



A Comparative Analysis of Feature Extraction Techniques for Fake Reviews Detection

Zahraa Fadhel ^{1*}, Hussien Attia¹, Yossra Hussain Ali²

¹Department of Computer Sciences, College of Science for Women, University of Babylon, Babylon, Iraq

²Department of Computer sciences, University of Technology, Baghdad, Iraq

Email: zahraa.alkhafaji.jsci140@student.uobabylon.edu.iq; wsci.husein.attia@uobabylon.edu.iq;
Yossra.H.Ali@uotechnology.edu.iq

Abstract

The current Internet era is characterized by the widespread circulation of ideas and viewpoints among users across many social media platforms, such as microblogging sites, personal blogs, and reviews. Detecting fake reviews has become a widespread problem on digital platforms, posing a major challenge for both consumers and businesses. Due to the ever-increasing number of online reviews, it is no longer possible to manually identify fraudulent reviews. Artificial intelligence (AI) is essential in addressing the problem of identifying fake reviews. Feature extraction is a crucial stage in detecting fake reviews, and successful feature engineering techniques can significantly improve the accuracy of opinion extraction. The paper compares five feature extraction methods for multiple opinion classification using Twitter on airline and Borderland game reviews. FastText with X-GBoost classifier outperformed all other techniques, achieving 94.10% accuracy on the airline dataset and 100% accuracy in Borderland game reviews.

Keywords: Feature extraction; Fake reviews; Natural language processing; FastText; X-GBoost

1. Introduction

The proliferation of social media platforms has provided individuals with the opportunity to publicly participate and express their ideas [1]. Consequently, online reviews have become more prevalent. Some of these reviews intentionally mislead readers or opinion mining systems, either by giving undeserved positive reviews to certain entities to promote their products or by giving unfair or harmful negative reviews to other entities to damage their reputation. Incorrect reviews are also referred to as fake or false reviews [2]. Natural Language Processing (NLP) focuses on the detection of bogus reviews. The objective is to distinguish, evaluate, and categorize any given review in the form of positive, negative, or natural[3]. Twitter data is a compelling choice for sentiment analysis due to the vast amount of information it offers. The volume of data available is truly immense, with over 500 million tweets sent

daily [4]. Twitter attracts a diverse group of individuals from various age groups, including a significant proportion of corporate executives from multiple countries who actively use the platform for social media purposes [5]. Given the escalating quantity of online reviews, it has become impracticable to manually detect deceptive reviews. Artificial intelligence and machine learning offer automatic and scalable techniques for detecting bogus reviews using algorithmic and data analysis capabilities [7]. The sentiment analysis and assessment of tweets or classified opinions is contingent on the quantity of data and the kind of document [3]. Feature engineering techniques play a vital role in machine learning by facilitating the identification and detection of fraudulent reviews[4]. These techniques aid in the recognition and extraction of pertinent information from review materials. In this paper, The analysis primarily focuses on various feature engineering techniques and their functionality in different classification algorithms, as feature extraction significantly impacts accuracy.

The primary contribution of this paper is:

1. Comprehensive evaluation of various machine learning and feature extraction techniques for fake reviews detection.
2. Emphasizes the significance of employing feature extraction methodologies to enhance sentiment analysis and opinion classification.
3. An explanation of what features are and how are they extracted and calculated. We also evaluate the effectiveness of the features for the existing methods to identify the best features in the fake reviews detection.

The remaining sections of the paper are organized in the following manner: section1:

Introduction; Section 2: Literature Review; Section 3: Connecting the various theoretical backgrounds Section 4 presents a system framework; and Section 5 presents the model outline and work. Section six is the performance evaluation Section seven shows the conclusion and future work.

2. Related Work

Usman Naseem et al [5] . attempted for multiclass sentiment analysis on COVID-19, The researchers examined the efficacy of several characteristics and supervised learning methods for their datasets., their results in comparison to DT, SVM, and NB, RF showed better performance for fast- Text. Convolutional Neural Network (CNN) [6] has been used for text classification on six benchmark datasets including Ag News, Amazon Full and Polarity, Yahoo Question Answer, Yelp Full, and Polarity with FastText model. Muhammad Umer et al. conducted experiment and it is obvious that Fasttext comes out with highest result by accuracy of 0.96 on Ag News and 0.90 f1-score.

