

Multilingual Sentiment Analysis: A Machine Learning Approach with a Focus on Malayalam

Shubham Ojha
School of Computer Science and Engineering
Vellore Institute of Technology Chennai, India
shubham.ojha2020@vitstudent.ac.in

Rajalakshmi R
Student School of Computer Science and Engineering
Vellore Institute of Technology Chennai, India
rajalakshmi.r@vit.ac.in

Abstract

In this study, Multilingual Sentiment Analysis serves as a robust computational tool for discerning and categorizing a spectrum of emotions embedded within texts across various languages. Emphasizing the significance of the Malayalam language, this research sheds light on the intricate emotional nuances prevalent within this linguistic domain. Through the application of machine learning algorithms, the study delves into the identification of genuine versus fabricated sentiments expressed in social media posts and comments. Furthermore, it extends its analytical prowess to the classification of news articles, meticulously categorizing them into five distinct truthfulness categories: false, mainly false, true, half true, and mostly true. This experimental endeavour exemplifies the capacity of machine learning to transcend linguistic boundaries, offering deeper insights into the dynamics of international digital communications. By bridging the gap between language and sentiment analysis, this research contributes to a more nuanced understanding of the complexities inherent in diverse linguistic contexts, thereby facilitating more accurate and insightful interpretations of textual data.

Introduction

Multilingual sentiment analysis stands as a pivotal frontier in computational linguistics, facilitating the interpretation of emotions across diverse linguistic landscapes. It serves as a cornerstone for understanding public sentiment and trends in social media on a global scale. In this research endeavour, our focus is directed towards leveraging the power of multilingual sentiment analysis within the context of Malayalam, an important South Indian language. We present a comprehensive analysis of two distinct Malayalam datasets: one tasked with discerning between original and fabricated YouTube comments, and the other dedicated to classifying news content into five distinct veracity levels.

Our study embarks on a thorough exploration of machine learning and deep learning methodologies to unravel the intricate nuances embedded within Malayalam textual data. Traditional classifiers such as Support Vector Machines (SVM), Random Forest (RF), Logistic Regression, and Naive Bayes are rigorously evaluated alongside state-of-the-art techniques employing BERT classifiers and the Indic BERT tokenizer. By juxtaposing these diverse approaches, we aim to not only enhance the accuracy of sentiment analysis in Malayalam but also showcase the efficacy of integrating traditional and cutting-edge methodologies for multilingual sentiment analysis.

This interdisciplinary approach underscores the significance of bridging conventional techniques with advanced computational models, thereby paving the way for more robust and nuanced sentiment analysis across languages. Through this research, we aim to contribute to the ongoing discourse on multilingual sentiment analysis, shedding light on the challenges and opportunities inherent in analyzing emotions within diverse linguistic contexts.

Literature Review

Salvador Contreras Hernández, María Patricia Tzili Cruz, and José Martín Espínola explored sentiment analysis of COVID-19-related tweets in Mexico using BERT-based models. Their research, focusing on semi-supervised learning with Spanish language models, demonstrated superior precision over multilingual BERT and traditional classifiers. This underscores the effectiveness of language-specific models in capturing public sentiment, offering significant implications for public health decision-making during the pandemic. [1]

George Manias, Argyro Mavrogiorgou, and Athanasios Kiourtis investigated multilingual sentiment analysis on Twitter, emphasizing the importance of language- and domain-agnostic approaches. Their study assessed the efficacy of four BERT-based classifiers against a zero-shot classification method. The findings suggest that while BERT-based classifiers are highly effective, zero-shot classification stands out as an innovative and scalable strategy, even though it may not reach the fine-tuned accuracy of its counterparts. [2]

Amina Amara, Mohamed Ali Hadj Taieb, and Mohamed Ben Aouicha spearheaded a study on COVID-19 trend analysis through the lens of Facebook data across seven languages using Latent Dirichlet Allocation (LDA). The research stands out by leveraging an underexplored platform for multilingual topic modelling, using graph visualization to trace the progression of public interest in the pandemic. The outcomes present unique insights into global sentiment and conversational trends on Facebook, mapping the chronological growth of discourse surrounding COVID-19. [3]

