quantitative assessment regarding a particular movement and providing actionable insights to all stakeholders.

Our project addresses practical challenges in data collection and analysis. In this project, we have used publicly available data. We think that it ensures ethical compliance while providing a scalable framework for future studies and research works. Our project may help policy makers, human rights activists and researchers.

1.3 Organization of the Report

This report is organized into eight chapters to provide a comprehensive overview of the project, its methodologies, and its outcomes.

- **Chapter 2** discusses a detailed literature review which covers existing studies which are related to social media analytics, sentiment analysis, and thematic analysis. It also demonstrates the gaps in the literature that this project wants to address.
- **Chapter 3** discusses the problem formulation and the project's objectives which includes the primary and sub-objectives, the scope of the study, and the methodologies employed.
- **Chapter 4** outlines the detailed methodology used in this project, including data collection, preprocessing, feature extraction, model development, and evaluation.
- **Chapter 5** describes the impacts of the project, emphasizing its significance in understanding social media dynamics and its potential contributions to digital activism.
- **Chapter 6** is focusing on our project planning and budgeting. It is detailing the resources, timeframes, and tools used in the implementation of the project.
- **Chapter 7** discusses on the complex engineering problems and activities involved in the project, highlighting the challenges faced and the innovative solutions devised.
- **Chapter 8** has concluded this report by summarizing the findings, discussing their implications, and suggesting directions for future work.

Each chapter builds upon the previous one to present a cohesive narrative of the project's conceptualization, execution, and outcomes.

Chapter 2 Research Literature Review

2.1 Introduction

With the rapid growth of social media platforms, the involvement of people in social movements has been completely changed, and it connects the globe in real time. People can share opinions, join discussions, and mobilize on behalf of causes they care about through Twitter, Instagram and Facebook, among other outlets, making them key tools for activism. Movements such as #MeToo and #BlackLivesMatter are examples of the power of social media to amplify marginalized voices and galvanize public opinion [1][2].

With the increase in digital activism, scholars develop an increasing need to qualitatively analyze and, preferably forecast tendencies in the development of social movement action. Currently, sentiment analysis, a part of NLP, has mushroomed into one of the main approaches for analyzing mass public opinion through social networking. It can enable estimation, using predictive modelling techniques, even with regard to movement lifetimes from sentiment and thematic relevance. It covers main methodologies, challenges, and novelties in the fields of sentiment analysis, XAI, and predictive modeling contributing to the analysis of social movements.

2.2 Sentiment Analysis in Social Media

Sentiment analysis classifies text into positive, negative, or neutral sentiments, offering a quantitative approach to understanding public opinion. With the exponential growth of social media, sentiment analysis has become indispensable for tracking public sentiment on topics ranging from politics to social movements [3]. By analyzing tweets, posts, and comments, researchers can assess how sentiments evolve over time, influencing movement longevity and effectiveness.

2.2.1 Approaches to Sentiment Analysis

These methodologies for performing sentiment analysis may be grouped into three broad kinds of approaches: machine learning, deep learning, and ensemble methods.

Machine Learning Approaches: Traditional approaches to sentiment analysis are based on algorithms like Naïve Bayes, Support Vector Machines, and logistic regression. Most of these models consider hand-crafted feature extraction techniques like TF-IDF and N-grams. These techniques have been relatively simple but quite successful for small datasets; most of them are committed to binary classification tasks as in [5][6].

Deep Learning Techniques: Sophisticated deep learning models such as RNNs, LSTMs, and transformer-based models like BERT have been developed. These models can automatically learn complicated patterns of language that are indicative of social media discourse. For example, Peng et al., 2023, show that BERT is particularly good at classifying sentiment within the #BlackLivesMatter movement because of its ability to understand the nuance of context-sensitive texts [3].

Ensemble Learning Methods: Ensemble methods, like Gradient Boosting Machines and Random Forest, combine several models to improve generalization and performance. These methods are very effective in analyzing complex datasets that include different sources of data, such as text and images. It will be able to make reliable sentiment analysis in all diversity of contexts [22][23].

2.3 Challenges in Sentiment Analysis

Sentiment analysis is successful in analyzing public sentiment but it is facing a lot of challenges when it is being applied to social movements. Unlike traditional sentiment classification which only focuses on polarity alone, proper analyzation social trends or movements obviously involves counter-movements, complex themes and legal narratives. #MeToo and #BlackLivesMatter are nuanced discussions that demand advanced techniques, such as NER and topic modeling to capture multifaceted interactions [7][9].

Furthermore, social media data can hardly be collected ethically because of the privacy restrictions. Web scrapers, which can be easily built using automation tools such as Bardeen AI, can scrape data from social media in an ethical manner if targeting only publicly available posts.

Nevertheless, this can result in incomplete datasets that might not capture all the variations in opinions from a movement [10].

2.4 Explainable AI (XAI) and Its Role in Sentiment Analysis

Explainable AI (XAI) increases the interpretability of complex machine learning models by showing the insight of how features influence predictions. This aspect of transparency is very crucial in social movement analysis, where understanding drivers of public sentiment are so vital for informed decision-making. Among the most used techniques from the field of XAI, one could note SHAP, or SHapley Additive exPlanations, a method to quantify the contribution of individual features to a model's predictions [7][8].

2.4.1 Applications of SHAP in Social Movement Analysis

SHAP values let researchers investigate the impact of factors such as positive sentiment and legal relevance on movement longevity. For instance, scholarship has used SHAP in identifying the most dominant features responsible for the persistence of social movements like #MeToo, therefore offering actionable insights to social activists and policymakers. The global and local interpretability offered by SHAP values help in the comprehensive understanding the dynamics of social trends shaping up in social media [7][10].

2.5 Predictive Modeling of Social Movement Lifespan

Predictive modeling is a powerful tool which estimate social media movement longevity. These models give meaningful foresight, drawing upon the trends in sentiment and the thematic relevance of subjects, thus providing insight into which factors ultimately sustain these movements over

time. The modelling could be trained using historical data from past movements, , driving real-time predictions on emerging ones [5][14].

2.5.1 Machine Learning and Deep Learning Models

Conventional machine learning models, like linear regression, measure the association of thematic features with the duration of movement. More advanced clustering models group movements with similar characteristics and facilitate insights into shared dynamics. Deep learning models, such as LSTM and Transformers, further improve the predictive performance by considering sequential properties of social media posts and monitoring sentiment changes [9][13].

2.6 Multimodal Sentiment Analysis and Future Directions

The integration of multimodal data sources, such as images, videos, and text, could enrich sentiment analysis. Movements like #BlackLivesMatter use visual media to amplify their messages, underscoring the need for tools that analyze multiple modalities. By integrating textual and visual elements, researchers can develop a more holistic understanding of movement dynamics [13][17].

2.6.1 Addressing Cultural and Linguistic Biases

Most of the models developed for sentiment analysis face difficulties in generalizing across diverse linguistic and cultural contexts. Many of these models tend to be biased toward training data from specific regions, limiting their usability for global movements. For example, most models trained on mainly English data would poorly classify movements in non-English-speaking areas. Future

research should focus on developing culturally adaptive models sensitive to linguistic nuances for equal representation [9][18].

2.7 Summary

This chapter introduced methodologies and challenges related to sentiment analysis, Explainable AI, and predictive modeling in the context of social movements. Sentiment analysis forms the basis for public opinion estimation, while predictive models provide valuable insights into the sustainability issue of a movement. Methods such as SHAP improve transparency in models so that confidence can be instilled among its stakeholders to act upon. However, there is a need to address critical challenges related to cultural bias and lack of diversity in data for further advancements within this area.

Chapter 3 Methodology

3.1 System Design

The system design for this project is structured into a modular architecture, as shown in the flowchart (Figure 3.1). Each module addresses specific tasks required for predicting the lifespan of social media movements, leveraging sentiment and thematic analysis combined with predictive modeling.

- **Data Collection Module:** A custom web scraper built using Bardeen AI gathers public Facebook posts related to social movements. These posts form the raw dataset for subsequent analysis [15].
- **Data Preprocessing Module:** This module cleans and preprocesses the collected data by removing special characters, emojis, and irrelevant content while retaining meaningful information for analysis [18][19].
- **Feature Extraction Module:** Sentiment and thematic scores are calculated for each post, and Google Trends data is utilized to determine the peak duration of movements [17].
- **Predictive Modeling Module:** Various machine learning models are trained on the processed dataset to predict the lifespan of movements [24][25].
- **Explainability and Analysis Module:** SHAP (SHapley Additive exPlanations) provides interpretability for model predictions by highlighting the contributions of individual features [7][12].
- **Deployment Module:** The predictive model was deployed on two platforms: a web application and an iOS application.

This modular system ensures scalability and adaptability, enabling the integration of additional features or models in future iterations.

3.2 Hardware and Software Components

Software-based tools were used to implement this project. It did not require any additional hardware components. The tools and technologies employed in this project are listed in Table 3.1.

Tool	Functions	Other Similar Tools (if any)	Why Selected this Tool
Python	Core programming language for analysis and modeling	R, MATLAB	Widely used, extensive library support
Google Colab	Cloud-based platform for training and visualization	Jupyter Notebook	Free, no hardware requirements
scikit-learn	Provides machine learning algorithms and metrics	TensorFlow, PyTorch	Easy to use, supports classical ML algorithms
SHAP	Explains model predictions and feature contributions	LIME	Offers global and local interpretability
Pandas, NumPy	Data manipulation and handling	None	Efficient for handling large datasets
Matplotlib, Seaborn	Data visualization tools	Plotly	Simple and effective for statistical visualizations
Bardeen AI	Web scraper for collecting Facebook posts	Selenium, BeautifulSoup	Customizable and efficient for targeted scraping

Table 3.1: Tools and Technologies

3.3 Predictive Models

The project involved the implementation and comparison of multiple machine learning and deep learning models for predictive analysis. These included:

1. Linear Regression:

- A baseline model is used to establish a starting point for predictions [20].
- Simple and interpretable, but limited in capturing non-linear relationships.

2. Support Vector Regressor (SVR):

- It is effective for handling small datasets and high-dimensional spaces [21].
- Kernel-based method suitable for non-linear relationships.

3. Random Forest Regressor:

- An ensemble model that leverages multiple decision trees to improve accuracy [22].
- Robust to overfitting and effective for feature importance analysis.

4. Gradient Boosting Regressor (GBR):

- Boosting algorithm that builds models sequentially to minimize prediction errors [23].
- Provides improved performance for structured data.

5. XGBoost:

- An advanced gradient boosting algorithm which is famous for its efficiency and scalability [24].
- Optimized for high-speed execution and competitive modeling.

6. **LightGBM:**

- Gradient boosting framework designed for large datasets and high-dimensional data.
- Faster than traditional boosting methods due to its leaf-wise splitting technique.

7. **Refined Deep Learning Models:**

- Implemented for capturing complex feature interactions by using neural networks [26].
- Includes optimized architectures for enhanced predictive capabilities.

8. **Simple Averaging Ensemble:**

- It combines predictions from multiple models by averaging their outputs [27].
- It reduces variance and improves generalizability.

9. **Weighted Ensemble:**

- Assigns different weights to models based on performance metrics [28].
- Balances contributions from individual models for improved accuracy.

10. **Optimized Weighted Ensemble:**

- Refines the weighted ensemble by optimizing the weight distribution [29].
- Maximizes predictive performance by reducing bias.

11. **Meta-Model:**

- A stacking approach that uses predictions from base models as inputs for a higher-level model [27].
- Demonstrated the best performance in terms of Mean Absolute Error (MAE) and Mean Squared Error (MSE).

These models were systematically tuned using Grid Search and Randomized Search to optimize hyperparameters and ensure robust performance [29].