SVM has been used for sentiment analysis on SatuSehat app user reviews with TF-IDF as feature extraction. Shahmirul Hafizullah Imanuddin et al. [3] their study indicates that accuracy testing of SVM yielded 91% with a positive sentiment of 92% precision, 71% recall, and 80% F1-score. In contrast, negative sentiment had 90% precision, 98% recall, and 94% F1-score. Koosha Sharif ani et al. [7] presents various ensemble techniques to perform the binary classification of news articles along with Naïve Bayes, Passive Aggressive, and SVM Classifiers, TF-IDF feature extraction has been used. According to their developed system that Passive Aggressive classifier outperforms with accuracy up to 93%. Saima Sadiq et al. [8] Introduced a simple deep learning approach that incorporates word embedding by utilizing the publicly available Tweepfake database. A conventional Convolutional Neural Network (CNN) approach is created, utilizing FastText word embedding to facilitate the task of identifying deep bogus tweets. This investigation included many machine learning approaches as a fundamental approach for comparison. These fundamental techniques utilize distinct characteristics, such as subwords FastText, FastText, Term Frequency-Inverse, Term Frequency, and Document Frequency. In 2023 [9] Afifah Mohd Asri et al. employed the Text Blob feature extraction technique and ACO, GA, and PSO to classify text in the Drug reviews datasets. The results showed that PSO demonstrated the highest levels of performance, with an average of 49.3% for precision, 73.6% for recall, 59% for F- score, and 57.2% for accuracy. Logistic Regression (LR) and Neural Network (NN) classifiers has been used by Kona Thala Karthikeya et al. [10]. For classifying tweets of Borderlands game reviews with FastText techniques. Their results showed that The LR model achieved an accuracy of 95%, while the NN model achieved an accuracy of 93%.

3. Method and material

Figure 1 represents the entirety of the study process. Initially, two datasets comprising reviews from both Airline and Borderlands games are applied. These datasets were subsequently subjected to pre-processing techniques in order to obtain sanitized text. Several key preprocessing stages include reducing the text to lowercase, encoding it, eliminating stop words, and transforming it into numerical numbers. These words can also be quantitatively represented using five distinct methods: TF-IDF, N-Gram, W2V, Glove, and FastText. Three standard supervised learning models, including LG, ETC, and X-GBoost, were utilized to determine the polarity of the two datasets.