Rami Mohawesh, Sumbal Maqsood, and Althebyan's research offers a novel semantic solution to multilingual fake news detection by utilizing capsule neural networks. Their framework incorporates word embeddings and Qutaibah n-gram features, significantly enhancing fake news detection across languages. The results show a marked improvement over existing methods, highlighting the capability of capsule neural networks to adeptly manage the intricacies of multilingual text analysis. [4]

Anjum and Rahul Katarya have developed HateDetector, an advanced technique tailored for detecting hate speech on social media in various languages. Incorporating an improved seagull optimization algorithm and a hybrid diagonal-gated recurrent neural network, their method shows notable enhancements in accuracy, precision, recall, and F-measure. These results position HateDetector as a promising tool for

effectively monitoring and curbing hate speech across multiple languages and social media platforms. [5]

Caio Mello, Gullal S. Cheema, and Gaurish Thakkar examine the construction of Olympic legacy narratives through multilingual sentiment analysis of news articles on the London 2012 and Rio 2016 Olympics. Their methodology intertwines four sentiment analysis (SA) algorithms with explainable AI techniques to scrutinize methodological constraints. While recognizing SA's value in content analysis, the research reveals complexities inherent in multilingual and specialized domains. A blend of leading classifiers paired with clear AI practices promises improvements, and an intriguing utopian versus dystopian narrative dichotomy in Olympic legacy portrayal is disclosed. [6]

Purbani Kar and Swapan Debbarma tackle the challenge of detecting hate speech and analyzing sentiments in multilingual code-mixed social media texts. They present an enhanced seagull optimization algorithm coupled with a novel hybrid diagonal gated recurrent neural network. The proposed methodology has shown considerable improvement in precision, recall, and F-measure over traditional approaches, asserting its effectiveness as a powerful tool for hate speech and sentiment analysis in diverse linguistic settings. [7]

Simran Sidhu, Surinder S. Khurana, Munish Kumar, and Parvinder Singh provide a thorough review of sentiment analysis techniques in Hindi, focusing on negation handling and the development of Hindi SentiWordNet. They explore a range of methodologies, including both semantic and machine learning approaches, and assess tools such as lexicons, stemmers, and morphological analyzers. The paper concludes with a call for more advanced sentiment analysis research for Hindi, highlighting its vast native-speaking community and expanding online footprint, while also pointing out potential areas for future investigation. [8]

Christian E. Lopez and Caleb Gallemore introduce a sizable multilingual Twitter dataset designed to support research on COVID-19 social discourse. With over 2.2 billion tweets enriched by sentiment analysis and named entity recognition, this resource permits an in-depth examination of discussions surrounding the pandemic. Their conclusion affirms the dataset's importance as a tool for tracking the progression of public sentiment and conversational patterns about COVID-19, enabling diverse analyses of social media data. [9]

Siroos Rahmani Zardak, Amir Hossein Rasekh, and Mohammad Sadegh Bashkari address the gap in sentiment analysis tools for Persian text by customizing the BERT algorithm for this context. Their work benchmarks the BERT algorithm's performance against former approaches across various datasets, concluding that BERT excels in analysing Persian sentiment, evidenced by superior accuracy and F1 scores. This breakthrough underscores BERT's adaptability and potency in processing Persian language datasets for sentiment analysis. [10]

Research Objective

In this study, we aim to implement and evaluate a range of traditional machine learning algorithms, including Support Vector Machines (SVM), Random Forest (RF), Logistic Regression, and Naive Bayes, for sentiment classification on Malayalam datasets. Specifically, our focus lies in discerning between original and fake sentiments expressed in YouTube comments and categorizing news articles into five distinct veracity levels: false, mainly false, true, half true, and mostly true. Through this approach, we seek to assess the performance and efficacy of these conventional algorithms in capturing the nuanced emotions prevalent within Malayalam text.

