

# Sentiment Analysis of Tweets Using Traditional and Deep Learning Approaches

Estelle CADENE, Maela BRELIVET, Youri HALMEART, Lelio GUALINO

## Abstract

This study investigates sentiment classification in tweets through a comparative analysis of traditional machine learning algorithms and deep learning models. We examine the performance of Logistic Regression, Naive Bayes, and a fine-tuned BERT model on a dataset composed of 12,000 labeled tweets. Our results show that BERT significantly outperforms traditional approaches in terms of accuracy and macro-averaged F1-score, underscoring the advantages of contextualized language representations in sentiment analysis. These findings contribute to the growing body of literature that demonstrates the importance of deep learning models in natural language processing (NLP) tasks, particularly for nuanced textual data such as social media content.

## 1 Introduction

Sentiment analysis, also known as opinion mining, is a subfield of natural language processing that involves identifying and categorizing opinions expressed in text to determine the writer’s attitude toward a particular topic. With the explosion of user-generated content on social media platforms like Twitter, sentiment analysis has become a critical tool for understanding public opinion in real time. Applications span across domains including marketing, political forecasting, product review aggregation, and public health monitoring.

Twitter presents a unique set of challenges for sentiment analysis. Tweets are short, often informal, and filled with linguistic anomalies such as abbreviations, emoticons, hashtags, and user mentions. These characteristics make traditional NLP methods less effective, pushing the field toward the adoption of deep learning models that are better equipped to handle the complexity and variability of social media language.

In this paper, we focus on the task of multi-class sentiment classification of tweets, aiming to classify each tweet into one of three categories: positive, neutral, or negative. We explore both traditional machine learning models—Logistic Regression and Naive Bayes—and a state-of-the-art deep learning model, BERT. Our objective is to assess the performance trade-offs between these approaches and to understand the extent to which modern deep learning methods improve sentiment classification accuracy in this domain.

## 2 Related Work

Early work in sentiment analysis primarily relied on lexicon-based or statistical models. Pang and Lee (2008) utilized Naive Bayes and Support Vector Machines (SVM) trained on bag-of-words features, establishing baseline performances for text classification tasks. Similarly, Wilson et al. (2005) proposed rule-based approaches using sentiment lexicons, which were interpretable but limited in handling context or sarcasm.

With the rise of deep learning, models such as LSTM (Long Short-Term Memory networks) began to dominate sentiment analysis. Tang et al. (2014) showed that sentiment-specific word embeddings improved classification performance by capturing the emotional tone of words more effectively than general-purpose embeddings.

The advent of transformer-based models marked a significant leap in the capabilities of NLP systems. Devlin et al. (2019) introduced BERT (Bidirectional Encoder Representations from Transformers), which leverages bidirectional attention mechanisms and pre-training on massive corpora to capture rich linguistic context. BERT has since become a foundational model for numerous NLP tasks, including sentiment analysis.

Hybrid approaches have also emerged, combining rule-based methods with statistical learning to capitalize on the strengths of each. For instance, Mohammad and Turney (2013) demonstrated that blending lexicon-derived features with machine learning classifiers could yield improvements over either method in isolation.

### 3 Dataset and Preprocessing

Our analysis utilizes a dataset of approximately 12,000 tweets, each labeled with a sentiment category: positive, neutral, or negative. The class distribution is moderately imbalanced, with 40.5% of tweets labeled as neutral, 31.2% as positive, and 28.3% as negative. This distribution reflects the natural skew found in real-world social media datasets.

To prepare the data for model training, we applied several preprocessing steps. Tweets were first normalized by converting all text to lowercase. URLs, user mentions (e.g., "@username"), hashtags, and punctuation were removed to reduce noise. We also eliminated stop words—common words with little semantic value—and applied lemmatization to reduce words to their base or dictionary forms.

For traditional machine learning models, feature extraction was performed using Term Frequency–Inverse Document Frequency (TF-IDF) vectorization. We limited the number of features to the top 5,000 by frequency to manage dimensionality. For BERT, we utilized pre-trained bert-base-uncased embeddings and tokenized input sequences to a maximum length of 128 tokens, consistent with the model’s architecture.

## 4 Methodology

### 4.1 Traditional Models

We evaluated two widely used traditional classifiers: Logistic Regression and Multinomial Naive Bayes.

Logistic Regression is a linear classifier that models the probability of a class using the logistic function. It is particularly suitable for multi-class classification and offers the advantage of interpretability through analysis of feature weights. However, it assumes a linear relationship between input features and the log-odds of the target variable.

Naive Bayes, specifically the Multinomial variant, is a probabilistic classifier based on Bayes’ Theorem with the assumption of feature independence. While this assumption rarely holds in practice, the model has proven effective in text classification due to its efficiency and surprisingly robust performance on sparse feature representations.