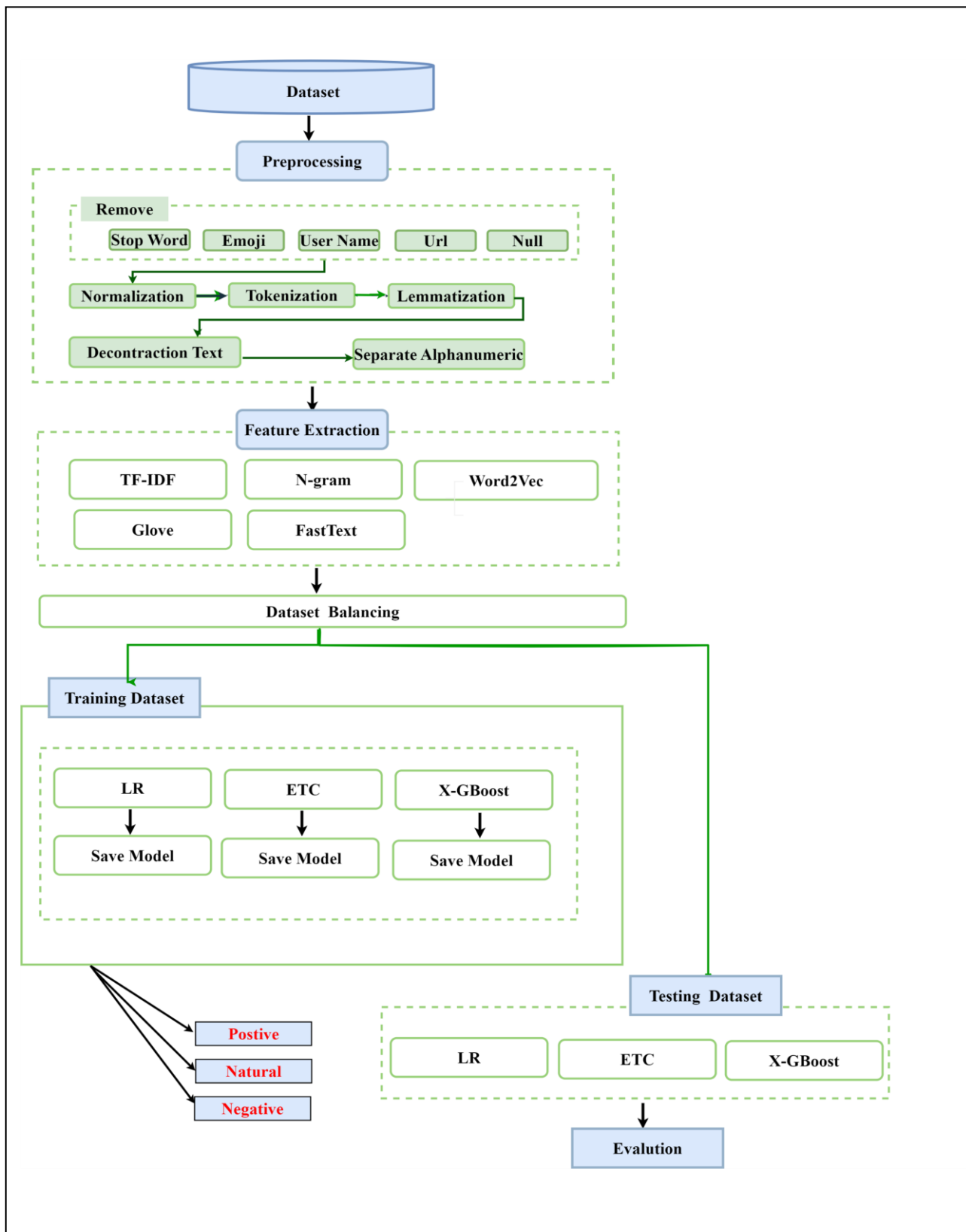


Figure 1. An overview of our approach.

3.1 Dataset collection

To conduct our research, two different datasets are used: first, The Kaggle dataset "CrowdFlower Twitter US Airline Sentiment" [20], consisting of 14,640 tweets. 2,363 positive values 9,178 negative values, and 3,099 neutral values, second Borderlands game reviews [21]. Consisted of 74682 tweets. The emotional tone conveyed in the text was classified as "Positive," 28%, "Negative," 30% or "Neutral," 42%.

3.2 Preprocessing

3.2.1 Removal of punctuation

In order to preserve the integrity of the data, it is necessary to remove any instances of null, missing values, and special characters such as punctuation from the dataset. By utilizing the drop null on the borderland dataset [14]. The dataset originally contained 74,681 instances. After removing null values, the dataset was reduced to 73,995 instances. While twitter airline lack of null value. None of these factors are essential for opinion classification [11].

3.2.2 Normalization

This technique employs the consolidation of many instances of a certain letter by eliminating all symbols and digits, and converting all letters to either uppercase or lowercase. [12]. The described operations entail the elimination of all numerical values. - Eliminate all punctuation marks and symbols, including exclamation marks, question marks, hashtags, and dollar signs, at marks, asterisks, square brackets, curly braces, equal signs, percentage signs, ampersands, parentheses, hyphens, underscores, quotation marks, plus signs, single quotes, and forward slashes. - Transform each word into lowercase (LC).

3.2.3 Tokenization

It is a basic step in preprocessing [13] which is in fact a process of segmenting an entire raw text document into smaller parts and each of these parts is known as tokens.

3.2.4 Lemmatization

It denotes to the process of reducing the inflected and derived forms of a word to its base or stem of the word. In this lemmatization process [14], The term used to denote the original form of a word is called the lemma. This procedure is more accurate than stemming, which is not reliant on dictionary definitions..

3.2.5 DE contraction text

Expanding contractions in text data ensures consistency in analysis by converting shortened forms to their full forms [15]. For instance, "can't" can be expanded to "cannot" to maintain uniformity and accuracy in natural language processing tasks.

3.2.6 Separate alphanumeric

Separating alphanumeric characters into separate entities improves text data analysis and processing [16]. For instance, in the example given, the original text is "The price of the product is \$99.99 and it received a 5-star rating." The resulting text is "\$ 99.99" and the "5-star" rating is "5 - star." This separation enhances clarity and facilitates better analysis of text data.

3.3 Feature extraction

Feature extraction is a vital process in NLP and ML. technologies such as TF-IDF, N-gram, Word to Vector, Glove, and FastText are applied to convert raw text data into parseable formats [17]. These techniques help capture the semantic relationships, context, important features of textual data, and enhance machine learning models performance [1].

3.3.1 Term Frequency-Inverse Document Frequency (TF-IDF)

is a metric used to determine the relevance or uniqueness of a word in a document, relative to a collection of documents [4]. Term frequency refers to the frequency of a term occurring in a document, while inverse document frequency measures the commonness of a word across all documents. Collection of documents [10]. The desired TF-IDF is obtained by multiplying these two-frequency metrics. The objective of this model is to retrieve papers that contain comparable words.

3.3.2. N-gram

An N-gram is a sequence of n words, letters, or tokens which describes contextual information and associations by examining their frequency distributions. [18].

$$N - \text{gram}(w_1, w_2, \dots, w_n) = \frac{\text{Frequency of } N - \text{gram}(w_1, w_2, \dots, w_n)}{\text{Total number of } N - \text{grams}} \dots (1)$$

N-gram (w_1, w_2, \dots, w_n) represents a sequence of N consecutive words.

(Frequency of N-gram (w_1, w_2, \dots, w_n)) denotes the number of times the specific N-gram (w_1, w_2, \dots, w_n) occurs in the text or dataset.

(Total number of N-grams) represents the overall count of N-grams present in the text or dataset.

3.3.3 Word2vec (Word to Vector)

Google introduced Word2vec in 2013, a neural network-based word embedding technology that uses continuous bag of words (CBOW) and skip-gram learning models to predict words and phrases [19]. The system uses a large collection of text to create a vector space with multiple dimensions, assigning a unique vector to each word, and enhancing the probability of predicting context words [20].

$$\text{maximize} \sum_{t=1}^T \sum_{c \leq j < c+j \neq 0} \log p(w_{t+j} | w_t) \dots (2)$$

Where:

w_t represents the target word.