Additionally, we endeavor to adapt and fine-tune BERT-based deep learning models to advance sentiment analysis of Malayalam YouTube comments and news articles. This entails leveraging the powerful capabilities of BERT (Bidirectional Encoder Representations from Transformers) to extract intricate emotional nuances embedded within the Malayalam language. Moreover, we explore different methodologies for incorporating BERT into our analysis, including direct BERT classification, translation of Malayalam text via the Helsinki-NLP pipeline, and utilization of the Indic BERT tokenizer for language-specific tokenization.

Our study is designed to comprehensively compare the effectiveness of these approaches in sentiment analysis, evaluating metrics such as accuracy, precision, recall, and F1-score on validation datasets. By rigorously assessing the performance of each model, we aim to ascertain which method yields the most reliable sentiment analysis results for Malayalam text. Furthermore, we delve into the efficacy of machine learning versus deep learning approaches in detecting fake comments and classifying news articles into nuanced truthfulness categories, providing insights into the strengths and limitations of each approach.

Ultimately, our research contributes to the broader field of multilingual sentiment analysis by addressing the unique challenges associated with language-specific sentiment assessment. Through comparative analyses of different algorithmic approaches, we seek to advance our understanding of the complexities inherent in multilingual sentiment analysis, thereby enhancing the accuracy and effectiveness of sentiment classification in diverse linguistic contexts.

Proposed Work

Our project embarks on advancing multilingual sentiment analysis with a concentrated focus on the Malayalam language, engaging with a pair of distinct datasets sourced from the CodaLab Fake News Detection in Dravidian Languages competition (Dravidian-LangTech@EACL 2024). [11] The first dataset comprises YouTube comments, each meticulously categorized as original or fake, serving as a testament to the intricacies of digital discourse.[11] The second dataset encompasses a diverse collection of news items, each painstakingly classified into one of five truthfulness categories: false, mainly false, true, half true, and mostly true.[11] This

granular classification scheme presents a nuanced spectrum of information authenticity, critical for the discerning algorithms we deploy.

The analytical journey of the project begins with the application of four classical machine learning classifiers; Support Vector Machine (SVM), Random Forest (RF), Logistic Regression, and Naive Bayes then executed on both datasets. The objective is to compare and contrast their validation performance rigorously, thereby determining the most efficacious model. The model that prevails in validation is then subjected to the crucible of testing within the datasets to ascertain its generalizability and robustness.

Venturing into the realm of deep learning, the project explores three distinct methodologies. Initially, we implement a BERT classifier to delve into the sentiment analysis directly. Subsequently, we enhance our linguistic reach by employing the Helsinki-NLP pipeline to translate the Malayalam text into English, thereafter applying the BERT classifier to this translated corpus.[12] The final stride in our deep learning endeavour employs the Indic BERT tokenizer, specially designed to improve tokenization of the Malayalam script, thus tailoring the BERT classifier to the linguistic nuances of Dravidian syntax and semantics. Through these varied approaches, we seek to construct a comprehensive analysis mechanism that stands at the vanguard of sentiment analysis for Malayalam language data, aiming for the highest echelons of accuracy and interpretability.