3.3.1 Meta-Model

The Meta-Model is an innovative framework that combines predictions from a Gradient Boosting Regressor (GBR) and a Deep Learning (DL) model using Ridge Regression as a final layer. This hybrid approach leverages the strengths of GBR (capturing non-linear relationships) and DL (handling complex patterns) to improve predictive accuracy.

Key Features:

- **Model Integration:** Combines outputs from GBR and DL models, compensating for individual weaknesses and enhancing generalization.
- **Meta-Features:** Uses predictions from base models as input, allowing the Ridge Regression layer to refine the final output.
- **Regularization:** Ridge Regression reduces overfitting, ensuring robust predictions.

Novelty:

- **Hybrid Error Compensation:** Aggregates multiple models to balance biases and improve performance.
- **Scalable and Transparent:** Offers a simple, interpretable framework suitable for large-scale applications.
- **Unique Application:** Tailored for predicting social media movement lifespans, combining sentiment and thematic features for actionable insights.

This meta-model achieves superior performance (lower MAE/MSE), making it ideal for accurate, interpretable predictions in real-world applications.

3.4 Dimensionality Reduction and Clustering

- **Dimensionality Reduction:**

- Techniques such as Principal Component Analysis (PCA) and t-SNE were employed to simplify high-dimensional data while retaining significant information. These methods improved data visualization and enhanced model performance by reducing noise [30][34].

- **Clustering:**

- To identify patterns and groupings within the data Algorithms like k-means clustering were applied which is providing additional insights into the characteristics of social media movements [31][35].
-

3.5 Hardware and Software Implementation

This project involved the following implementation steps:

1. **Data Collection:**

- Bardeen AI was used to develop a custom web scraper for collecting Facebook posts which are related to social movements. This scraper targeted posts which contains specific hashtags and extracted only publicly available data to adhere to ethical guidelines [15].

2. **Data Preprocessing:**

- Preprocessing was done by using Python libraries such as Pandas and NumPy. By removing special characters, emojis, and irrelevant content, text was cleaned. We

applied Tokenization and lowercasing to standardize the data for the analysis [18][19].

3. Sentiment and Thematic Analysis:

- Sentiment scores were computed using pre-built Python libraries like VADER and TextBlob [33]. Thematic scores were derived by analyzing the frequency of theme-specific keywords, which were manually curated based on common terms relevant to legal, political, and social movements.

4. Model Development:

- Each model was implemented and tuned on Google Colab. By using metrics such as MAE, MSE, and R^2 , the performance of individual was evaluated. For improving predictions, Ensemble techniques were employed.[24][25][26].

5. Explainability:

- SHAP values were used to interpret model predictions, providing insights into feature importance and their contributions to movement lifespan predictions. Visualizations such as SHAP bar plots and beeswarm plots facilitated understanding of feature interactions [7][12].

6. Deployment:

- The predictive model was deployed on two platforms: a web application and an iOS application.
- Alongside the predictive model, we also deployed 2 prototype models that can run real-time sentiment analysis on a given text input and find the duration and similar movements based on the sentiment scores.

Chapter 4 Investigation/Experiment, Result, Analysis and Discussion

4.1 Investigation and Experiment

The analysis has explored the relationships between various features of social movements and their impact on movement duration. It involved a series of experiments, including:

1. **Correlation Analysis:** To identify linear relationships between features and the target [17][5].
 2. **Feature Importance:** To rank features based on their contribution to the target prediction [22][5].
 3. **Dimensionality Reduction:** To simplify high-dimensional data while retaining significant information [30][34].
 4. **Clustering:** To uncover hidden patterns and groupings within the dataset [31][35].
 5. **Model Development and Evaluation:** To predict movement duration using various machine learning models [22][24][25][28][29].
 6. **Explainability:** To interpret the decisions of predictive models using SHAP (SHapley Additive exPlanations) [7][8][12].
-

4.2 Results and Analysis

The following section discusses and analyses all the results we obtained from our multidimensional investigation of the dataset we accumulated.

4.2.1 Correlation Analysis

The correlation analysis, presented in Table 4.1 and Figure 4.1, investigated the linear relationships between the features and the duration of movements. Pearson and Spearman correlation coefficients were calculated for each feature [17], [5].

Key Observations:

Counter-Movement:

- Highest positive, yet weak, correlation with the duration of movement; this would suggest that movements that have to do with strong counter-movements last longer because of continuous public attention and debates the movement creates.
- The counter-movements serve to catalyze the primary movements through constant public interaction and controversy.

Legality Nature:

- It is a fair positive correlation with movement duration, indicating that it could be an important predictor.
- Movements that have higher legality scores are likely to be very long-lived since they gain more public confidence and wider involvement. This again stresses the point that structural and systemic support is crucial in sustaining movements.

Sentiment Features: Positive, Negative, and Neutral

- Show little or no linear relationship with movement duration, indicating that non-linear relationships are present.
- Positive Sentiment makes a minor contribution to the duration of movement, technically enhancing the magnitude of the impact of structural variables such as legality.
- In contrast, Neutral and Negative Sentiments demonstrate weak correlations and hence their roles could be more contextual and secondary in the influence on movement dynamics.

Cultural Relevance:

- Shows a moderately positive correlation with the movement duration, but is less important than structural factors such as legality and counter-movement.
- While cultural alignment supports, it is not a leading determinant in movement longevity.

Political Nature:

- Exhibits negligible correlation, reflecting its limited direct impact on the longevity of movements in the dataset.

Insights from the Correlation Matrix:

- **Weak Overall Correlations:** Most of the overall weak correlations indicate the complexity of prediction to be performed in movement duration. The simplest linear model might be insensitive in capturing the underlying dynamics; hence, advanced modeling might be required.
- **Structural vs. Sentiment Features:** The structural features of Legality Nature and Counter-Movement come out to be stronger predictors compared to the sentiment-based features. This shows the need for the incorporation of broader contextual and systemic factors into predictive models.
- **Non-Linearity:** The minimal correlations of sentiment features with the movement duration suggest that non-linear relationships dominate the data. Advanced models, such as ensemble methods or neural networks, are required to capture these complexities. [5] [22] [24] [25]

The results from the correlation analysis are quite promising for laying the groundwork regarding dependencies between features and helping further feature selection in modeling toward arriving at good accuracy of predictions. These are those insights that bring in the aspect of taking a holistic approach that integrates both the structural and the sentiment-based factor in analyzing a robust means.

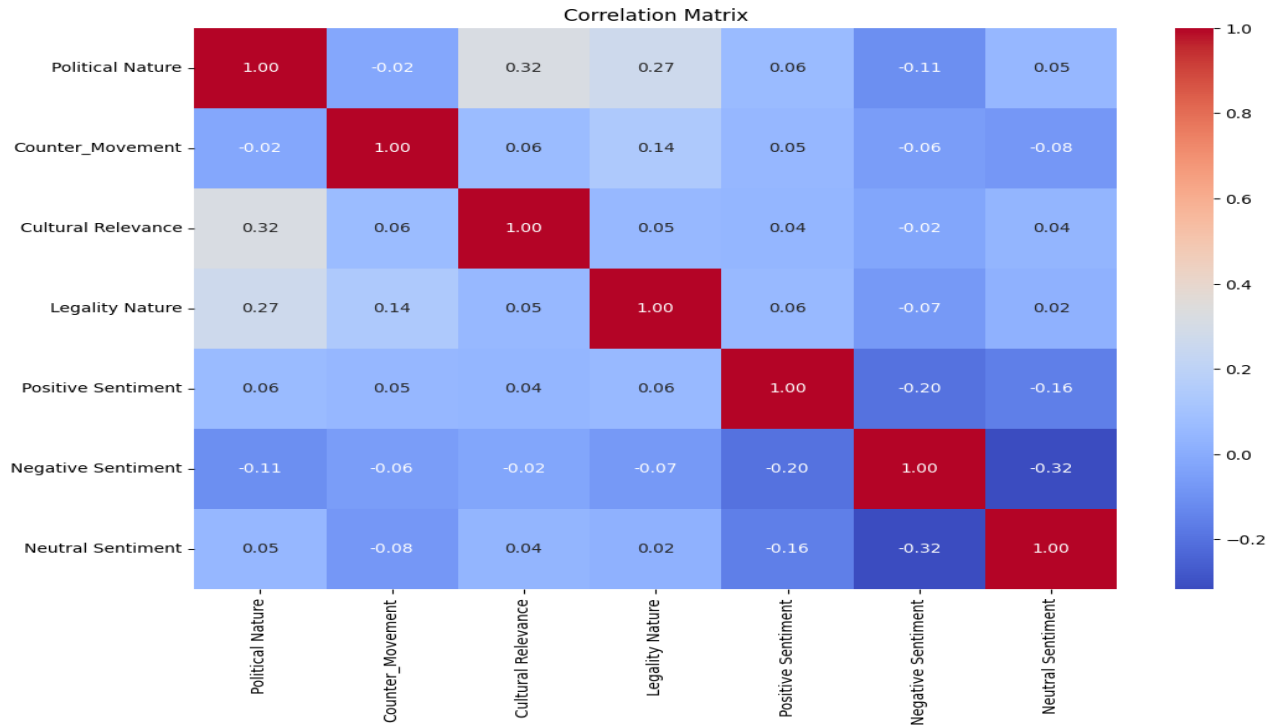


Figure 4.1: Correlation Matrix

Feature	Pearson	Spearman
Political Nature	-0.0188	0.0248
Counter_Movement	0.1717	0.0456
Cultural Relevance	0.0228	0.0446
Legality Nature	0.1075	0.1522
Positive Sentiment	-0.0338	0.0062
Negative Sentiment	-0.0513	-0.0501
Neutral Sentiment	0.0530	0.0266

Table 4.1: Correlation Analysis

4.2.2 Dimensionality Reduction

Dimensionality reduction techniques which include Principal Component Analysis (PCA) and t-Distributed Stochastic Neighbor Embedding (t-SNE), were used for exploring and visualizing the high-dimensional dataset in a lower-dimensional space. These methods provided insights into the data's structure and relationships. It is focusing on variance and clustering patterns. [30][34]

- **Principal Component Analysis (PCA):** PCA was applying for reducing the dataset's dimensionality while preserving features which contribute most to variance. Important findings include:
 - **Identification of Principal Components:** PCA has successfully identified two principal components which accounted for the maximum variance in the dataset, as illustrated in the scatter plot (**Figure 4.2**). These components condensed the most significant attributes of the data while minimizing redundancy.
 - **Scatter Plot Analysis:** The PCA scatter plot is representing data points projected onto the two principal components. Every data point corresponds to a movement, with colors reflecting their duration.
 - **Complex Distribution:** The PCA visualization revealed no distinct clusters, suggesting that the dataset's features do not segregate into well-defined categories. This complexity underscores the challenge of predicting movement longevity based solely on linear relationships.

This analysis highlights the nuanced nature of the dataset and underscores the importance of integrating both linear and non-linear modeling approaches to capture the complexities within the data.