$w_{(t+j)}$ represents the context word.

c represents the context window size.

$\log p(w_{(t+j)} | w_t)$ represents the conditional probability of observing the context word $w_{(t+j)}$ given the target word w_t .

3.3.4 Global Vectors for Word Representation (Glove)

Glove is an unsupervised learning algorithm that generates vector representations of words based on global word frequency statistic[20]. It generates word embeddings by combining a global word-word co-occurrence matrix from a corpus, exhibiting linear substructures. The model can be trained rapidly on larger data sets and minimizes a loss function [21].

$$J = \sum_{i,j=1}^V f(P_{ij})(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log P_{ij})^2 \dots (3)$$

Where:

J represents the loss function.

f is a weighting function.

P_{ij} represents the probability of word i appearing in the context of word j

w_i and \tilde{w}_j are the word vectors.

$b_i + \tilde{b}_j$ are the bias terms.

V represents the vocabulary size.

3.3.5 FastText

Facebook's FastText algorithm uses n-grams to create vector representations for words not in the lexicon, producing 300-dimensional token vectors for tweets [22]. It expands Word2Vec capturing morphological details by expressing words as vector sums of n-grams, with supplementary components for sub words[23].

3.4 Data balancing:

class imbalance refers to an uneven distribution of data between classes [30], which can lead to biased behaviour in machine learning models. SMOTE is a synthetic oversampling technique that aims to address class imbalance by generating synthetic examples for the minority class. Instead of simply replicating existing instances, SMOTE selects minority class examples and creates synthetic instances along line segments connecting their nearest

neighbours. By oversampling the minority class, SMOTE helps create a more balanced class distribution, ensuring fair treatment of all classes and preventing models from ignoring the under-represented class [24].

3.5 Classifier

Our research methodology involved the utilization of three distinct supervised classifiers, namely X-Boost, ETC, and LR, to explore the effectiveness of various feature extraction techniques.

3.5.1. Extreme Gradient Boosting (X-Boost)

is a powerful machine learning technique used for predictive modelling. It combines gradient boosting, regularization methods, tree construction improvements, and parallel computation. It uses weak decision tree models to create a powerful ensemble, aiming to reduce prediction errors [12]. X-Boost can adjust learning rates, apply regularization techniques, and manage missing values. It's widely used for classification, regression, and ranking problems due to its efficiency and scale, resulting in high predicted accuracy. It's also applicable in fields like banking, healthcare, and natural language processing [25].

3.5.2. Extra Tree Classifier (ETC)

ETC utilizes a meta-estimator to train a specific number of weak learners (randomized decision trees) on different samples of the dataset[17]. This process enhances the accuracy of the predictions. It is a type of ensemble learning model that is used for categorization, similar to RF (Random Forest). The sole distinction between ETC and RF is in the manner in which trees are formed inside the forest. ETC utilizes the original training sample to produce decision trees, whereas RF constructs decision trees using bootstrap samples selected from the original dataset [26]. At every test node, each tree is given a random subset of k characteristics selected from the feature-set. Every decision tree must choose the optimal feature to divide the data, using a mathematical criterion, usually the Gini Index. This arbitrary selection of characteristics results in the formation of several decision trees that are not correlated with each other[27]. ETC classifier is used with two main parameters, $n_estimators = 100$ and $max_depth = None$.

3.5.3. Logistic Regression (LR):

LR is a linear model that applies the concepts of linear regression analysis. Logistic regression is a statistical model that utilizes a logistic sigmoidal function to convert its output into a discrete probability value[28]. The sigmoid function ensures that its output is constrained to a range of 0 to 1. The maximum likelihood estimation (MLE) approach is employed to determine the optimal fit line for coefficients, and the process is iterated until the Log Likelihood (LL) stabilizes.[29].

3.6 Evaluation metrics

The F1 accuracy is used to assess the efficacy of TF-IDF, N-gram, W2V, Glove, and FastText. The accuracy and F1 formulae are represented as:

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (4)$$

where, TP, TN, FP, and FN indicate true positive, true negative, false positive and false negative respectively.

$$\text{F1} = 2 * \frac{((\text{TP}/(\text{TP} + \text{FP})) * (\text{TP}/(\text{TP} + \text{FN})))}{((\text{TP}/(\text{TP} + \text{FP})) + (\text{TP}/(\text{TP} + \text{FN})))} \quad (5)$$

4. Experiment result

This section provides empirical results that specifically determine the performance of various feature extraction methods on the opinion classification of airlines and Borderlands 'game reviews. The accuracy is provided in Table I. However, the graphical representation of the accuracy is presented in Fig. 2 and Fig. 3 for the airline and borderlands' game reviews datasets respectively.