### 4.2 Deep Learning Model

The deep learning model employed in this study is BERT (Devlin et al., 2019), a transformer-based architecture designed to pre-train deep bidirectional representations. BERT incorporates an attention mechanism that enables the model to weigh the importance of each word in a sentence relative to others, allowing it to capture subtle linguistic and semantic nuances.

We fine-tuned a pre-trained bert-base-uncased model on our tweet dataset. Fine-tuning involved re-training the final classification layer while preserving the pre-trained weights of earlier layers. Training was conducted over three epochs with a batch size of 16 and a learning rate scheduler for dynamic adjustment. The maximum sequence length was set to 128 tokens to accommodate tweet length while maintaining computational efficiency.

The training process was executed on a high-performance workstation equipped with a modern GPU, allowing us to complete the fine-tuning of the BERT model in approximately three hours. This highlights the computational demands of transformer-based models compared to traditional classifiers, which trained in under a minute.

## 5 Experimental Results and Discussion

### 5.1 Performance Metrics

Model evaluation was based on accuracy and macro-averaged F1-score, providing a balanced measure of performance across classes despite the mild class imbalance.

Model	Accuracy	Macro F1-score
BERT	76.0%	0.76
Logistic Regression	69.0%	0.69
Naive Bayes	64.0%	0.63

Table 1: Model performance comparison

The BERT model significantly outperformed the traditional approaches, achieving a 7-point improvement in accuracy and F1-score over Logistic Regression. This result highlights BERT’s ability to model complex language patterns, particularly in noisy and context-dependent domains such as Twitter.

### 5.2 Confusion Matrix

An analysis of the BERT confusion matrix reveals that the model is particularly effective at identifying neutral sentiments, with relatively higher misclassifications occurring between positive and negative categories. This behavior suggests that while BERT handles ambiguity better than traditional models, it still encounters challenges in resolving subtle differences in tone, especially when sarcastic or ironic language is used.

Actual \ Predicted	Positive	Neutral	Negative
Positive	1080	204	258
Neutral	498	1523	983
Negative	384	501	983

Table 2: BERT confusion matrix

### 5.3 Feature Interpretability

While deep learning models are often criticized for their lack of interpretability, traditional models provide insight into influential features. For instance, the most significant positive words identified by Logistic Regression included "love," "thanks," and "awesome," whereas "sad," "hate," and "sorry" were associated with negative sentiment. Neutral sentiment was harder to define but often included words like "happy," "good," and "enjoy," indicating possible misclassifications due to subtle positive connotations.

### 5.4 Limitations

Traditional models exhibited several limitations. Their reliance on bag-of-words representations renders them unable to capture word order, syntactic structure, or contextual meaning. As a result, they fail to detect sentiment in tweets that use sarcasm or irony, or that depend on inter-word relationships for interpretation. Moreover, their performance is heavily dependent on feature engineering and preprocessing quality.

BERT, by contrast, requires minimal feature engineering and is capable of contextual understanding. However, it is computationally intensive, requiring GPU resources and longer training times. Moreover, it remains a black-box model, limiting its use in applications that demand transparency and explainability.

## 6 Conclusion and Future Work

In this study, we demonstrated that deep learning models, specifically BERT, significantly outperform traditional machine learning approaches for sentiment analysis on Twitter data. The ability to understand context, handle non-standard language, and capture semantic relationships gives BERT a decisive edge in this domain. Nonetheless, traditional models remain useful for resource-constrained environments or applications requiring interpretability.

For future work, several extensions can be pursued. First, we aim to experiment with alternative transformer architectures such as RoBERTa, DistilBERT, and DeBERTa, which may offer improved performance or efficiency. Secondly, data augmentation techniques could be employed to expand the training set and mitigate class imbalance. This may include paraphrasing, synonym substitution, or back-translation.

Further, incorporating additional features such as emoji analysis, syntactic parsing, and temporal information may improve classification accuracy. We also plan to explore ensemble methods that combine the strengths of multiple models, potentially yielding better generalization.

Finally, we envision deploying this system in a real-time setting to monitor sentiment dynamics on social media. Integration with a streaming data pipeline and the development of a dashboard for visualization could provide valuable tools for marketers, policy analysts, and researchers alike.

## References

- Pang, B., & Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and Trends® in Information Retrieval*, 2(1-2), 1–135.
- Wilson, T., Wiebe, J., & Hoffmann, P. (2005). Recognizing contextual polarity in phrase-level sentiment analysis. *Proceedings of the Human Language Technology Conference and the Conference on Empirical Methods in Natural Language Processing*, 347–354.
- Tang, D., Qin, B., & Liu, T. (2014). Learning sentiment-specific word embedding for Twitter sentiment classification. *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, 1555–1565.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805*.
- Mohammad, S. M., & Turney, P. D. (2013). Crowdsourcing a word–emotion association lexicon. *Computational Intelligence*, 29(3), 436–465.