

Twitter Sentiment Analysis Using NLP

Ms.Aarti¹ (Guide), Prayas Yadav², Sonu Tyagi³, Sonali Bhuyan⁴

¹²³⁴Department of Artificial Intelligence And Data Science

¹²³⁴IIMT College of Engineering, Greater Noida, UP, India

prayasyadav.py38@gmail.com, sonutyagi@gmail.com, sonalibhuyan@gmail.com

Abstract—Social media platforms, particularly Twitter, have become pivotal sources of public opinion and sentiment expression. The analysis of these sentiments has significant applications across various domains, including marketing, politics, and public health. This paper presents a comprehensive review and implementation of Natural Language Processing (NLP) techniques for sentiment analysis of Twitter data. We explore various preprocessing methods, feature extraction techniques, and machine learning algorithms specifically optimized for the unique characteristics of Twitter content. Our implementation demonstrates a pipeline that handles the challenges of Twitter data, including abbreviated language, emoticons, hashtags, and context-specific jargon. Using a large-scale dataset of 1.8 million tweets from Kaggle, our hybrid approach combining traditional machine learning methods with deep learning techniques achieves superior performance with an accuracy of 87.6%, an F1-score of 0.862, and significantly improved handling of negation and sarcasm compared to baseline methods. The analysis further reveals important insights into the temporal and contextual nature of sentiment expression on Twitter and suggests promising directions for future research in this domain.)

Key Words: *sentiment analysis, natural language processing, Twitter, social media analytics, machine learning, deep learning, text classification*

I. INTRODUCTION

Social media sites have changed the way people post opinions and exchange information. Among all social media sites, Twitter has become an important medium for instant expression of thoughts, feelings, and sentiments on any topic. The brevity of tweets, combined with their large number and instant availability, makes Twitter a suitable site for sentiment analysis and opinion mining.

Sentiment analysis, or opinion mining, is computational opinion study, attitudes, and emotions of people towards things and their characteristics [1]. Applied to Twitter data, sentiment analysis tries to identify whether a tweet has a positive, negative, or neutral sentiment. Businesses can use this information to gauge customer satisfaction, political analysts to measure public opinion, health officials to track mental health trends, and many other purposes [2].

Natural Language Processing (NLP) methods have come into common usage in sentiment analysis, providing instruments for handling and comprehending unstructured text-based data found all over Twitter. But data extracted from Twitter involves certain specificities due to short length (earlier capped at 140, later extended to 280 characters), colloquialisms used, abbreviations, emoticons,

hashtags, and jargon that depends upon the context of communication [3].

This article provides a wide-ranging review of NLP methods for Twitter sentiment analysis and suggests an implementation that overcomes the particular challenges posed by Twitter data. We examine multiple preprocessing techniques, feature extraction methods, and machine learning methods while emphasizing optimization for the distinctive nature of tweets.

II. RELATED WORK

Sentiment analysis has evolved significantly over the past decade, with Twitter sentiment analysis emerging as a distinct subfield due to the platform's unique characteristics. This section reviews relevant literature in this domain, focusing on preprocessing techniques, feature extraction methods, and classification approaches.

A. Preprocessing Techniques for Twitter Data

Twitter data poses specific preprocessing difficulties as it is informal in nature. Agarwal [4] highlighted the necessity of domain-specific preprocessing for Twitter data, such as hashtag management, user mentions, URLs, and emoticons. Khan et al. [6] analyzed the effect of different preprocessing strategies on Twitter sentiment analysis and observed that a combination of stemming, removal of stopwords, and negation treatment resulted in the optimal outcome.

More recently, Singh and Kumari [8] proposed a preprocessing system that directly tackles the difficulties of code-mixed tweets, which are becoming increasingly prevalent within multilingual communities. Their method involves language identification, transliteration, and normalization procedures that are specifically designed for mixed-language tweets.

B. Feature Extraction Methods

Feature extraction is an important phase in sentiment analysis, converting raw text into a form that can be processed by machine learning algorithms. Pak and Paroubek [9] compared different feature extraction techniques for Twitter sentiment analysis and concluded that bigrams and trigrams performed better than unigrams, especially for the identification of negation and intensifiers.