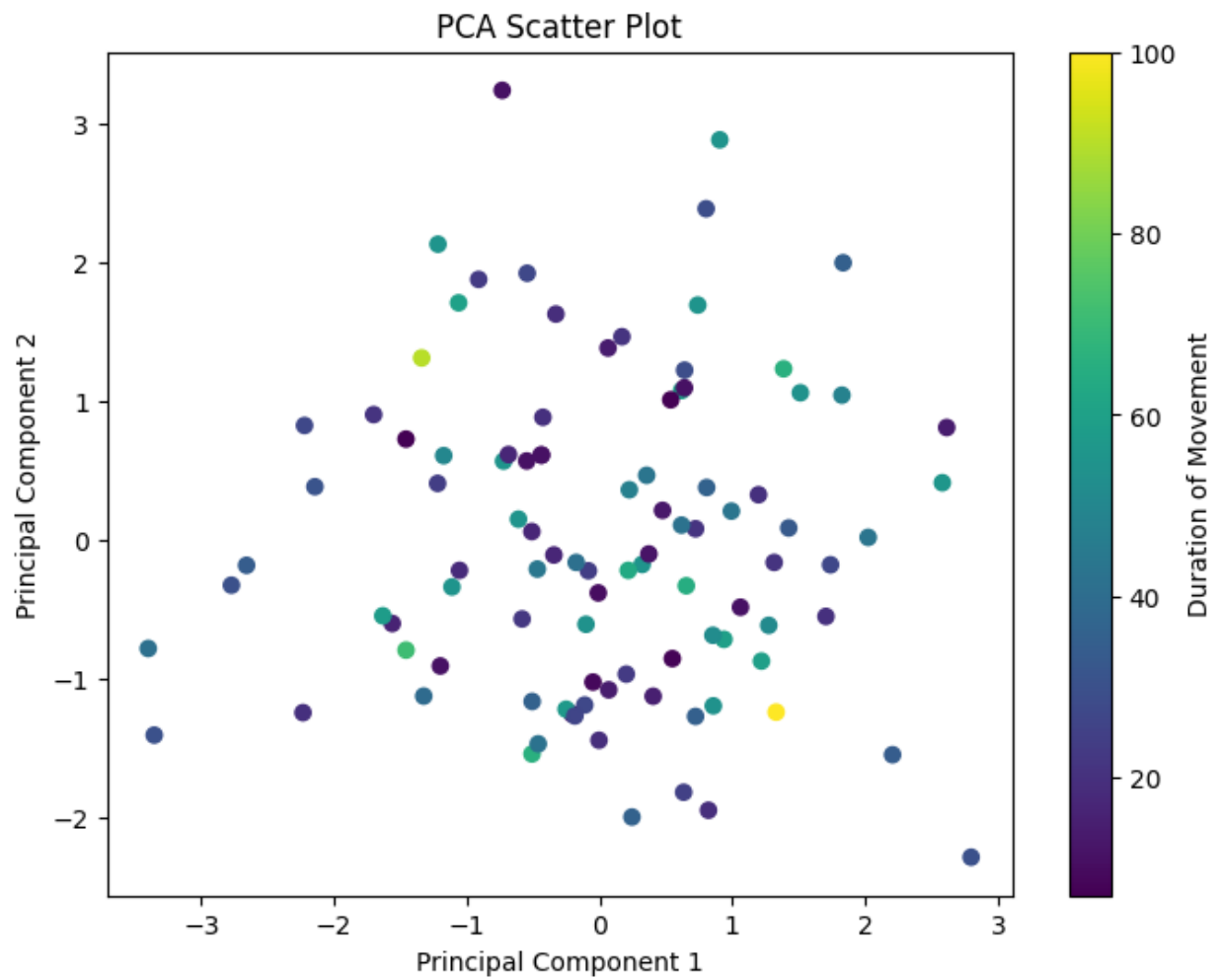


Figure 4.2: PCA Scatter Plot

- **t-SNE Visualization:** t-SNE is a non-linear dimensionality reduction technique. It was employed to supplement PCA by uncovering local relationships among data points. This approach provided deeper insights into the dataset's structure, especially where non-linear patterns dominate.
 - **Clarity of Local Relationships:** Unlike PCA, which emphasizes global variance, t-SNE excels in capturing finer, local structures. In Figure 4.3, movements with similar feature characteristics are grouped more closely which is highlighting localized patterns within the dataset.
 - **Non-Linear Patterns:** The t-SNE visualization reveals intricate non-linear relationships that are not apparent through linear techniques like PCA. This indicates the presence of subtle and complex feature interactions that contribute to the dataset's overall structure.
 - **Grouping Analysis:** While some movements with shared characteristics form localized clusters in the t-SNE plot, significant overlap and dispersal remain. This persistence of overlap reinforces the dataset's complexity and the challenges associated with clearly separating movement characteristics.

The findings from the t-SNE analysis (**Figure 4.3**) show the necessity for advanced modeling techniques capable of handling these non-linear and overlapping patterns for achieving more accurate predictions.

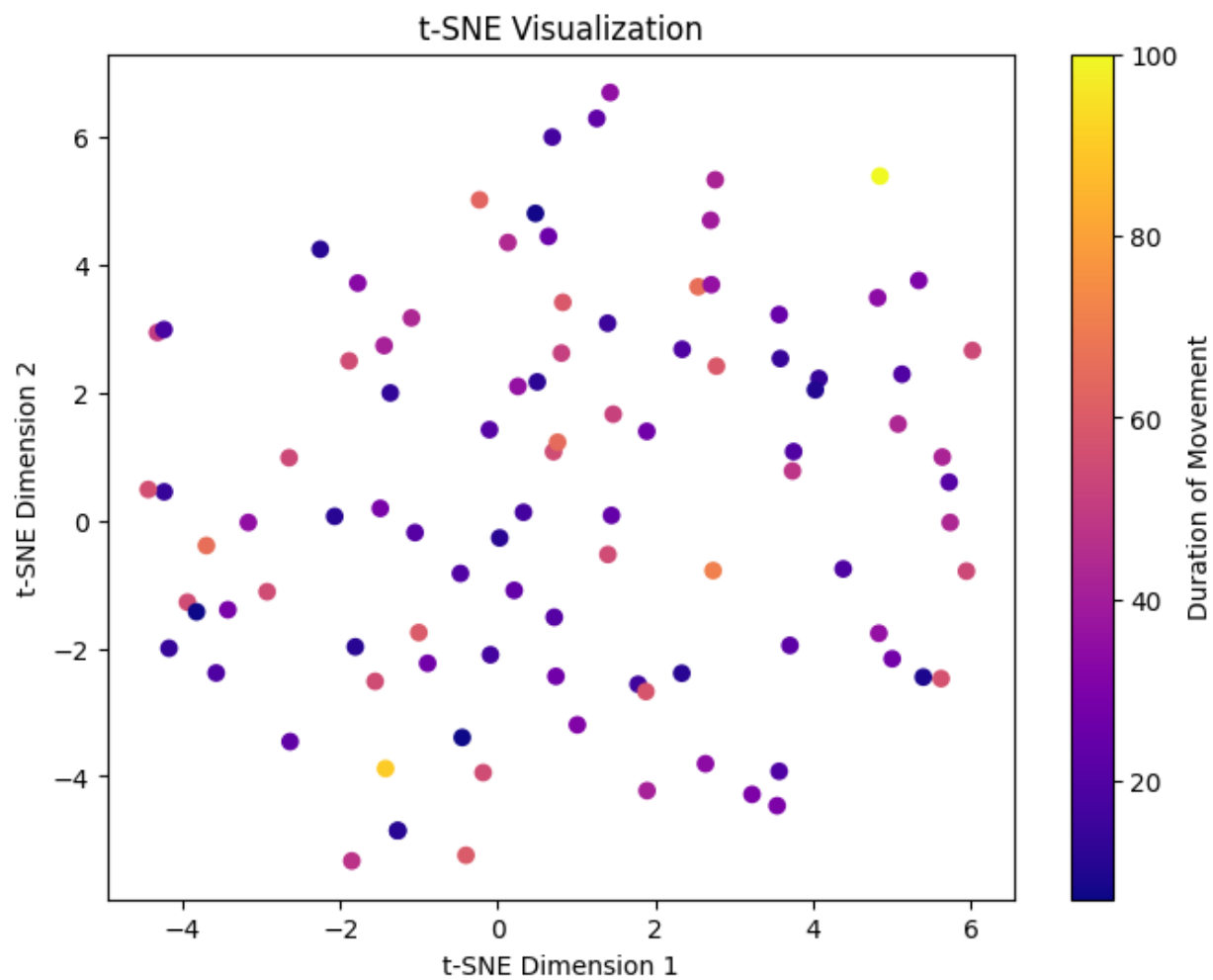


Figure 4.3: t-SNE Visualization

Insights:

- **Need for Advanced Models:** The lack of distinct clusters in both PCA and t-SNE visualizations is showing the complexity of the feature space. This calls for advanced machine learning models capable of capturing non-linear interactions and subtle patterns, as simpler models may fail to adequately capture the intricacies of the data.
- **Complementary Roles of PCA and t-SNE:** PCA provides a macroscopic perspective by analyzing global variance, while t-SNE focuses on localized relationships within the data. Together, these techniques offer complementary insights, helping to better understand the dataset's overall structure and localized nuances.
- **Challenges in Prediction:** The overlapping data points in both PCA and t-SNE visualizations reflect the complexity of the dataset. These overlaps indicate that the relationships between features are neither linear nor easily separable, posing significant challenges for predictive modeling, especially for tasks like predicting movement duration.

These dimensionality reduction techniques have provided essential insights into the dataset's structure, guiding the development of suitable predictive models. The findings highlight the importance of combining linear and non-linear methods to effectively address the complexities inherent in analyzing social media movement dynamics. [5] [30] [34]

4.2.3 Clustering Analysis

Clustering techniques were used for grouping similar movements based on their features and uncover patterns within the data and the analysis utilized three algorithms—K-Means, DBSCAN, and Gaussian Mixture Models (GMM)—each is offering distinct insights into the dataset's structure and relationships. [31][35]

1. **K-Means Clustering:** K-Means is a partition-based clustering method. It groups data points into clusters by minimizing intra-cluster variance.
 - **Outcome:** This dataset was divided into three clusters. It was visualized in **Figure 4.4**, based on PCA-reduced components.
 - **Observations:**
 - Significant overlap between clusters was evident. It is indicating weak separability within the data.
 - The lack of clear boundaries shows the complexity and interconnectedness of social movement characteristics. It is suggesting that the features of our dataset do not form distinct groups.