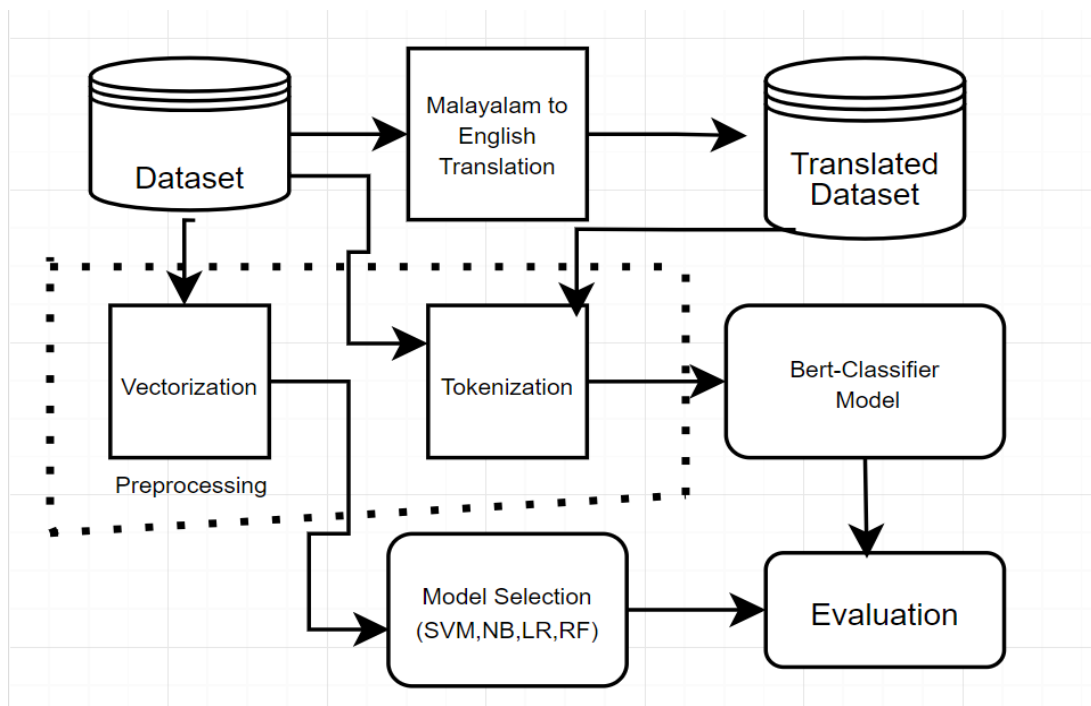


Figure 1 – Flowchart Diagram of Proposed System

Methodology

Traditional ML Models Steps

1. Data Preparation: Import the necessary libraries, including scikit-learn and pandas. Load the training dataset for model training.
2. Vectorization of Text: To convert the textual data into numerical vectors, use the Term Frequency-Inverse Document Frequency (TF-IDF) vectorizer.
3. The Logistic Regression Model: To find the ideal regularisation parameter (C), use crossvalidation and Logistic Regression using a hyperparameter grid search. Utilize the determined ideal hyperparameters to train the classifier.
4. The Random Forest Model: To determine the optimal set of hyperparameters, such as the number of estimators, maximum depth, minimum samples split, and minimum samples leaf, use Random Forest in conjunction with a grid search. Utilize the optimized hyperparameters to train the Random Forest model.
5. The Support Vector Model (SVM): Use the Support Vector Machine (SVM) technique for categorization. Utilizing the TF-IDF vectorized training data, train the SVM model.
6. The Naive Bayes Model: Multinomial Naive Bayes should be used for categorization. Utilizing the TF-IDF vectorized training data, train the Naive Bayes model.
7. Verification Assessment of the dataset: To evaluate the model, load the validation dataset. Take care of missing values and vectorize the text column. Evaluate each classifier's accuracy score on the validation set.
8. Ideal Model Choice: Determine which classifier on the validation set has the best accuracy. Declare this classifier to be the best model to test further. Evaluation of Testing Dataset.
9. Testing: Load the testing dataset. Utilizing the TF-IDF vectorizer, vectorize the text column of testing data. To predict labels for the test dataset, use the best classifier available.
10. Conclusion: Conclude the study by summarizing the chosen classifier's performance on the testing dataset.

Deep Learning Steps

1. Data Preparation: Begin by importing essential libraries such as transformers for the BERT model, torch for the deep learning framework, and pandas for data manipulation. Load the training dataset into a DataFrame for further processing.
2. Tokenization: Utilize the BERT tokenizer to convert the text into tokens that are understandable by the model. This step will include converting the Malayalam sentences into a format with token IDs, attention masks, and segment IDs suitable for BERT.

3. **Model Configuration:** Choose a pre-trained BERT model that is optimized for the Malayalam language or the multilingual version that includes Malayalam. Configure the BERT model with appropriate parameters, paying attention to the number of epochs, learning rate, and batch size for training.
4. **Fine-Tuning:** Using the tokenized text data, fine-tune the BERT model on the Malayalam sentiment analysis task. This involves training the model on the dataset, adjusting weights, and ensuring the model learns the context of the dataset effectively.
5. **Validation Assessment:** After the model has been fine-tuned, evaluate its performance on a separate validation set to gauge its effectiveness. Process the validation dataset similarly to the training set, with tokenization followed by the creation of DataLoader objects for the BERT model.
6. **Testing:** Load the testing dataset and process it through the BERT models using the same tokenization and DataLoader creation steps as the validation set. Apply the selected model to obtain predictions for sentiment classification.
7. **Conclusion:** Conclude the process by summarizing the performance of the BERT model on the testing dataset. Discuss the effectiveness of the model, its ability to generalize to unseen data, and any observations regarding its performance on multilingual sentiment analysis.