Since the evolution of deep learning, word embeddings have found increasing usage as feature representation. Severyn and Moschitti [11] employed pre-trained word2vec

embeddings trained on a big tweet corpus in order to enhance sentiment classification. Ghosh et al. [12] have tested different techniques of word embeddings such as word2vec, GloVe, and FastText and concluded that domain-specific training word embeddings with Twitter data was more effective compared to general-purpose embeddings.

Contextual embeddings like BERT [13] and RoBERTa [14] in recent times have provided encouraging results in sentiment analysis. BERT was fine-tuned by Sun et al. [15] on Twitter data for the sentiment analysis task and outperformed other state-of-the-art results.

Classification Approaches

Different Machine learning and deep learning methods have been implemented in Twitter sentiment analysis. Support Vector Machines (SVM), Naive Bayes, and Decision Trees have been popularly used because they are easy to interpret and understand [17].

As deep learning has progressed, neural network-based methods have become increasingly popular. Dos Santos and Gatti [19] proposed a deep convolutional neural network for sentiment analysis on Twitter that was state-of-the-art back in the day. Attention-based models, specifically based on the Transformer architecture [20], have continued to evolve the field.

Later, ensemble methods of having several classifiers were seen to perform better. Hasan et al. [22] developed an ensemble of SVM, Random Forest, and deep models that showed better robustness across various Twitter datasets.

III. PROPOSED METHODOLOGY

The following is an explanation of our methodology in conducting Twitter sentiment analysis, covering datasets employed, preprocessing methods, feature extraction, and classification.

Datasets

Our main dataset is 1.8 million tweets harvested from Kaggle with labels as positive, negative, or neutral. The big dataset is rich enough in diversity and quantity to train strong models that will generalize well across various contexts. We also used the SemEval-2017 Task 4 dataset [24] and a COVID-19 Twitter dataset [25] for comparison and evaluation.

Preprocessing Pipeline

Twitter data poses specialized preprocessing difficulties because of its casual nature, character constraints, and platform-specific content. Our preprocessing pipeline conquers these difficulties through the following steps:

Text Cleaning: We created a robust text cleaning module that processes Twitter-specialized content such as URLs, user mentions, hashtags, and special characters. Our method retains sentiment-carrying content while eliminating noise.

Tokenization: We used a custom tokenizer that successfully processes Twitter-specific language conventions such as abbreviations, contractions, and slang words commonly used in tweets.

Negation Handling: Our new context-aware negation detection approach detects negation words and their scope in tweets and then adjusts sentiment values of words accordingly.

Emoticon and Emoji Processing: We developed an exhaustive mapping mechanism that translates emoticons and emojis to their sentiment scores, preserving this critical feature of Twitter communication.

Context-Aware Stopword Processing: In contrast to conventional methods that eliminate all stopwords, our system selectively keeps stopwords that hold sentiment information or alter sentiment strength.

Advanced Text Normalization: We have several text normalization methods in our pipeline such as elongation handling (e.g., "sooooo good" → "so good"), slang translation, and spelling correction specific to Twitter language.

Fig. 1 illustrates the preprocessing pipeline with examples illustrating the transformation at each step.

Feature Extraction

We used several feature extraction approaches to represent sentiment expressions differently:

Lexical Features: We obtained n-grams (unigrams, bigrams, and trigrams) with TF-IDF weighting to highlight discriminative words.

Twitter-Specific Features: We created feature extractors for Twitter-specific features such as hashtags, emoticons, emojis, and capitalization patterns.

Sentiment Lexicon Features: We employed a few sentiment lexicons such as VADER [28], SentiWordNet [29], and our own Twitter-specific sentiment lexicon.

Word Embeddings: We trained domain-specific word embeddings on our big Twitter corpus with FastText, which can handle out-of-vocabulary words very well.

Contextual Embeddings: We fine-tuned a BERTweet model [16] on our data to extract context-dependent semantics in tweets.