Table 1: Accuracy result in percentage for airline and borderlands' game dataset.

Feature extraction	Airline			Boderland's 'game		
	X-Gboost	ETC	LG	X-Gboost	ETC	LG
TF-IDF	0.90	0.94	0.86	0.72	0.90	0.67
N-gram	0.44	0.44	0.43	0.45	0.445	0.442
W2V	0.939	0.932	0.83	0.82	0.87	0.64
Glove	0.85	0.90	0.69	0.71	0.81	0.58
FastText	0.94	0.93	0.83	1.00	0.99	1.00

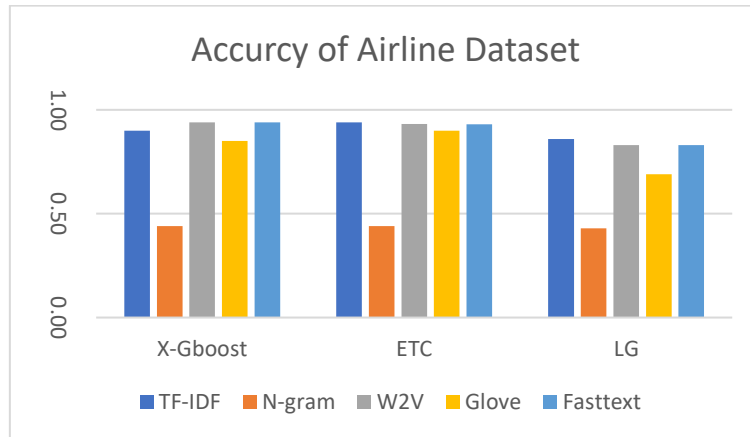


Figure 2. Accuracy for airlines review dataset.

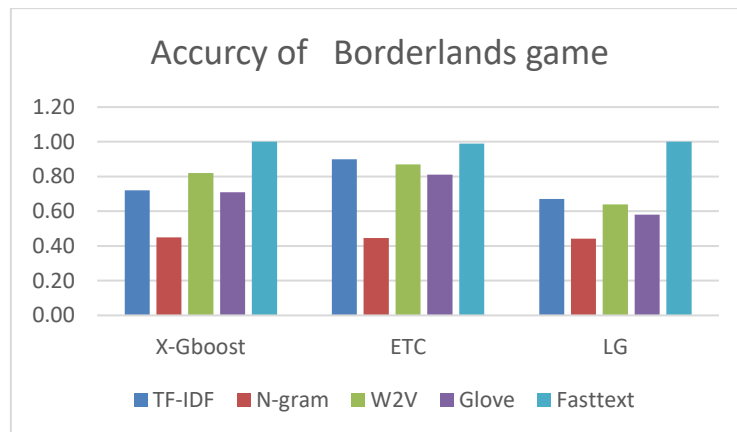


Figure 3. Accuracy for borderland's' game review dataset.

4.1. TF-IDF

Initially, the TF-IDF approach is explored, which yielded 300 features for both the airline and borderlands' review datasets. It has been confirmed that ETC achieved the highest accuracy rates of 94% and 90% for airline and borderlands' review datasets, respectively.

4.2. N-gram

N-gram feature extraction technique has been applied for both datasets which extract 1 categorical feature and it has been explored that ETC and X-GBoost achieve the highest accuracy of 0.44% for airline and Borderlands 'game reviews datasets, whereas LG achieves the best accuracy of 43% for airline review dataset and 0.44 for Borderlands 'game dataset.

4.3. W2V

In this time, W2V which extract 300 features was indicated for both airline and borderlands 'game review datasets. it has been inspected that X-Boost achieves the best accuracy of 0.939 % for airline dataset. While ETC classifier attains the best accuracy of 0.87.

4.4. Glove

Glove feature extraction is considered for both airline and Borderlands' game datasets yielded 25 features .it is confirmed that ETC classifier has been able to gain the best accuracy by securing 0.90% and 0.81 respectively.

4.5. FastText

Finally, the utilization of FastText has resulted in the extraction of 300 features. The Boost classifier achieved the highest accuracy of 0.100 for both datasets, while LG also achieved 0.100 on the Borderlands dataset.