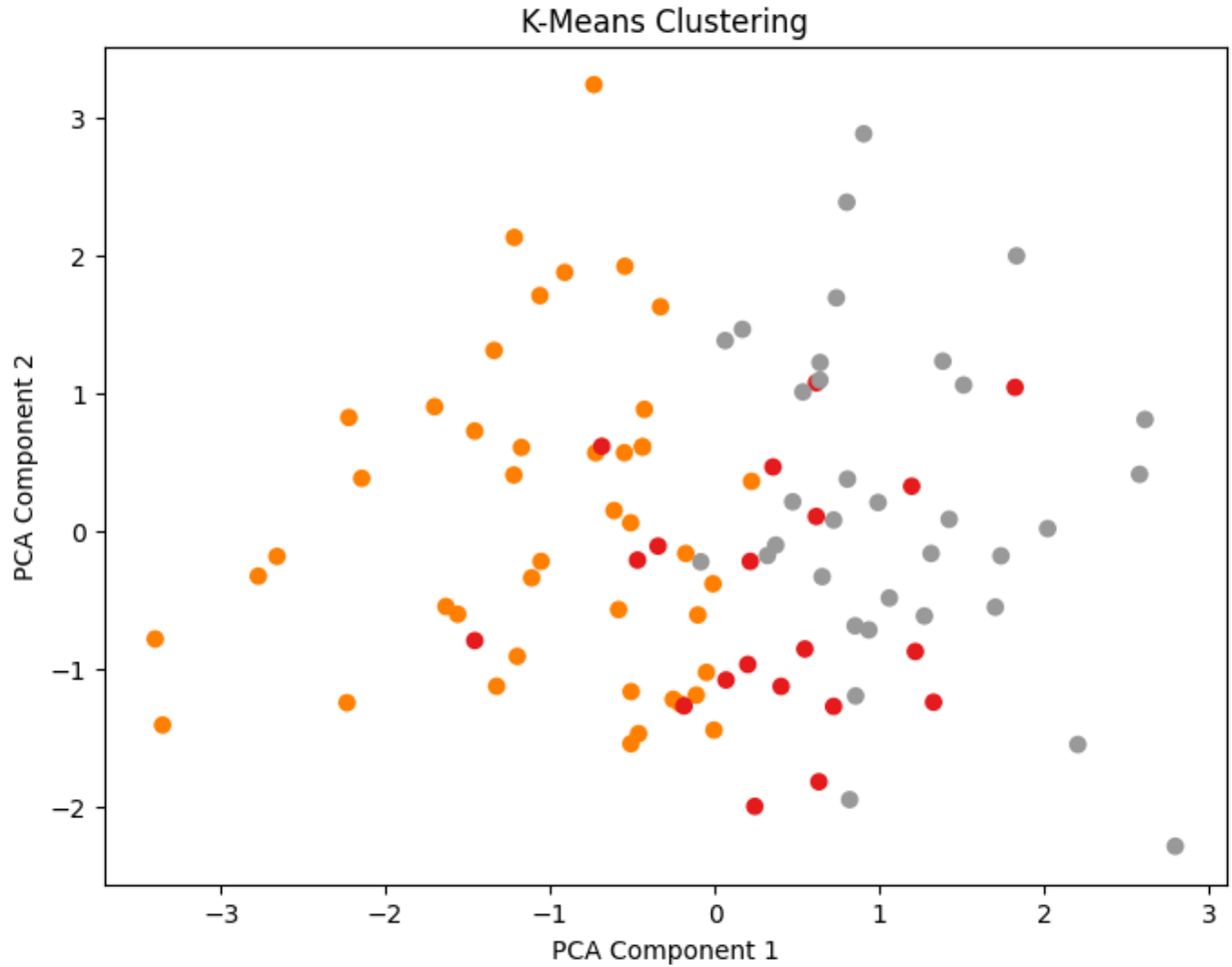


Figure 4.4: K-Means Clustering Results

2. **DBSCAN Clustering:** DBSCAN is a density-based clustering method. It identifies dense regions and isolates outliers.

- **Outcome:** It successfully identified dense clusters and flagged sparse data points as outliers. It is shown in **Figure 4.5**.
- **Observations:**
 - It revealed clusters that K-Means failed to capture, particularly those with irregular shapes.
 - DBSCAN excelled at handling non-linear relationships, making it a valuable tool for datasets with high variability.
 - The identification of outliers highlights movements with unique or atypical feature combinations, providing insights into rare or exceptional cases.

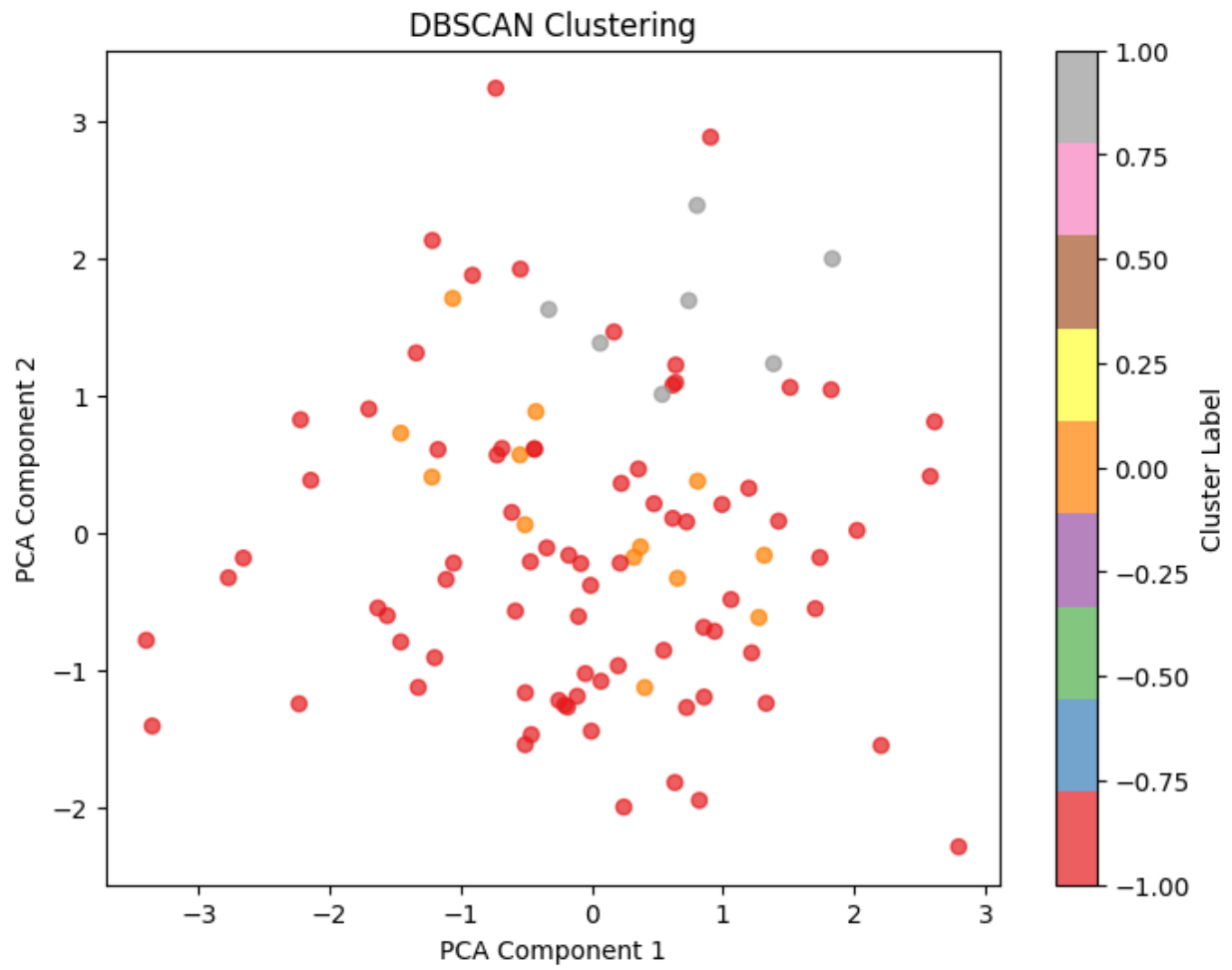


Figure 4.5: DBSCAN Clustering Results

3. **Gaussian Mixture Models (GMM):** GMM is a probabilistic approach. It has assigned data points to clusters based on their likelihood of belonging to each cluster.

- **Outcome:** GMM generated overlapping clusters with soft boundaries. It is seen in **Figure 4.6**.
- **Observations:**
 - GMM captured the dataset's complexity. It was done by allowing data points to belong to multiple clusters with varying probabilities.
 - Flexibility of it in modeling overlapping spaces makes it particularly effective for datasets with nuanced relationships and no clear separations.

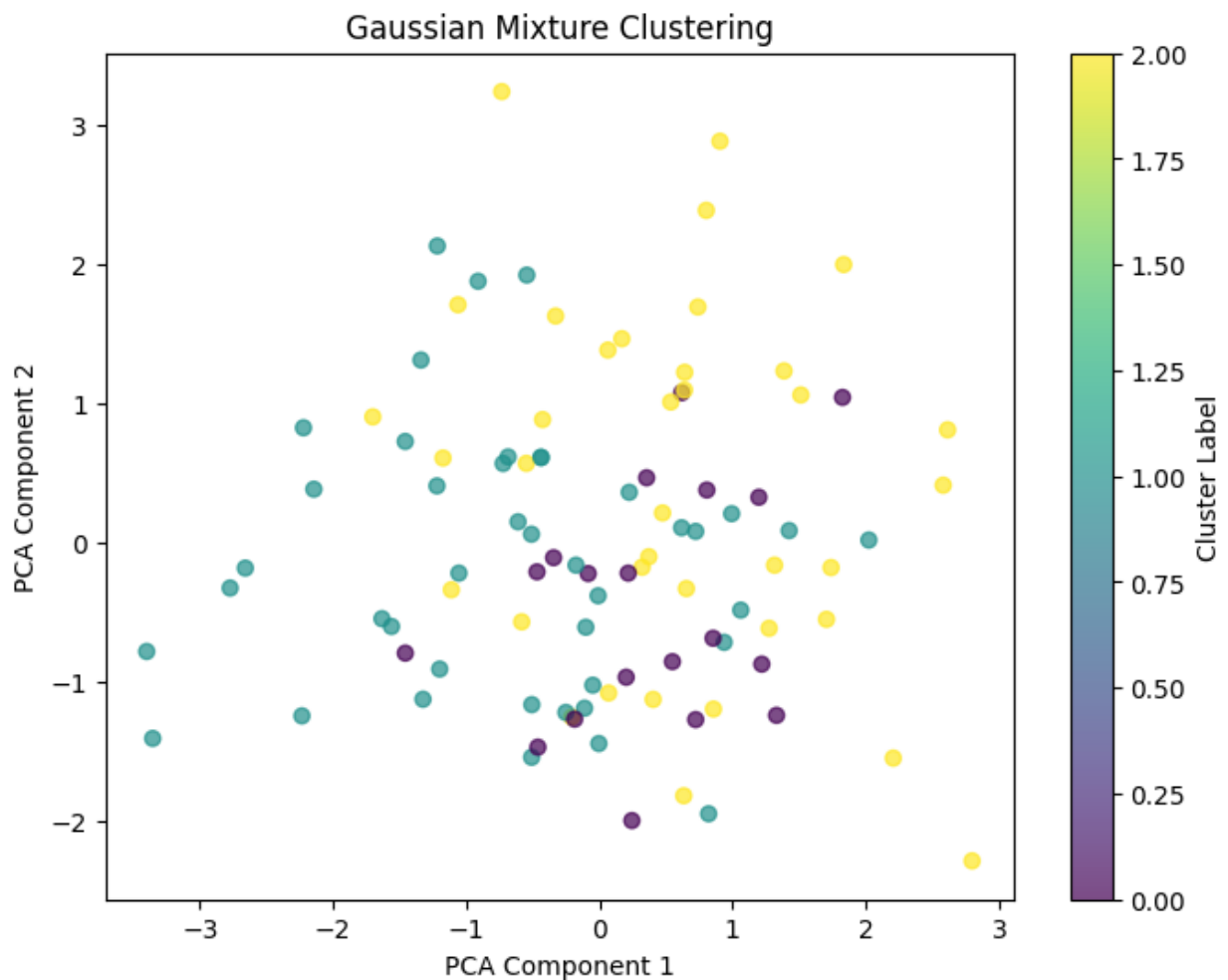


Figure 4.6: Gaussian Mixture Clustering Results

Insights:

- **Dataset Complexity:** All clustering methods showed that the dataset is lacking well-defined clusters. It is confirming its high dimensionality and complexity. [5][31][35]
- **Method Comparisons:**
 - K-Means: Basic partitioning. But it has struggled with overlapping data points.
 - DBSCAN: It is absolutely best at identifying outliers. It is handling non-linear relationships which highlight unique movement features.
 - GMM: Most nuanced. It is capturing overlapping clusters and allowing soft memberships. It has made it highly adaptable for this dataset.
- **Relevance to Predictive Modeling:** The overlapping clusters across all methods emphasize the need for advanced predictive models that account for non-linear relationships and subtle feature interactions.

By integrating insights from these clustering techniques, the analysis provides a comprehensive understanding of the dataset's dynamics. These findings help us to learn the design of robust machine learning models. It underscores the importance of advanced clustering methods for exploring and modeling complex datasets like those used in social movement and trend analysis.

4.2.4 Feature Importance

Random Forest was used for ranking features based on their importance in predicting movement duration. Table 4.2 and Figure 4.7 are summarizing the results [22][5].