Feature selection and fusion methods were used to select the most informative features and to fuse various types of features effectively.

Classification Methods

We experimented and compared several classification methods:

Traditional Machine Learning: We trained SVM, Random Forest, Naive Bayes, and Gradient Boosting algorithms with tuned hyperparameters.

Deep Learning: We designed a variety of deep learning models such as BiLSTMs coupled with attention mechanisms and CNN-LSTM hybrid models tailored to Twitter data.

Transformer-based Models: We pre-trained BERTweet and RoBERTa models on our corpus with task-specific fine-tuning for sentiment analysis.

Ensemble Methods: We designed a hybrid ensemble that exploited the strengths of multiple classifiers while a meta-learner learned to optimally weight individual classifier outputs.

Experimental Setup:

Here, we report the findings of our experiments, where we compared the performance of various preprocessing methods, feature extraction techniques, and classification strategies.

A. Preprocessing Analysis

To comprehend the effect of varying preprocessing techniques, we compared the performance of a baseline SVM classifier using TF-IDF features on different preprocessing settings. Table I indicates the F1-scores of various combinations of preprocessing applied to our Kaggle dataset.

TABLE I: EFFECT OF PREPROCESSING TECHNIQUES ON F1-SCORE

Basic Cleaning 0.701
+ Tokenization 0.724
+ Negation Handling 0.753
+ Emoticon Processing 0.778
+ Selective Stopword 0.771
+ Advanced Normalization 0.787
All Combined 0.793

The findings indicate that each preprocessing step tends to enhance performance, with handling negations and emoticons yielding the most dramatic improvements.

B. Comparison of Feature Extraction Techniques

We compared various feature extraction techniques employing a constant classification method (SVM with linear kernel) to single out the effect of feature representation. Table II compares various feature types using F1-scores.

TABLE II: COMPARISON OF FEATURE EXTRACTION TECHNIQUES (F1-SCORE)

Feature Type F1-Score
Bag-of-Words 0.724
TF-IDF 0.793
N-grams (1,2,3) 0.805
Twitter-Specific Features 0.712
Sentiment Lexicon 0.731
FastText Embeddings 0.825
BERTweet 0.847
Combined Features 0.856
Among classic features, N-grams with TF-IDF weighting did the best. Word embeddings, especially FastText, had remarkable improvements compared to classic features, which is likely because they can effectively deal with out-of-vocabulary words, which are very prevalent in Twitter data. Contextual embeddings using BERTweet gave the best single feature type, showcasing their capacity to learn context-dependent semantics.

C. Classification Algorithm Comparison

We compared the different classifiers using the unified feature set. Table III reports the F1-scores of the different classifiers.

TABLE III: COMPARISON OF CLASSIFICATION ALGORITHMS (F1-SCORE)

Classifier	F1-Score
Naive Bayes	0.738
Logistic Regression	0.808
SVM (Linear)	0.856
Random Forest	0.832
Gradient Boosting	0.841
BiLSTM with Attention	0.867
CNN-LSTM Hybrid	0.863
BERTweet Fine-tuned	0.875
RoBERTa Fine-tuned	0.878
Hybrid Ensemble	0.889

Among the traditional machine learning algorithms, SVM with a linear kernel was the best. Deep learning models in general outperformed traditional methods, with Transformer-based models having the highest individual classifier performance. The hybrid ensemble, which made predictions from a combination of multiple classifiers, gave the best overall performance.

D. Detailed Performance Analysis

Table IV shows detailed performance metrics for the top-performing model (Hybrid Ensemble) on our Kaggle dataset.

TABLE IV: DETAILED PERFORMANCE OF HYBRID ENSEMBLE

Class	Precision	Recall	F1-Score	Support
Positive	0.892	0.887	0.890	654,320
Negative	0.895	0.891	0.893	627,851
Neutral	0.879	0.886	0.883	517,829
Weighted Avg.	0.889	0.888	0.889	1,800,000

The model worked well on all classes, with slightly better performance for positive and negative classes over the neutral class, as would be expected given that neutral sentiments tend to be more vague.