Result

The evaluation matrix used here are accuracy, training time, bar graph and Confusion matrices.

Youtube Comment Originality Detection

YouTube Comment Originality Detection

Model	Accuracy (%)	Training Time
Naive Bayes	78.8	0.51 sec
SVM	50.1	0.27 sec
Random Forest	75.5	12.61 sec
Logistic Regression	78.6	1.57 sec
BERT	75.2	1 hr 30 min 22 sec
EnglishTranslation+BERT	74.0	3 hr 56 min 27 sec
IndicBERT	75.1	1 hr 29 min 27 sec

Figure 2 - Table for YouTube Comment Originality Detection Accuracy

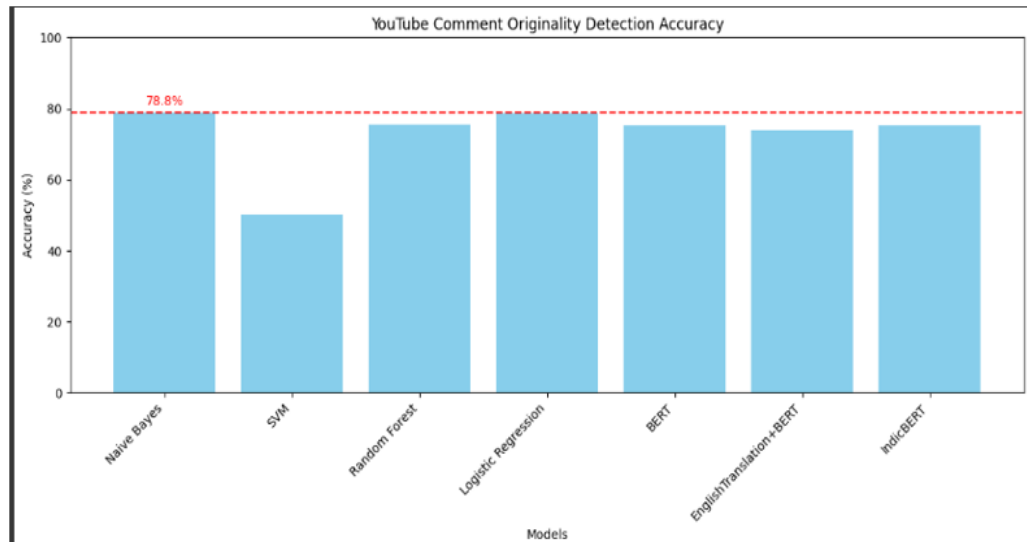


Figure 3 - Bar Graph showing accuracy of all models for YouTube Comment Originality Detection Accuracy

Confusion Matrix Table:

	True Positives (TP)	True Negatives (TN) \
Naive Bayes	333	344
SVM	338	336
Random Forest	336	368
Logistic Regression	376	398
BERT	321	305
EnglishTranslation+BERT	369	353
IndicBERT	375	368

	False Positives (FP)	False Negatives (FN) \
Naive Bayes	113	230
SVM	248	98
Random Forest	82	234
Logistic Regression	34	212
BERT	393	1
EnglishTranslation+BERT	293	5
IndicBERT	242	35

	Precision	Recall	Accuracy
Naive Bayes	0.746637	0.591474	78.8
SVM	0.576792	0.775229	50.1
Random Forest	0.803828	0.589474	75.5
Logistic Regression	0.917073	0.639456	78.6
BERT	0.449580	0.996894	75.2
EnglishTranslation+BERT	0.557402	0.986631	74.0
IndicBERT	0.607780	0.914634	75.1

Figure 4 – Table showing Confusion Matrix of all the models for Yobute Comment Originality Detection