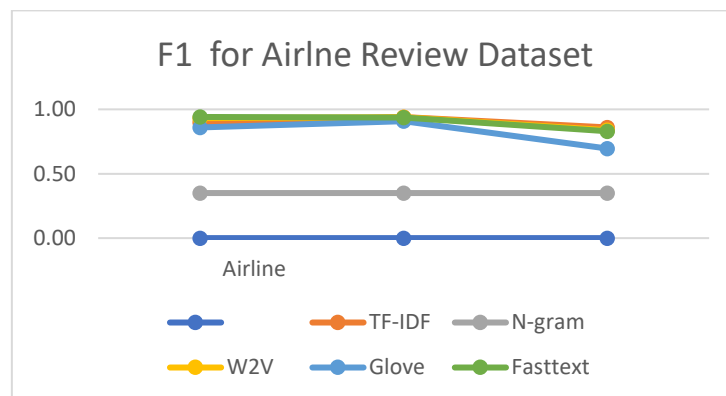


Figure 4. F1 for airline review dataset

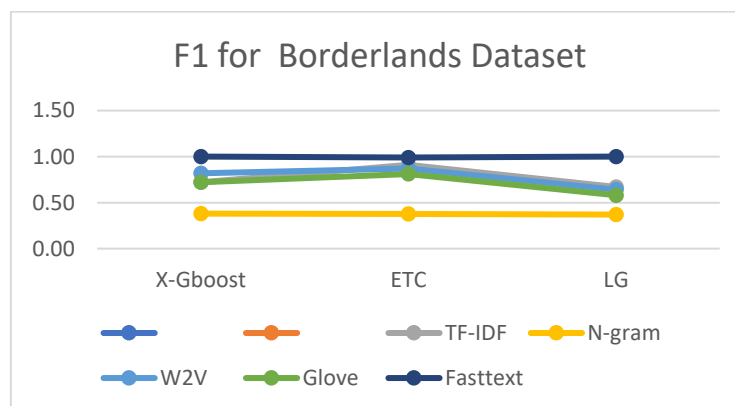


Figure 5. F1 for airline review dataset

The bar graph depicted in Figure 4 displays the F1 score for the airline review dataset, while Figure 5 represents the F1 score for the borderlands game review dataset. Based on the results, FastText achieves the maximum F1 score when using X-Boost and LR on the Borderlands game review dataset. Conversely, the lowest F1 score is obtained while using N-gram with LR. while evaluating airline reviews, FastText achieves the highest F1 score while using X-Boost, whereas N-gram yields the lowest F1 score when using LG.

Utilizing FastText technology. The analysis showed that the model achieved a precision rate of 95%, whilst the NN model achieved a precision rate of 93%.

Upon careful examination, it becomes evident that FastText demonstrated superior performance in terms of accuracy and F1 score for both the airline review dataset (using LR) and the Borderlands game review dataset. Additionally, FastText also outperformed other models in terms of accuracy and F1 score when using LG. Once again, the performance of both datasets confirms that the suggested technique does not suffer from either overfitting or underfitting.

5. Conclusion

Our study utilized five extraction strategies (TF-IDF, N-gram, W2V, Glove, and FastText) on two distinct data sets. A comparison analysis was conducted to determine the suitability of these techniques for text categorization using machine learning models (X-GBoost, ETC, and LG). A proposed approach aims to efficiently search for the optimal feature selection and machine learning algorithms for multi-emotion classification, based on accuracy and F1 metrics. This is apparent from the outcomes, since certain methods of extracting features and classifiers have poor performance for both airline game evaluations and border regions. The empirical findings indicate that fastText outperformed the other evaluated approaches in effectively representing non-vocabulary words. Additionally, the XGBoost algorithm demonstrated its strength in handling complicated and high-dimensional datasets. This feature clearly accounts for its superior performance in comparison to logistic regression (LR) and additive trees (ET) classifiers. • Acknowledges the potential of artificial intelligence and machine learning in combating bogus reviews. The objective of this study is to assess the effectiveness of various feature sets by employing the suggested machine learning classifiers.

6. Comparison with previous works

The table below indicates that the suggested system outperforms a group of earlier studies that created the classification system using the some of the same techniques and datasets.

Table 2: Comparison with previous works.

No.	Ref & Year	Dataset	Feature extraction	Classifiers	Accuracy
1	[5] 2021	Tweets on COVID-19.	TF-IDF	SVM	0.845
				RF	0.841
				NB	0.773
				DT	0.794
			W2V	RF	0.769
				DT	0.766
			Glove	RF	0.726
				DT	0.701
			FastText	SVM	0.815
				NB	0.735
			RF	0.845	
2	[7] 2022	News	TF-IDF	Naïve Bayes, Passive Aggressive, and SVM	Passive Aggressive classifier outperforms with accuracy up to 93%.
3	[6] 2023	Ag News, Amazon Full, Polarity, Yahoo Question Answer, and Yelp Full.	FastText	CNN	accuracy of 0.96 and 0.90 f1-score on Ag News
4	[8] 2023	Twitter	TF, TF-IDF, FastText.	DT, LR, AC, SGC, RF, GBM, ETC, NB, CNN, LSTM, and CNN-LSTM.	CNN worked along with FastText. It exhibited exceptional performance compared to previous combinations of deep learning and machine learning models, achieving a precision of 0.93, a recall of 0.95, and an F1 score of 0.93.