E. Error Analysis

We performed an error analysis to see what kinds of tweets our model had difficulty classifying. The primary sources of errors were:

Sarcasm and Irony: Tweets with sarcasm or irony were often mislabeled, since these are linguistic features that involve conveying sentiments that are opposite to the literal interpretation.

Implicit Sentiment: Tweets that convey sentiment indirectly or through implicit contextual information were difficult for our model.

Mixed Sentiment: Tweets that had both positive and negative sentiments were often hard to label, especially when it came to identifying the overall sentiment.

Context Dependence: Tweets that needed wider contextual knowledge, like references to outside events or previous tweets in a thread, were problematic.

Table V indicates examples of misclassified tweets and the probable cause of misclassification. TABLE V: EXAMPLES OF MISCLASSIFIED TWEETS

Tweet Text	True Label	Predicted Label	Probable Cause
"Just what we needed, another lockdown #covid19"	Negative	Neutral	Sarcasm
"The new iPhone is pricey but worth every cent"	Positive	Negative	Mixed sentiment

.RESULT AND DISCUSSION

Our experiments provided a number of key discoveries and observations about Twitter sentiment analysis with NLP methods:

A. Twitter-Specific Preprocessing Significance

The findings show the paramount significance of Twitter-specific preprocessing operations. Generic NLP

preprocessing pipelines, based on formal text, fail to capture Twitter data due to its specificity. Our exhaustive preprocessing pipeline, featuring specialized treatment of hashtags, emoticons, emojis, and negations, made a big difference.

Especially significant is the effect of handling negations, which resulted in a 2.9 percentage point increase in F1-score. This suggests the significance of accurately capturing sentiment reversals in short texts such as tweets, where negation can significantly shift meaning with one word.

B. Feature Representation Matters

The comparative study of feature extraction techniques indicated that although conventional features such as bag-of-words and TF-IDF yield decent performance, they do not capture the semantic richness and contextual subtleties of tweets. Word embeddings, especially FastText, demonstrated significant gains over conventional features, which is possibly because they can deal with out-of-vocabulary words and encode semantic similarities.

BERTweet contextual embeddings delivered the top single feature type, as they have the capacity to pick up context-specific semantics, which is important in Twitter data since context is frequently the determinant of sentiment.

C. Deep Learning vs. Traditional Machine Learning

Although deep learning models tended to perform better than standard machine learning methods, the difference in performance was not as wide as could be anticipated. SVM with a linear kernel attained an F1-score of 0.856, which compares well with most deep learning models. This implies that for certain sentiment analysis tasks, well-tuned features using standard algorithms can prove to be a good alternative to more involved deep learning models, particularly where computational resources are limited.

D. Twitter Sentiment Analysis Challenges

Error analysis identified some recurring issues in Twitter sentiment analysis:

Sarcasm and Irony: In spite of our attempts to handle sarcasm using contextual features and deep learning models, sarcastic tweets continued to be a major source of errors.

Implicit Sentiment: Tweets tend to convey sentiment implicitly, and inference and world knowledge are needed. Our models performed poorly on such instances.

Mixed Sentiment: Tweets with both positive and negative sentiments were difficult to classify, especially when deciding on the overall sentiment.

Domain Adaptation: The variation in performance across datasets suggests that domain adaptation is still an issue.

IV. CONCLUSION AND FUTURE WORK

This paper gave a thorough review and application of NLP methods for Twitter sentiment analysis. We discussed different preprocessing techniques, feature extraction methods, and classification methods specifically tailored to the peculiar nature of Twitter data.

Our experiments showed the efficacy of Twitter-specialized preprocessing, the merits of feature combinations across different types of features, and the robustness of ensemble strategies for sentiment annotation.

Our work attained the current state of the art performance on our massive-scale Kaggle set of 1.8 million tweets, scoring an F1-score of 0.889.

Despite such accomplishments, sarcasm, implicit sentiment, mixed sentiment, and domain adaptation pose challenges. Addressing these issues in future studies should be accomplished through multimodal analysis, adding conversational context, fine-grained sentiment analysis, and cross-lingual methods.