Fake News Detection

Fake News Detection

Model	Accuracy (%)	Training Time
Naive Bayes	63.5	0.62 sec
SVM	65.8	0.37 sec
Random Forest	68.6	20.76 sec
Logistic Regression	61.7	12.07 sec
BERT	60.52	1 hr 56 min 12 sec
EnglishTranslation+BERT	65.7	4 hr 22 min 40 sec
IndicBERT	62.31	1 hr 50 min 18 sec

Figure 5 - Tabe for Fake News Detection Accuracy

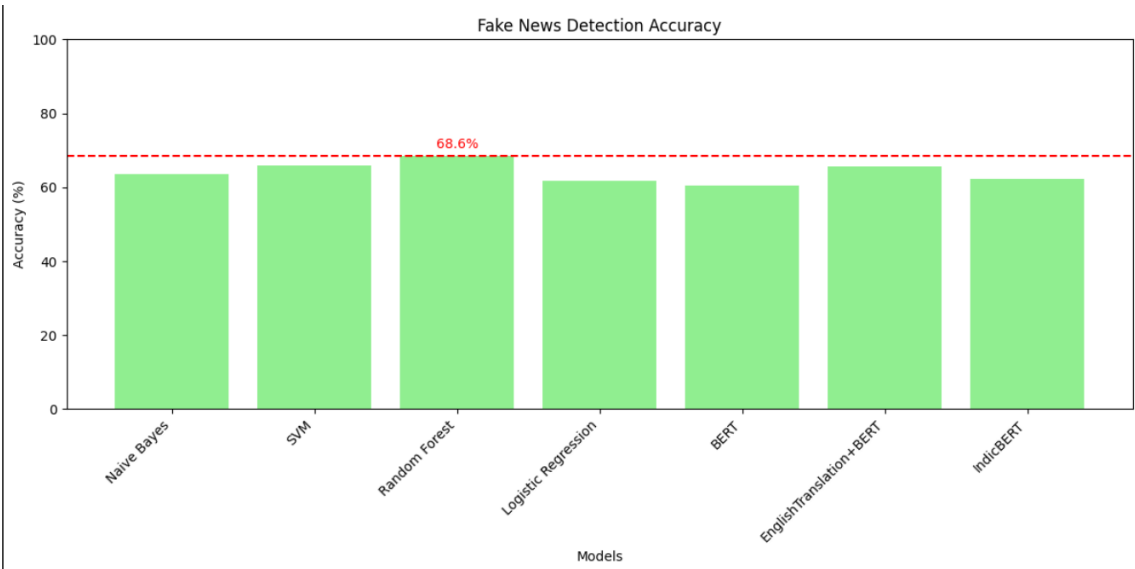


Figure 6 - Bar Graph showing accuracy of all models for Fake News Detection Accuracy

Confusion Matrix Table:

	True Positives (TP)	True Negatives (TN)	\
Naive Bayes	93	98	
SVM	94	93	
Random Forest	103	100	
Logistic Regression	102	89	
BERT	94	91	
EnglishTranslation+BERT	91	96	
IndicBERT	101	85	

	False Positives (FP)	False Negatives (FN)	\
Naive Bayes	11	48	
SVM	62	1	
Random Forest	11	36	
Logistic Regression	59	0	
BERT	52	13	
EnglishTranslation+BERT	45	18	
IndicBERT	28	36	

	Precision	Recall	Accuracy
Naive Bayes	0.894231	0.659574	63.50
SVM	0.602564	0.989474	65.80
Random Forest	0.903509	0.741007	68.60
Logistic Regression	0.633540	1.000000	61.70
BERT	0.643836	0.878505	60.52
EnglishTranslation+BERT	0.669118	0.834862	65.70
IndicBERT	0.782946	0.737226	62.31

Figure 7 – Table showing Confusion Matrix of all the models for Fake News Detection

Conclusion

The comparative analysis of sentiment analysis models on Malayalam datasets underscores the dominance of Naive Bayes for YouTube comment classification and Random Forest for fake news detection in terms of accuracy. Despite longer training times, BERT-based models offer deep linguistic insights, albeit not surpassing traditional models in accuracy.