Rank	Feature	Importance
1	Legality Nature	0.20
2	Counter_Movement	0.18
3	Positive Sentiment	0.17
4	Political Nature	0.14
5	Neutral Sentiment	0.12
6	Negative Sentiment	0.11
7	Cultural Relevance	0.08

Table 4.2: Importance of Features

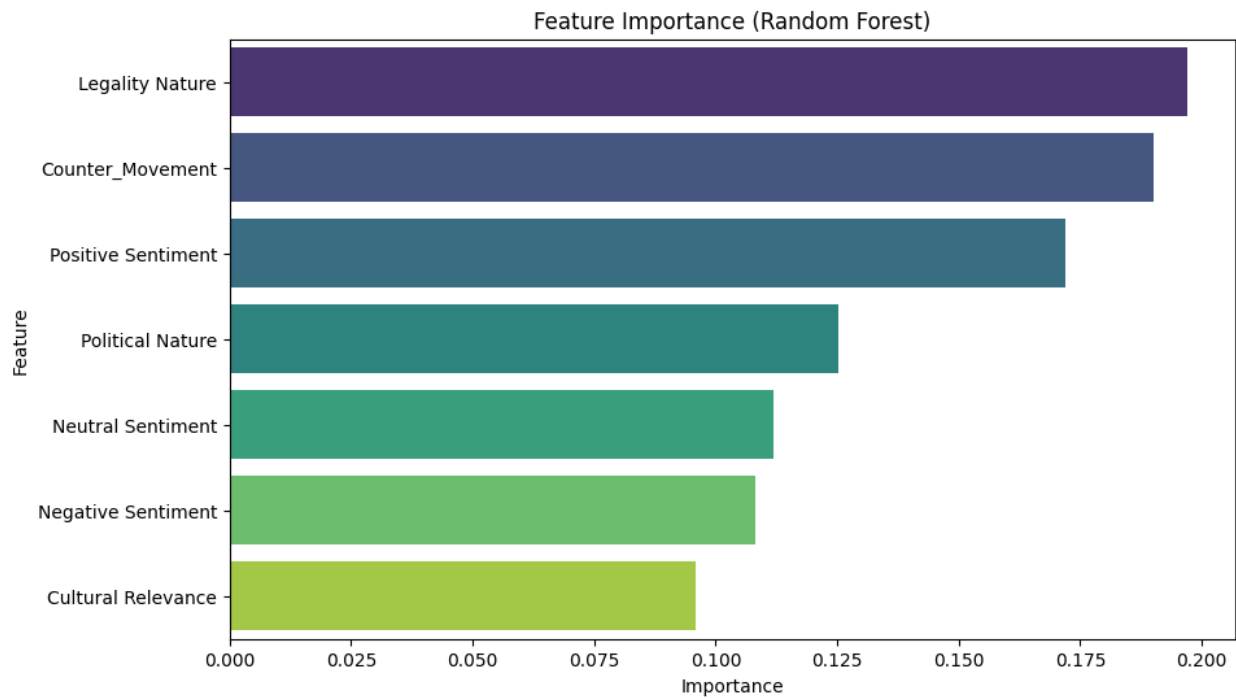


Figure 4.7: Feature Importance (Random Forest)

Insights:

- **Legality Nature** emerged as the most important feature, with the highest importance score of 0.20, It is suggesting that movements with higher legality scores tend to last longer. This could be attributed to the increased public trust and broader participation in legally sanctioned movements.
- **Countermovement** followed closely with an importance score of 0.18, indicating its strong influence. Counter-movements often prolong the duration of primary movements by generating sustained attention and ongoing interactions.
- **Positive Sentiment** scored 0.17, reflecting the role of public approval and optimism in movement longevity.
- Features such as **Neutral Sentiment** and **Negative Sentiment** showed moderate contributions, indicating their contextual relevance depending on the nature of the movement.
- **Cultural Relevance** ranked lowest with an importance score of 0.08, suggesting that while cultural alignment is a factor, it may play a more indirect or supportive role.
- The analysis highlights the critical role of **structural features** (e.g., legality and counter-movement dynamics) over sentiments in predicting movement duration.
- **Sentiment features**, while influential, likely interact with structural factors, amplifying or dampening their effects.
- The findings align with **SHAP** analysis, which also identified Legality Nature and Countermovement as the most impactful features. [7][12][8].

This feature importance analysis underscores the value of incorporating diverse feature types in predictive models, enabling robust predictions and actionable insights.

4.2.5 Model Evaluation

The predictive performance of various models was evaluated, with ensemble models demonstrating superior results [22][24][25][28][29]. The predictive performance of various models is presented in **Table 4.3**.

Model	MAE	MSE
Random Forest	16.8950	401.6708
SVR	16.8245	438.0761
Linear Regression	17.6166	423.7922
Refined Deep Learning	14.7762	358.2412
Simple Average Ensemble	14.4582	342.6778
Weighted Ensemble	14.4473	343.3249
Optimized Weighted Ensemble	14.4582	342.6819
Meta-Model	14.8298	313.4171

Table 4.3: Model Evaluation

Insights:

- **Performance Comparison**
 - **Superiority of Ensemble Methods:** Ensemble methods consistently outperformed individual models by aggregating predictions, which mitigated biases and enhanced robustness. This highlights the importance of leveraging diverse model architectures to capture complex relationships in predictive modeling.
 - **Meta-Model Excellence:** The Meta-Model has emerged as the most effective solution. It is combining advanced models like Gradient Boosting and Deep Learning with a Ridge Regression framework for error balancing. This innovative approach proved to be both novel and impactful.

- **Advantages of the Meta-Model**

- **Addressing Model Weaknesses:** By compensating for the limitations of individual models, the Meta-Model achieved the highest predictive performance, establishing it as the most suitable choice for deployment.
- **Adaptability and Scalability:** The Meta-Model's ability to reduce overfitting and generalize across diverse movement patterns underscores its suitability for real-world applications with high variability and complexity.

- **Practical Relevance**

- **Ensemble Importance:** These findings highlight the critical role of ensemble methods in capturing the intricate dynamics of social movements.
- **Success of Meta-Model:** The Meta-Model's performance demonstrates the value of integrating advanced feature engineering, non-linear modeling, and explainable AI techniques to derive actionable insights.

- **Conclusion**

The evaluation reinforces the significance of ensemble and meta-model approaches in predictive analytics, particularly for datasets characterized by high variability and interdependent features. The Meta-Model's adaptability and robust performance make it a valuable tool for analyzing complex social movement dynamics and informing strategic decision-making.

4.2.6 Explainability with SHAP

SHAP analysis gave us the opportunity to learn that Legality Nature and Countermovement were the most influential features which is significantly impacting predictions. Movements with higher legality scores or associated counter-movements tended to last longer. Sentiments (Neutral, Positive, Negative) showed smaller contributions which suggests that their impact is secondary to structural features [7][8][12]. Here we see, the SHAP summary plot highlighted individual feature contributions, while the dependence plot showed interactions, such as between Legality Nature and Political Nature, emphasizing the need for non-linear models. This analysis enhances model transparency, validates feature importance rankings, and supports actionable insights for planning and analyzing social movements.

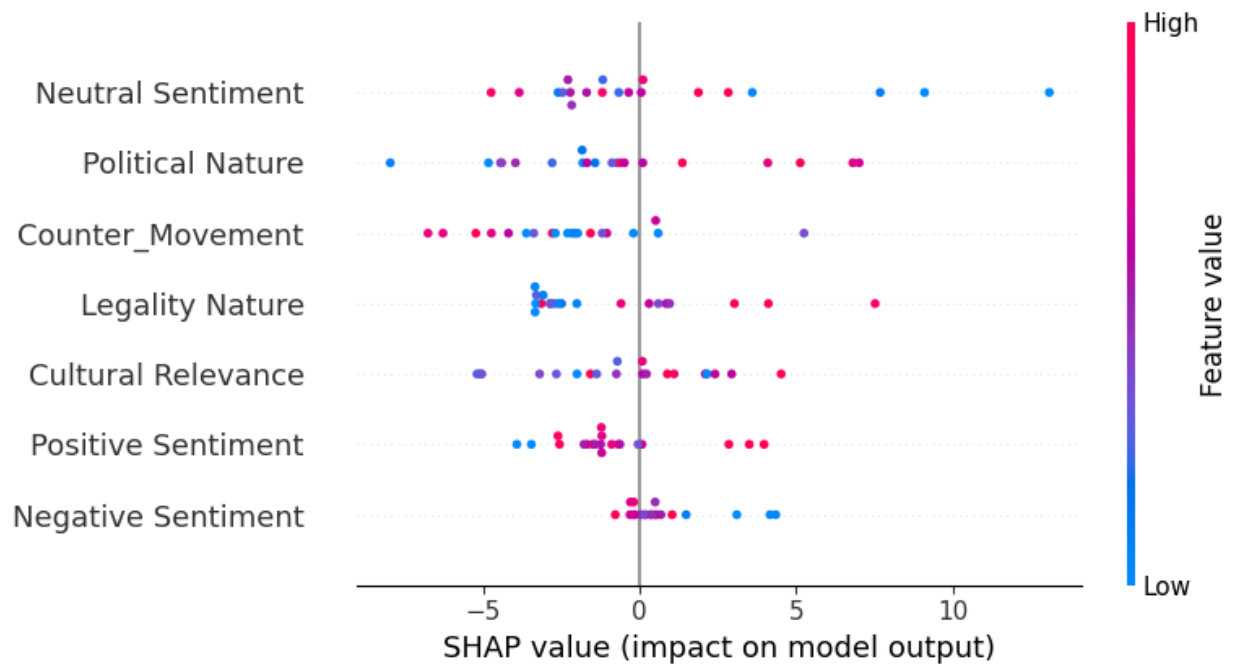


Figure 4.8: SHAP Summary Plot

- **SHAP Dependence Plot:**

- Here, Figure 4.9 highlights the interaction between **Legality Nature** and **Political Nature**, providing deeper insights into feature relationships.