5	[9] 2023	Drug reviews.	Text Blob.	ACO, GA, and PSO.	(PSO) algorithm demonstrated superior performance, achieving an average precision of 49.3%, recall of 73.6%, F-score of 59%, and precision of 57.2%.
6	[3] 2023	SatuSehat app user reviews	TF-IDF	SVM	(SVM) achieved a positive sentiment accuracy of 91% with a precision of 92%, a recall of 71%, and an F1 score of 80%. Conversely, negative emotions exhibited a 90% accuracy rate, 98% recall rate, and 94% F1 score.
7	[10] 2023	Tweets of Borderlands game reviews.		LR, NN.	95% 93%
Proposed	2024	DS1	TF-IDF	ETC	0.94
			N-gram	RC	0.45
			W2V	XGBoost	0.93
			Glove	ETC	0.94
			Fasttext	XGBoost	0.94
		DS2	TF-IDF	ETC	0.90
			N-gram	XGBoost	0.45
			W2V	Knn	0.88
			Glove	Knn	0.81
			Fasttext	XGBoost	0.100

7. Future Works

In order to optimize the performance of the proposed system for future endeavours the following tactics can be employed:

1. Develop immediate implementation of the suggested technology for utilization on social media and various platforms.
2. Employing deep learning techniques on a vast dataset to enhance the process of extracting features or classifying data.
3. The objective is to develop a hybrid feature model by merging multiple extraction methods in order to generate a versatile model with a wide range of characteristics. Our objective is to evaluate tweets that are neutral in nature, devoid of any positive or negative sentiments. We primarily focus on Twitter data but may also include data from other social networking platforms.

References

- [1] "Comparative Analysis of Feature Extraction," pp. 1–13, 2022.
- [2] B. Liu, "Sentiment analysis and subjectivity," *Handb. Nat. Lang. Process. Second Ed.*, no. January 2010, pp. 627–666, 2010.
- [3] S. H. Imanuddin, K. Adi, and R. Gernowo, "Sentiment Analysis on Satushat Application Using Support Vector Machine Method," vol. 5, no. 3, pp. 143–149, 2023.

- [4] M. Syamala and N. J. Nalini, "A filter based improved decision tree sentiment classification model for real-time amazon product review data," *Int. J. Intell. Eng. Syst.*, vol. 13, no. 1, pp. 191–202, 2020, doi: 10.22266/ijies2020.0229.18.
- [5] U. Naseem, I. Razzak, M. Khushi, P. W. Eklund, and J. Kim, "COVIDSenti: A Large-Scale Benchmark Twitter Data Set for COVID-19 Sentiment Analysis," *IEEE Trans. Comput. Soc. Syst.*, vol. 8, no. 4, pp. 976–988, 2021, doi: 10.1109/TCSS.2021.3051189.
- [6] M. Umer *et al.*, "Impact of convolutional neural network and FastText embedding on text classification," *Multimed. Tools Appl.*, vol. 82, no. 4, pp. 5569–5585, 2023, doi: 10.1007/s11042-022-13459-x.
- [7] K. Sharifani, M. Amini, Y. Akbari, and J. A. Godarzi, "Operating Machine Learning across Natural Language Processing Techniques for Improvement of Fabricated News Model," *Int. J. Sci. Inf. Syst. Res.*, vol. 12, no. 9, pp. 20–44, 2022, [Online]. Available: <https://www.researchgate.net/publication/364340252>
- [8] S. Sadiq, T. Aljrees, and S. Ullah, "Deepfake Detection on Social Media: Leveraging Deep Learning and FastText Embeddings for Identifying Machine-Generated Tweets," *IEEE Access*, vol. 11, no. September, pp. 95008–95021, 2023, doi: 10.1109/ACCESS.2023.3308515.
- [9] A. M. Asri, S. R. Ahmad, and N. M. M. Yusop, "Feature Selection using Particle Swarm Optimization for Sentiment Analysis of Drug Reviews," *Int. J. Adv. Comput. Sci. Appl.*, vol. 14, no. 5, pp. 286–295, 2023, doi: 10.14569/IJACSA.2023.0140530.
- [10] K. Karthikeya, "Sentiment Analysis of Tweets Using Logistic Regression and Neural Networks with Emojis and Emoticons Sentiment Analysis of Tweets Using Logistic Regression and Neural Networks with Emojis and Emoticons," pp. 0–6, 2023.
- [11] D. Sameh, G. Khoriba, and M. Haggag, "Behaviour analysis voting model using social media data," *Int. J. Intell. Eng. Syst.*, vol. 12, no. 2, pp. 211–221, 2019, doi: 10.22266/IJIES2019.0430.21.
- [12] W. H. Asaad, R. Allami, and Y. H. Ali, "Fake Review Detection Using Machine Learning," *Rev. d'Intelligence Artif.*, vol. 37, no. 5, pp. 1159–1166, 2023, doi: 10.18280/ria.370507.
- [13] T. Hasan, A. Matin, M. Kamruzzaman, S. Islam, and M. O. F. Goni, "A Comparative Analysis of Feature Extraction Methods for Human Opinion Grouping Using Several Machine Learning Techniques," *Proc. 2020 IEEE Int. Women Eng. Conf. Electr. Comput. Eng. WIECON-ECE 2020*, pp. 272–275, 2020, doi: 10.1109/WIECON-ECE52138.2020.9398025.
- [14] M. Qorib, T. Oladunni, M. Denis, E. Ososanya, and P. Cotae, "Covid-19 vaccine hesitancy: Text mining, sentiment analysis and machine learning on COVID-19 vaccination Twitter dataset," *Expert Syst. Appl.*, vol. 212, no. September 2022, p. 118715, 2023, doi: 10.1016/j.eswa.2022.118715.
- [15] M. Siino, I. Tinnirello, and M. La Cascia, "Is text preprocessing still worth the time? A comparative survey on the influence of popular preprocessing methods on Transformers and traditional classifiers," *Inf. Syst.*, vol. 121, no. December 2023, p. 102342, 2024, doi: 10.1016/j.is.2023.102342.
- [16] S. Bird, E. Klein, and E. Loper, *LIVRO: cookbook Natural Language Processing with Python*, vol. 28, no. 4. 2009. [Online]. Available: <https://www.oreilly.com/library/view/natural-language-processing/9780596803346/>
- [17] K. L. ANUSHA, P. and PRASAD, "Survey on Fake Online Reviews Using Machine Learning Algorithms," *J. Crit. Rev.*, vol. 7, no. 18, pp. 2752–2758, 2020.
- [18] S. N. Alsubari *et al.*, "Data analytics for the identification of fake reviews using supervised learning," *Comput. Mater. Contin.*, vol. 70, no. 2, pp. 3189–3204, 2022, doi: 10.32604/cmc.2022.019625.
- [19] L. Ma and Y. Zhang, "Using Word2Vec to process big text data," *Proc. - 2015 IEEE Int. Conf. Big Data, IEEE Big Data 2015*, pp. 2895–2897, 2015, doi: 10.1109/BigData.2015.7364114.
- [20] HANSON ER, "Musicassette Interchangeability. the Facts Behind the Facts," *AES J. Audio Eng. Soc.*, vol. 19, no. 5, pp. 417–425, 1971.
- [21] A. Yadav and D. K. Vishwakarma, "Sentiment analysis using deep learning architectures: a review,"