A critical factor influencing performance disparities is dataset size. Deep learning models like BERT thrive on vast data, whereas the modest datasets of 3,200 comments for YouTube and 1,700 articles for news may limit their potential. Traditional models exhibit notable efficiency, possibly due to their compatibility with smaller datasets and less complex feature spaces.

The restricted dataset size emerges as a pivotal variable, potentially constraining deep learning algorithms' effectiveness. Expanding datasets or employing data augmentation strategies could enhance deep learning models' performance, leveraging their ability to grasp complex language patterns.

This research illuminates current sentiment analysis model capabilities for Malayalam texts while emphasizing the necessity for larger datasets to maximize deep learning

techniques' potential. Future endeavours should focus on dataset expansion or data augmentation strategies to enhance deep learning models' performance, striking a balance between model choice, dataset size, and computational efficiency in multilingual sentiment analysis research.

References

- [1] Contreras Hernández, S., Tzili Cruz, M. P., Espínola Sánchez, J. M., & Pérez Tzili, A. (2023). Deep learning model for covid-19 sentiment analysis on twitter. *New Generation Computing*, 41(2), 189-212.
- [2] Manias, G., Mavrogiorgou, A., Kiourtis, A., Symvoulidis, C., & Kyriazis, D. (2023). Multilingual text categorization and sentiment analysis: a comparative analysis of the utilization of multilingual approaches for classifying twitter data. *Neural Computing and Applications*, 35(29), 21415-21431.
- [3] Amara, A., Hadj Taieb, M. A., & Ben Aouicha, M. (2021). Multilingual topic modeling for tracking COVID-19 trends based on Facebook data analysis. *Applied Intelligence*, 51, 3052-3073.
- [4] Mohawesh, R., Maqsood, S., & Althebyan, Q. (2023). Multilingual deep learning framework for fake news detection using capsule neural network. *Journal of Intelligent Information Systems*, 60(3), 655-671.
- [5] Anjum, & Katarya, R. (2023). HateDetector: Multilingual technique for the analysis and detection of online hate speech in social networks. *Multimedia Tools and Applications*, 1-28.
- [6] Mello, C., Cheema, G. S., & Thakkar, G. (2023). Combining sentiment analysis classifiers to explore multilingual news articles covering London 2012 and Rio 2016 Olympics. *International Journal of Digital Humanities*, 5(2), 131-157.
- [7] Kar, P., & Debbarma, S. (2023). Multilingual hate speech detection sentimental analysis on social media platforms using optimal feature extraction and hybrid diagonal gated recurrent neural network. *The Journal of Supercomputing*, 79(17), 19515-19546.
- [8] Sidhu, S., Khurana, S. S., Kumar, M., Singh, P., & Bamber, S. S. (2023). Sentiment analysis of Hindi language text: a critical review. *Multimedia Tools and Applications*, 1-30.
- [9] Lopez, C. E., & Gallemore, C. (2021). An augmented multilingual Twitter dataset for studying the COVID-19 infodemic. *Social Network Analysis and Mining*, 11(1), 102.
- [10] Zardak, S. R., Rasekh, A. H., & Bashkari, M. S. (2023). Persian Text Sentiment Analysis Based on BERT and Neural Networks. *Iranian Journal of Science and Technology, Transactions of Electrical Engineering*, 47(4), 1623-1634.
- [11] Subramanian, M., Chakravarthi, B. R., Shanmugavadivel, K., Pandiyan, S., Kumaresan, P. K., Palani, B., Singh, M., Raja, S., Vanaja, & S, Mithunajha. (2023). Overview of the Shared Task on Fake News Detection from Social Media Text. In *Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages*. Varna, Bulgaria: Recent Advances in Natural Language Processing.

[12] Kakwani, D., Kunchukuttan, A., Golla, S., Gokul, N. C., Bhattacharyya, A., Khapra, M. M., & Kumar, P. (2020, November). IndicNLP Suite: Monolingual corpora, evaluation benchmarks and pre-trained multilingual language models for Indian languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020* (pp. 4948-4961).