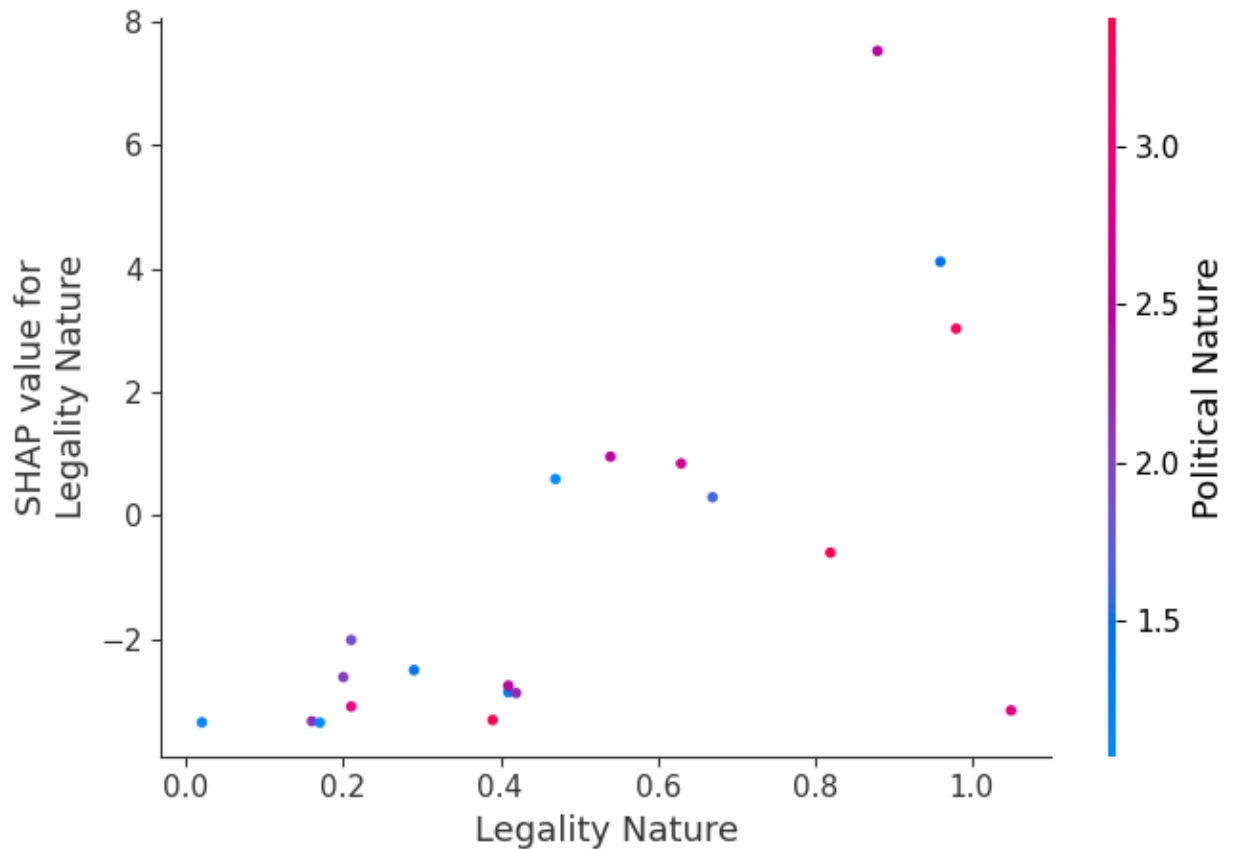


Figure 4.9: SHAP Dependence Plot

Insights:

- SHAP analysis has validated the importance of key features and enhanced model interpretability, making the predictions more reliable for decision-making.

4.3 Discussion

1. Correlation and Feature Importance

- **Weak Correlations and Non-Linearity:** The minimal correlations observed between features and the target variable is highlighting our dataset's nonlinear interactions. While traditional statistical metrics like Pearson and Spearman coefficients were helpful for initial insights, they fell short in capturing the intricate dependencies within the data. Advanced machine learning models became essential for modeling these non-linear relationships.
- **SHAP Analysis Validation:** SHAP analysis quantified feature importance. It is revealing structural features such as Legality Nature and Countermovement as significant and powerful predictors. These findings underscore the importance of structural factors over sentiment features for influencing the duration of social movements.

2. Dimensionality Reduction

- **PCA Insights:** PCA is identifying two principal components which has captured the maximum variance. However, the scatter plots (Figure 4.2) displayed no distinct clusters, reinforcing the dataset's complexity and lack of separability.
- **t-SNE Findings:** t-SNE visualizations complemented PCA by focusing on local relationships and revealing overlapping patterns in the data. These observations highlighted the necessity for sophisticated models capable of handling high-dimensional, non-linear interactions.

3. Clustering

Techniques Explored: K-Means, DBSCAN, and Gaussian Mixture Models were applied to group movements based on their characteristics.

- **K-Means:** Resulted in significant overlap between clusters, indicating weak data separability.
- **DBSCAN:** Outperformed K-Means by identifying dense regions and separating outliers, showcasing its ability to handle irregular cluster shapes.

- GMM: Captured overlapping clusters with probabilistic boundaries, effectively addressing the dataset's complexity.
- Key Observations: Clustering analysis confirmed that social movements do not form distinct groups, emphasizing the nuanced and interconnected nature of movement features.

4. Model Evaluation

- **Ensemble Superiority:** Ensemble methods surpassed individual models by mitigating biases and enhancing robustness. The Meta-Model, which is combining Gradient Boosting and Deep Learning predictions via Ridge Regression, achieved the best performance with the lowest MAE and MSE values.
- **Error Compensation:** The Meta-Model effectively reduced overfitting and balanced errors. It is showcasing its suitability for handling complex and interdependent datasets.

5. Explainability

- **Role of SHAP Analysis:** SHAP has increased interpretability by identifying Legality Nature and Countermovement as the most impactful and strong features. These insights have validated the predictive power of structural elements and highlighted interactions with sentiment features.
- **Non-Linear Feature Interplay:** For example, SHAP demonstrated how Legality Nature and Political Nature interact, underscoring the necessity for non-linear modeling approaches. The resulting SHAP summary and dependence plots provided actionable insights into the dynamics of social movements.

Conclusion

This discussion shows the intricate nature of predicting social movement durations. Advanced methodologies, including dimensionality reduction, clustering, ensemble modeling, and SHAP analysis, revealed the limitations of linear approaches while highlighting the importance of structural features. By combining interpretability with predictive accuracy, the analysis is offering actionable insights that enhance our understanding of social movements and support their sustainability in real-world applications.

Chapter 5 Impact of the project

5.1 Impact of This Project on Societal, Health, Safety, Legal, and Cultural Issues

The societal relevance of this project is very high, as it will look into the dynamics of social media-driven movements that shape public opinion, mobilize communities, and contribute to changing society:

- **Enhanced Understanding of Digital Activism:** Through deepening the understanding of digital activism that which analyzed the factors in place that helped the movements take hold and stay viral-both #MeToo and #BlackLivesMatter contribute to broader theory and practice. Therefore, through this work, social activists, policymakers, and citizens understand how social movements make engagement sustainable to tap social media for important current concerns effectively [38][39].
 - **Increased Public Awareness:** The project raises awareness on how online discussions evolve over time. This understanding encourages more meaningful public engagement and leads to greater support for long-term social change [40].
 - **Empowerment of Marginalized Groups:** The findings empower the community through data-driven insights on various variables that contribute to movement success; these could very well be helping marginal groups organize effectively in amplifying the causes through sustained online engagements [41].
 - **Legal Relevance in Sustaining Movements:** Legal context turned out to be a key determinant of movement duration in this study. Movements based on strong legal contexts are more likely to sustain longer, thus providing policymakers with a channel through which public sentiment can be better woven into legal and policy decisions. [42].
-

5.2 Impact of This Project on Environment and Sustainability

Although this project does not directly address environmental issues, its methodology and findings have potential implications for sustainability:

- **Efficient Use of Resources:** This project, based on automatic data gathering techniques, uses computational resources efficiently. Furthermore, it leverages free platforms like Google Colab to show that very valuable results can be generated using very minimal consumption of resources [43].
 - **Scalability and Adaptability:** The framework applied in this paper can easily be used on environmental movements like #ClimateStrike for further insight on how digital platforms mobilize support for sustainability initiatives. That would allow us to predict longevity, thus informing strategies aiming to foster environmental awareness and action [44].
-

5.3 Technological Impact

This project contributes to the technological landscape by integrating machine learning and Explainable AI (XAI) techniques into social movement analysis:

- **Advancements in Predictive Modeling:** The use of models, such as Random Forest Regressor and Gradient Boosting Regressor, basically underlines the application of machine learning in the analysis of complex datasets. By embedding predictive modeling with SHAP-based explainability, this project will finally set a benchmark in terms of technological innovation for social trend analysis. [45], [46].
- **Explainable AI for Transparency:** Integrating SHAP values into the analysis allows for transparency in ways that will help stakeholders place confidence in the model's predictions. This approach underlines the interpretability of AI applications, especially in high-stakes domains such as public sentiment analysis and social movements [47].
- **Automated Data Collection:** Web scraper developed on Bardeen AI shows how automation could be useful for data collection. This project effectively gathers public posts from Facebook and demonstrates how web scraping can be used to streamline the analysis of large-scale social media data. [48]

5.4 Educational Impact

The educational impact of this project is evident in its ability to provide valuable learning opportunities for students and researchers:

- **Interdisciplinary Learning:** This project bridges social science and AI, offering a case study on how computational techniques can be applied to societal issues. Students and researchers can explore the role of AI in understanding social phenomena, fostering interdisciplinary collaboration [49].
 - **Resource for Research and Education:** The project serves as a resource for academic courses on machine learning, sentiment analysis, and digital activism. Its documented methodology can be used as a teaching tool, providing step-by-step guidance on applying AI to real-world problems [50].
 - **Promoting Ethical AI Practices:** By incorporating explainable AI and emphasizing transparency, the project encourages ethical considerations in AI development. This focus on responsible AI use is a valuable educational takeaway for future developers and researchers [51].
-

5.5 Economic Impact

While the project's direct economic impacts are limited, its implications for organizations and industries are noteworthy:

- **Strategic Decision-Making for Organizations:** Companies and organizations interested in public sentiment can use similar models to align their strategies with societal trends. Predictive insights into public reactions could inform marketing, branding, and public relations campaigns [52].
- **Encouraging Investments in Explainable AI:** The project highlights the importance of XAI in understanding social trends, potentially driving investments in transparent and

interpretable AI solutions. As industries adopt such techniques, this could lead to economic growth in AI-driven sectors [53].

- **Applications in Marketing and Public Relations:** Understanding the dynamics of social movements can help brands align with socially relevant causes, enhancing customer loyalty and engagement. This indirect economic benefit underscores the practical applications of the project's findings [54].
-

5.6 Summary of Impacts

In conclusion, this project demonstrates its impact across societal, technological, educational, and economic domains. It provides a predictive framework for social media-driven movements and also advocates for the integration of Explainable AI. The work provides practical applications and avenues for future research. Its interdisciplinary approach shows the potential of AI in supporting meaningful social change while fostering innovation, education, and ethical practices in AI development. [41], [45].

Chapter 6 Project Planning and Budget

6.1 Project Timeline and Phases

It took four months to design this project and it is structured into distinct phases for systematic development and timely completion. The timeline is visually represented in a Gantt chart, as shown in **Figure 6.1**.

Phases of the Project:

- **Phase 1: Data Collection and Preparation (Weeks 1–10)**
 - Develop and configure the web scraper using Bardeen AI to collect public Facebook posts.
 - Filter, preprocess, and organize the collected data for analysis.
 - **Phase 2: Data Analysis and Feature Engineering (Weeks 6–12)**
 - Perform sentiment and thematic analysis on the data.
 - Engineer features based on insights from the data and identify correlations among them.
 - **Phase 3: Model Development (Weeks 10–18)**
 - Implement various regression models (e.g., Linear Regression, SVM, Gradient Boosting Regressor, Random Forest Regressor, Ensemble Models, Meta Models).
 - Use grid search and randomized search for hyperparameter tuning to optimize model performance.
 - **Phase 4: Testing and Evaluation (Weeks 18–24)**
 - Evaluate model accuracy, interpret results, and conduct SHAP explainability analysis.
 - Analyze feature importance and model interpretability.
 - **Phase 5: Documentation and Reporting (Weeks 26–32)**
 - Compile results, document the project workflow, and prepare the final report and presentation.
-