- Artif. Intell. Rev.*, vol. 53, no. 6, pp. 4335–4385, 2020, doi: 10.1007/s10462-019-09794-5.
- [22] A. Joulin, E. Grave, P. Bojanowski, and T. Mikolov, “Bag of tricks for efficient text classification,” *15th Conf. Eur. Chapter Assoc. Comput. Linguist. EACL 2017 - Proc. Conf.*, vol. 2, pp. 427–431, 2017, doi: 10.18653/v1/e17-2068.
- [23] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, “Enriching Word Vectors with Subword Information,” *Trans. Assoc. Comput. Linguist.*, vol. 5, pp. 135–146, 2017, doi: 10.1162/tacl_a_00051.
- [24] J. Park, S. Kwon, and S. P. Jeong, “A study on improving turnover intention forecasting by solving imbalanced data problems: focusing on SMOTE and generative adversarial networks,” *J. Big Data*, vol. 10, no. 1, 2023, doi: 10.1186/s40537-023-00715-6.
- [25] J. Yao, Y. Zheng, and H. Jiang, “An Ensemble Model for Fake Online Review Detection Based on Data Resampling, Feature Pruning, and Parameter Optimization,” *IEEE Access*, vol. 9, pp. 16914–16927, 2021, doi: 10.1109/ACCESS.2021.3051174.
- [26] S. Hakak, M. Alazab, S. Khan, T. R. Gadekallu, P. K. R. Maddikunta, and W. Z. Khan, “An ensemble machine learning approach through effective feature extraction to classify fake news,” *Futur. Gener. Comput. Syst.*, vol. 117, pp. 47–58, 2021, doi: 10.1016/j.future.2020.11.022.
- [27] F. Rustam, M. Khalid, W. Aslam, V. Rupapara, A. Mehmood, and G. S. Choi, “A performance comparison of supervised machine learning models for Covid-19 tweets sentiment analysis,” *PLoS One*, vol. 16, no. 2, pp. 1–23, 2021, doi: 10.1371/journal.pone.0245909.
- [28] D. R. Nahma, “Patient Opinion Mining : Analysis of Patient Drugs Satisfaction using Support Vector Machine and Logistic Regression algorithm ضميرملا اضر ليلحت : ضميرملا يار نيدعت يتسجوللا رادحنلاا قيمزراوخو ” vol. 12, no. 2, pp. 164–171, 2020.
- [29] M. Avinash and E. Sivasankar, *A study of feature extraction techniques for sentiment analysis*, vol. 814. Springer Singapore, 2019. doi: 10.1007/978-981-13-1501-5_41.