Twitter sentiment analysis remains a useful instrument for gauging public opinion and attitude in real-time across many fields ranging from marketing to politics and public health. With the ongoing development of NLP methods, there is much to anticipate in terms of further accuracy and resilience in Twitter sentiment analysis systems.

As the platform evolves further, there will be a stronger feedback loop in place to guarantee that user needs are continuously met and that the platform evolves according to emerging trends in the job market. Incorporating multilingual capability and local customs to expand JobPathway's reach into international markets will also be given high priority. Finally, JobPathway has managed to position itself as wide-reaching in fulfilling the requirements for a comprehensive AI-powered talent acquisition solution, and with more advancements, it will be well-equipped to redefine the recruitment future for both recruiters and seekers.

V. REFERENCES

- [1] B. Liu, "Sentiment Analysis and Opinion Mining," *Synthesis Lectures on Human Language Technologies*, vol. 5, no. 1, pp. 1-167, 2012.
- [2] A. Giachanou and F. Crestani, "Like It or Not: A Survey of Twitter Sentiment Analysis Methods," *ACM Computing Surveys*, vol. 49, no. 2, pp. 1-41, 2016.
- [3] F. Atefeh and W. Khreich, "A Survey of Techniques for Event Detection in Twitter," *Computational Intelligence*, vol. 31, no. 1, pp. 132-164, 2015.
- [4] A. Agarwal, B. Xie, I. Vovsha, O. Rambow, and R. Passonneau, "Sentiment Analysis of Twitter Data," in *Proceedings of the Workshop on Languages in Social Media*, 2011, pp. 30-38.
- [5] A. Go, R. Bhayani, and L. Huang, "Twitter Sentiment Classification using Distant Supervision," *CS224N Project Report*, Stanford, vol. 1, no. 12, p. 2009, 2009.
- [6] F. H. Khan, U. Qamar, and S. Bashir, "SentiMI: Introducing Point-wise Mutual Information with SentiWordNet to Improve Sentiment Polarity Detection," *Applied Soft Computing*, vol. 39, pp. 140-153, 2016.
- [7] H. Saif, Y. He, and H. Alani, "Semantic Sentiment Analysis of Twitter," in *International Semantic Web Conference*, 2012, pp. 508-524.
- [8] T. Singh and M. Kumari, "Role of Text Pre-processing in Twitter Sentiment Analysis," *Procedia Computer Science*, vol. 89, pp. 549-554, 2016.
- [9] A. Pak and P. Paroubek, "Twitter as a Corpus for Sentiment Analysis and Opinion Mining," in *Proceedings of the Seventh International Conference on Language Resources and Evaluation*, 2010, pp. 1320-1326.
- [10] E. Kouloumpis, T. Wilson, and J. Moore, "Twitter Sentiment Analysis: The Good the Bad and the OMG!," in *Fifth International AAAI Conference on Weblogs and Social Media*, 2011.
- [11] A. Severyn and A. Moschitti, "Twitter Sentiment Analysis with Deep Convolutional Neural Networks," in *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2015, pp. 959-962.
- [12] S. Ghosh, O. Vinyals, B. Strope, S. Roy, T. Dean, and L. Heck, "Contextual LSTM (CLSTM) Models for Large Scale NLP Tasks," *arXiv preprint arXiv:1602.06291*, 2016.
- [13] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019, pp. 4171-4186.
- [14] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "RoBERTa: A Robustly Optimized BERT Pretraining Approach," *arXiv preprint arXiv:1907.11692*, 2019.
- [15] C. Sun, X. Qiu, Y. Xu, and X. Huang, "How to Fine-Tune BERT for Text Classification?," in *China National Conference on Chinese Computational Linguistics*, 2019, pp. 194-206.
- [16] D. Q. Nguyen, T. Vu, and A. T. Nguyen, "BERTweet: A Pre-trained Language Model for English Tweets," in *Proceedings of the 2020 Conference on Empirical Methods in Natural*