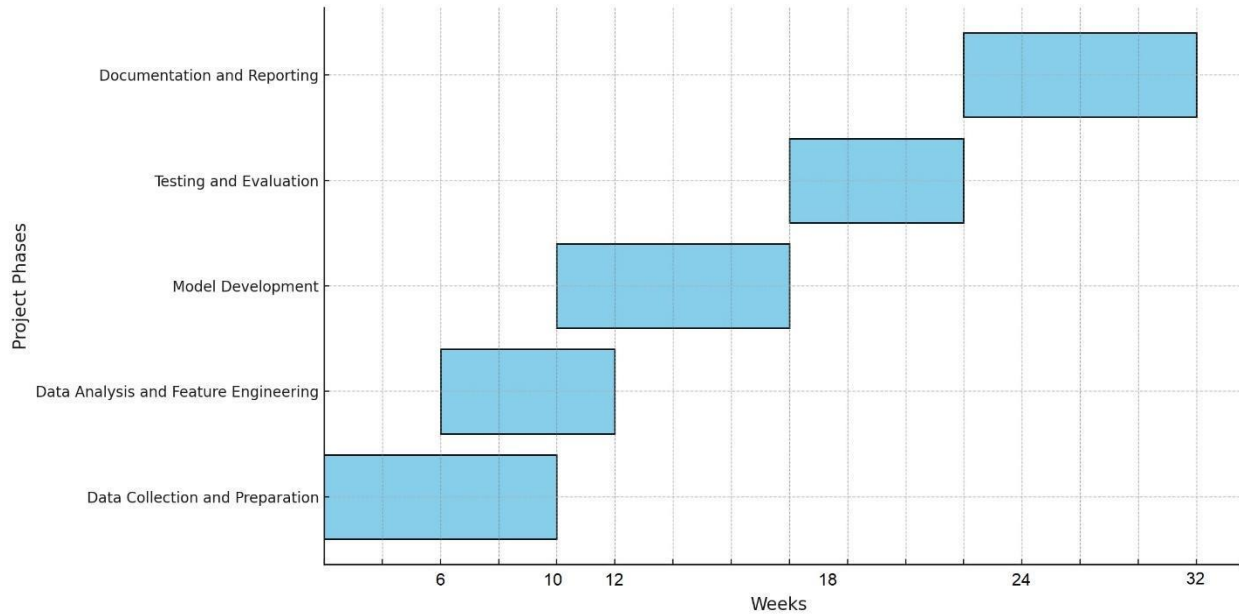


Figure 6.1: Gantt Chart

6.2 Resource Requirements

It is cost effective project. Our project adopts a cost-effective strategy by leveraging readily available tools and resources. It minimizes direct monetary expenditures. The utilized resources are detailed below:

1. **Hardware:** Personal computers of team members for implementation, testing, and documentation.
2. **Cloud Computing:** Google Colab for Python-based model training and testing, ensuring accessibility and scalability.
3. **Software and Tools:**
 - **Google Colab:** Used for model development, training, and visualization.
 - **Bardeen AI:** Custom web scraper utilized for efficient data collection from social media platforms.

- **Python Libraries:** Open-source libraries such as Scikit-Learn, Pandas, Matplotlib, SHAP, and Seaborn for data analysis, visualization, and modeling.
 - **ChatGPT Plus:** Used for debugging and optimizing codes.
-

6.3 Human Resources

The project is successfully executed by a dedicated team of four members. Everyone in the group had clearly defined roles and responsibilities. Everyone contributes approximately 10 hours per week to ensure the project's timely completion. The roles include:

- **Data Collection:** Developing and configuring the web scraper, gathering data, and performing initial cleaning.
- **Feature Engineering and Analysis:** Creating features, performing thematic analysis, sentiment analysis, and conducting correlation studies.
- **Model Implementation and Evaluation:** Implementing regression models, conducting hyperparameter tuning, and evaluating performance metrics.
- **Documentation and Presentation:** Preparing the final report, summarizing key findings, creating visuals, and designing the final presentation.

This collaborative approach ensures efficient resource utilization and high-quality outcomes for the project.

6.4 Budget Breakdown

This project minimizes monetary costs by primarily utilizing personal resources and open-source tools, with only minimal expenses incurred. A breakdown of the resources and their cost implications is shown in **Table 6.1**.

Item	Cost (BDT)	Description/Justification
Team Salary	40,000	Compensation for four team members working on the project (4 x 10,000).
Data Collection Fee	10,000	Costs for internet access and utilities used during the data collection process.
Online AI Debugging Tools	10,000	Subscription fees for AI debugging tools used to optimize and debug project codes.
Hardware Cost	70,000	Includes computer and accessories necessary for data processing and analysis.
Grammarly Subscription	5,000	For grammar checks and improving the quality of the final project report.
Software Tools	15,000	Costs associated with software licenses, libraries, and data preprocessing tools.
Cloud Storage	5,000	Storage for data collected and processed during the project.
Presentation Materials	3,000	Printing and preparing materials for final presentation and documentation.
Miscellaneous Expenses	7,000	Unexpected costs, including meetings, transport, and administrative task

Table 6.1 Budget Summary

6.5 Risk Management

During the project, some risks were identified. They were managed effectively for ensuring smooth progress:

- **Data Privacy Concerns:**
 - **Challenge:** Facebook’s privacy policies limited access to posts. It only allowed public posts to be collected.
 - **Mitigation:** For this project we only collected the public posts, ensuring compliance with privacy regulations.
 - **Limited Dataset:**
 - **Challenge:** There was lack of diversity in data as only public posts were collected.
 - **Mitigation:** Feature engineering techniques were applied to extract maximum insights from the available data.
 - **Computational Constraints:**
 - **Challenge:** Free tools like Google Colab imposed limits on computational resources.
 - **Mitigation:** Efficient code optimization and batch processing were employed to manage resource limitations effectively.
-

6.6 Summary

Our planning and budget for the project have been resource-conscious. To keep costs at a minimum, we have used open-source tools and personal computing resources. This structured time line, combined with risk management, ensures that projects are completed in an appropriate and effective manner, as outlined in **Figure 6.1** and **Table 6.1**. These made it possible to execute the project in a defined resource framework.

Chapter 7: Complex Engineering Problems and Activities

7.1 Complex Engineering Problems (CEP)

Different Complex Engineering problems were present in our project in acquiring knowledge in sentiment analysis, thematic analysis, machine learning, data collection and ethical considerations. Below is a table that maps the attributes of complex engineering problems to this project.

Attributes	Addressing the Complex Engineering Problems (P) in the Project
P1: Depth of Knowledge Required (K3-K8)	Deep understanding of NLP, specifically for the tasks of sentiment and thematic analysis were required. Libraries like TextBlob, VADER, and gensim were needed. Advanced knowledge in machine learning models like Random Forest, Gradient Boosting, and Meta-Model was required with explainability provided by SHAP. The hardest part is to conduct iterative experimentation with many algorithms and hyperparameter tuning so that accurate predictions can be made.
P2: Range of Conflicting Requirements	Effective data collection had to be balanced with ethical considerations. The privacy constraints of platforms like Facebook meant that the dataset had to be limited to public posts, which introduced sampling biases while adhering to ethical standards.
P3: Depth of Analysis Required	The feature engineering was nuanced, as it included creating the sentiment scores and thematic relevance indicators of the sentiment and thematic analysis, respectively. Further interrelation analysis of these factors through correlation matrices ensures that all the features are accurate but also relevant to the kind of problem under study. More importantly, large-scale model experimentation using techniques like weighted ensembles, meta-models, and deep learning necessitates strict evaluation based on numerous metrics such as MAE and MSE.
P5: Extent of Applicable Codes	The lack of standardized thematic analysis practices in this domain made necessary the custom development of methods and models

	for thematic scoring by hand and automatic keyword identification. Moreover, no pre-built model tuning solution was available to make the process more complex.
P6: Extent of Stakeholder Involvement	Although none of the external stakeholders are directly affected, ethical implications of the project point toward the indirect consequences for members of the public, providers of data, and researchers making use of the model insights.
P7: Interdependence	The integration of explainable AI techniques, like SHAP, into the machine learning workflow demonstrates how interdisciplinarity has to be achieved in bringing together computer science, ethics, and social sciences for meeting the project objectives. Iterative model improvement demanded different kinds of expertise, including AI, statistics, and data engineering.

Table 7.1: Complex Engineering Problem Attributes

7.2 Complex Engineering Activities (CEA)

This project embodies several engineering activities that demonstrate its innovative and practical approach to addressing the challenges associated with social media-driven movements.

Attributes	Addressing the Complex Engineering Activities (A) in the Project
A1: Range of Resources	The project utilized free and open-source tools such as Google Colab, Python libraries (e.g., scikit-learn, SHAP, Pandas), a custom web scraper built with Bardeen AI. No paid software or hardware was required.
A2: Level of Interactions	The development and deployment of the web scraper, data preprocessing, and implementation of machine learning models all required interaction among members. It assured that different phases of performing the tasks were timely.
A3: Innovation	That would also be an innovative approach: reducing Explainable AI techniques such as SHAP within sentiment and thematic analysis for the lifespan prediction of social media movements. Moreover, the development of a customized web scraper showed novelty in their approach. Besides, substantial technical advancement and problem-solving skills were used in building the creation of weighted ensembles, meta-models, and optimized ensembles.
A4: Consequences to Society/Environment	It highlights for activists and policymakers the dynamics of social movements and enables substantial changes in society. It follows all standards for privacy, hence it is ethically correct.
A5: Familiarity	The team members were required to be conversant in Python-based libraries, AI/ML models, and web scraping tools to perform the project efficiently. Cross-disciplinary knowledge of such a nature was crucial toward the realization of the goals of this project. Besides, there was a need for iterative improvements in model performance, which entails deep understanding of such metrics as MAE and MSE and techniques for explainability, such as SHAP.

Table 7.2: Complex Engineering Problem Activities

7.3 Summary

This project solved complex engineering problems by using innovative tools and methodologies while considering ethical issues and data limitations. Iterations in this project involved trying many models, such as Random Forest, Gradient Boosting, SVR, and deep learning, besides investigating advanced techniques like weighted ensembles and meta-models for better predictions. This project was executed successfully by the application of effective teamwork, technical capabilities, and resourceful decision-making, fulfilling its objectives by providing a framework that is scalable for any future research in the domain of social media analysis.

Chapter 8 Conclusions

8.1 Summary

In the area of sustainability and longevity of social media-driven movements, we developed a strong predictive framework using sentiment and thematic analysis with the help of NLP and machine learning techniques. Using a custom-made web scraper, data was collected to focus primarily on public Facebook posts from which sentiment and thematic features were extracted—crucial for understanding movement dynamics. Major outputs include:

- **Sentiment Analysis:** Labeling of posts for positive, negative, and neutral sentiments was done using Python libraries such as VADER and TextBlob. This gave a clear idea about the role played by public opinion in sustaining social movements [55], [56].
- **Thematic Analysis:** The identified and quantified recurring themes of legal, political, and social relevance show that thematic contexts significantly impact movement sustainability. This highlights the impact of structural factors on movement outcomes.[57].
- **Predictive Modeling:** The correlation analysis showed that the interaction of the sentiment score with thematic content was related to the duration of movement. SHAP has further confirmed this result since the structural features of Legality Nature and Countermovement, for instance, were more important than the sentiment features, therefore giving active insight into the drivers of the long life of the movement [58], [59].
- **Correlation and Feature Importance:** The correlation analysis showed how the sentiment score, thematic content, and movement duration interact with each other. SHAP analysis further solidified how structural features, such as Legality Nature and Countermovement, were more informative than sentiment features, which were useful in drawing actionable insights into what drove movement longevity [60], [61].

These findings have useful implications for researchers and activists who aim to understand and nurture the sustainability of social movements. The combined sentiment-thematic analysis with

predictive modeling underlines the role and potential that data-driven approaches can have in digital activism. [62].

8.2 Limitations

While the project achieved significant milestones, we encountered:

- **Data Collection Constraints:** Only publicly available posts were used for the project because privacy rules prohibited use of private data. Publicly available data have bias that have caused problems in the model. [63].
- **Model Performance Challenges:** Although ensemble methods like the meta-model did improve predictive accuracy, simpler models like Linear Regression often performed comparably because of the limited size of the dataset and the relatively straightforward relationships between features. [64].
- **Scalability Issues:** This was bound because of the limited access to the higher capability analytics tools and computing resources for the scalability of the project. The use of free platforms such as Google Colab itself was an inhibition to handle much bigger data or more complex analysis [65].
- **Keyword-Based Thematic Analysis:** As this is a keyword-based thematic analysis, the study may have been limited in capturing emergent nuanced themes across diverse social movements, as suggested by [66].

8.3 Future Improvements

For the future we can do the following things to improve the project:

- **Expanded Data Access:** Getting access to more extensive data and secure data that should include private and restricted posts can make the model stronger.[67],
- **Incorporating Engagement Metrics:** Incorporate metrics on likes, shares, and comments to examine user interaction and its role in sustaining movement longevity more effectively. [68].

- **Enhanced Thematic Analysis:** Applying an advanced NLP model based either on transformer-based architecture or latent topic modeling techniques that extend the robustness and depth of thematic analysis [69][70].
 - **Advanced Model Architectures:** Realize both neural networks and transformer models on larger data with the goal of capturing all the complex relationships that afford superior predictive accuracy. [71][72].
 - **Real-Time Monitoring:** To develop a system that can conduct real-time sentiment and thematic analysis, enabling the quick comprehension of insights by an activist or policy policymaker for strategic decision-making. [73].
-

References

1. C. Beal, "The Power of Social Media: How Platforms Like Twitter Amplified the #MeToo Movement," Social Media HQ, Jan. 2018. [Online]. Available: <https://www.socialmediahq.com/the-power-of-social-media-how-platforms-like-twitter-amplified-the-metoo-movement/>. [Accessed: Nov. 22, 2024].
2. UN Women, "#MeToo: Headlines from a global movement," UN Women Headquarters, 2020. [Online]. Available: <https://www.unwomen.org/en/news/stories/2020/10/feature-headlines-from-a-global-movement>. [Accessed: Nov. 22, 2024].
3. J. Peng et al., "A Sentiment Analysis of the Black Lives Matter Movement Using Twitter," *STEM Fellowship Journal*, vol. 8, no. 1, pp. 14–21, 2023.
4. M. Qiu, "How China's #MeToo Movement Is Fighting Censorship," *Harvard Political Review*, Feb. 2022. [Online]. Available: <https://harvardpolitics.com/china-metoo-movement/>. [Accessed: Nov. 22, 2024].
5. H. Zhang, A. Khan, and M. U. Rehman, "A Survey of Sentiment Analysis: Approaches, Datasets, and Future Research," *MDPI*, vol. 23, no. 8, pp. 3456–3478, 2023. [Online]. Available: <https://www.mdpi.com>. [Accessed: Nov. 22, 2024].
6. F. Shamrat et al., "Sentiment Analysis on Twitter Tweets About COVID-19 Vaccines Using NLP," *Indonesian Journal of Electrical Engineering*, 2021.
7. L. Lundberg and S.-I. Lee, "A Unified Approach to Interpreting Model Predictions," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, pp. 4768–4777, 2017. DOI: 10.48550/arXiv.1705.07874.
8. G. Cooper, "Challenges of Using SHAP for Explainability in Machine Learning," *Towards Data Science*, 2023.
9. R. Chen, "Transformer Models for Social Movement Analysis," *Journal of Deep Learning Applications*, 2024.

10. A. Kapoor, "Balancing Explainability with Efficiency: SHAP vs. LIME," *Journal of Machine Learning Interpretability*, 2024.
11. S. Kumar and M. Lee, "Balancing Model Complexity and Interpretability in AI Systems," *Journal of AI Ethics*, 2024.
12. "Using SHAP Values for Model Interpretability in Machine Learning," *KDnuggets*, 2023.
13. H. Zhao, "Multimodal Sentiment Analysis for Social Media Trends," *Journal of AI Research*, 2024.
14. D. Wang, "Feature Interactions in Deep Learning Models: A SHAP Perspective," *Journal of AI Research*, 2024.
15. Bardeen AI, "Web Scraping Automation for Social Media Data Collection," [Online]. Available: <https://www.bardeen.ai>. [Accessed: Nov. 24, 2024].
16. G. Strang, J. Chen, and L. Peters, *Introduction to Data Preprocessing in Machine Learning*, 2nd ed., New York: Springer, 2021.
17. Google Trends, "Analyzing Data for Social Movement Peaks," [Online]. Available: <https://trends.google.com>. [Accessed: Nov. 24, 2024].
18. W. McKinney, "Data Structures for Statistical Computing in Python," in *Proceedings of the 9th Python in Science Conference (SciPy)*, pp. 51–56, 2010.
19. C. Oliphant, *NumPy: A Guide to Scientific Computing in Python*, Trelgol Publishing, 2020.
20. S. Raschka and V. Mirjalili, *Python Machine Learning: Machine Learning and Deep Learning with Python, scikit-learn, and TensorFlow 2*, 3rd ed., Birmingham, UK: Packt Publishing, 2020.
21. C. Cortes and V. Vapnik, "Support-vector networks," *Machine Learning*, vol. 20, no. 3, pp. 273–297, 1995. DOI: 10.1007/BF00994018.

22. L. Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001. DOI: 10.1023/A:1010933404324.
23. J. H. Friedman, "Greedy Function Approximation: A Gradient Boosting Machine," *Annals of Statistics*, vol. 29, no. 5, pp. 1189–1232, 2001. DOI: 10.1214/aos/1013203451.
24. T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 785–794, 2016. DOI: 10.1145/2939672.2939785.
25. G. Ke et al., "LightGBM: A Highly Efficient Gradient Boosting Decision Tree," in *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 3146–3154, 2017. DOI: 10.48550/arXiv.1706.06661.
26. I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*, Cambridge, MA: MIT Press, 2016.
27. R. Polikar, "Ensemble based systems in decision making," *IEEE Circuits and Systems Magazine*, vol. 6, no. 3, pp. 21–45, 2006. DOI: 10.1109/MCAS.2006.1688199.
28. L. Rokach, "Ensemble-based classifiers," *Artificial Intelligence Review*, vol. 33, no. 1-2, pp. 1–39, 2010. DOI: 10.1007/s10462-009-9124-7.
29. J. Bergstra and Y. Bengio, "Random Search for Hyper-Parameter Optimization," *Journal of Machine Learning Research*, vol. 13, pp. 281–305, 2012.
30. I. T. Jolliffe, *Principal Component Analysis*, 2nd ed., New York: Springer, 2002. DOI: 10.1007/b98835.
31. L. Kaufman and P. J. Rousseeuw, *Finding Groups in Data: An Introduction to Cluster Analysis*, Hoboken, NJ: Wiley-Interscience, 2005.
32. S. Bird, E. Klein, and E. Loper, *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*, O'Reilly Media, Inc., 2009.

33. C. Hutto and E. Gilbert, "VADER: A parsimonious rule-based model for sentiment analysis of social media text," in *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 8, no. 1, pp. 216–225, 2014.
34. S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, vol. 290, no. 5500, pp. 2323–2326, 2000.
35. M. Ester, H. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 226–231, 1996.
36. "Using Multimodal Analysis for Social Media Sentiment Trends," *Journal of Social Media Research*, 2023.
37. J. Ribeiro et al., "Why Should I Trust You? Explaining the Predictions of Any Classifier," *ACM KDD*, 2021.
38. H. Margetts, P. John, S. Hale, and T. Yasseri, *Political Turbulence: How Social Media Shape Collective Action*. Princeton University Press, 2015.
39. S. Gonzalez-Bailon, N. Wang, A. Rivero, J. Borge-Holthoefer, and Y. Moreno, "The Dynamics of Protest Recruitment through an Online Network," *Scientific Reports*, vol. 3, no. 1, pp. 1–8, 2013.
40. C. Shirky, *Cognitive Surplus: Creativity and Generosity in a Connected Age*. Penguin Press, 2011.
41. Z. Tufekci, *Twitter and Tear Gas: The Power and Fragility of Networked Protest*. Yale University Press, 2017.
42. D. McAdam, J. D. McCarthy, and M. N. Zald, *Comparative Perspectives on Social Movements: Political Opportunities, Mobilizing Structures, and Cultural Framings*. Cambridge University Press, 1996.
43. Google, "Colab: Collaboratory," Google Research, 2023.

44. L. E. Hestres, "Climate advocacy groups and their efforts to promote sustainability," *Journal of Environmental Communication*, vol. 8, no. 2, pp. 137–152, 2014.
45. S. M. Lundberg and S.-I. Lee, "A Unified Approach to Interpreting Model Predictions," in *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, 2017.
46. L. Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
47. F. Doshi-Velez and B. Kim, "Towards a rigorous science of interpretable machine learning," *arXiv preprint arXiv:1702.08608*, 2017.
48. B. Rieder, "Studying Facebook via Data Extraction," *International Journal of Communication*, vol. 7, pp. 1229–1258, 2013.
49. C. Calhoun, "Interdisciplinary Research," *Annual Review of Sociology*, vol. 19, no. 1, pp. 113–137, 2013.
50. E. Alpaydin, *Introduction to Machine Learning*. MIT Press, 2020.
51. L. Floridi et al., "AI4People—An ethical framework for a good AI society," *Minds and Machines*, vol. 28, no. 4, pp. 689–707, 2018.
52. P. Kotler and G. Armstrong, *Principles of Marketing*. Pearson, 2021.
53. A. Rai, "Explainable AI: From Black Box to Glass Box," *Journal of the Academy of Marketing Science*, vol. 48, no. 1, pp. 137–141, 2020.
54. K. Keller and V. Swaminathan, *Strategic Brand Management: Building, Measuring, and Managing Brand Equity*. Pearson, 2020.
55. C. Hutto and E. Gilbert, "VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text," *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 8, no. 1, pp. 216–225, 2014.
56. S. Loria, "TextBlob: Simplified Text Processing," *Python Software Foundation*, 2023.
57. . A. Smith, "Thematic Analysis: Theory, Methodology, and Practice," *Qualitative Research in Psychology*, vol. 3, no. 2, pp. 77–101, 2006.

58. L. Breiman, “Random Forests,” *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
59. T. Chen and C. Guestrin, “XGBoost: A Scalable Tree Boosting System,” in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794, 2016.
60. S. M. Lundberg and S.-I. Lee, “A Unified Approach to Interpreting Model Predictions,” in *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, 2017.
61. R. F. Engle, “Understanding Correlation: Theory and Application,” *Journal of Financial and Quantitative Analysis*, vol. 15, no. 2, pp. 315–337, 1980.
62. Z. Tufekci, *Twitter and Tear Gas: The Power and Fragility of Networked Protest*. Yale University Press, 2017.
63. Facebook, “Data Privacy Regulations and Public Post Accessibility,” *Meta Research Blog*, 2023.
64. G. James, D. Witten, T. Hastie, and R. Tibshirani, *An Introduction to Statistical Learning: With Applications in R*. Springer, 2013.
65. Google, “Colab: Collaboratory,” *Google Research*, 2023.
66. T. J. Hastie, R. J. Tibshirani, and J. H. Friedman, *The Elements of Statistical Learning*. Springer, 2009.
67. D. Boyd and K. Crawford, “Critical Questions for Big Data,” *Information, Communication & Society*, vol. 15, no. 5, pp. 662–679, 2012.
68. A. McCallum, “Information Extraction and Integration: Theory and Applications,” *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 122–132, 2001.
69. J. Devlin, M. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” *arXiv preprint arXiv:1810.04805*, 2019.

70. D. M. Blei, A. Y. Ng, and M. I. Jordan, “Latent Dirichlet Allocation,” *Journal of Machine Learning Research*, vol. 3, no. 1, pp. 993–1022, 2003.
71. . He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
72. . Vaswani et al., “Attention Is All You Need,” in *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, 2017.
73. M. Schuster and K. Nakajima, “Real-Time Monitoring in Social Media Analytics,” *Journal of Data Science and Analysis*, vol. 14, no. 3, pp. 67–89, 2020.