Jana Publication & Research

Comparing Sentiment Analysis Methods: Flipkart Reviews



BioTech



Institut Seni Indonesia Surakarta

Document Details

Submission ID

trn:oid:::1:3419507024

Submission Date

Nov 21, 2025, 12:05 PM GMT+7

Download Date

Nov 21, 2025, 3:38 PM GMT+7

File Name

IJAR-54888.pdf

File Size

943.4 KB

16 Pages

5,309 Words

28,865 Characters



18% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

Filtered from the Report

- Bibliography
- Quoted Text

Match Groups

57 Not Cited or Quoted 14%

Matches with neither in-text cit:

Matches with neither in-text citation nor quotation marks

14 Missing Quotations 4%

Matches that are still very similar to source material

0 Missing Citation 0%

Matches that have quotation marks, but no in-text citation

• 0 Cited and Quoted 0%

Matches with in-text citation present, but no quotation marks

Top Sources

12% 📕 Publications

3% Land Submitted works (Student Papers)





Match Groups

57 Not Cited or Quoted 14%

Matches with neither in-text citation nor quotation marks

14 Missing Quotations 4%

Matches that are still very similar to source material

0 Missing Citation 0%

Matches that have quotation marks, but no in-text citation

• 0 Cited and Quoted 0%

Matches with in-text citation present, but no quotation marks

Top Sources

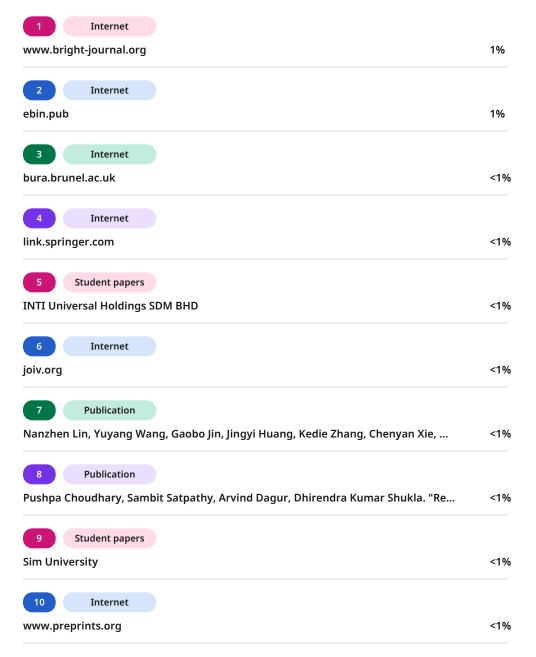
13% 🌐 Internet sources

12% 📕 Publications

3% Land Submitted works (Student Papers)

Top Sources

The sources with the highest number of matches within the submission. Overlapping sources will not be displayed.







11 Internet	
www.jatit.org	<1%
12 Internet	
www.mdpi.com	<1%
13 Publication	
Sushil Kamboj, Pardeep Singh Tiwana. "Innovations in Computing", CRC Press, 2025	<1%
14 Student papers	
University of South Wales - Pontypridd and Cardiff	<1%
15 Publication	
"Requirements Engineering: Foundation for Software Quality", Springer Science a	<1%
16 Publication	
Arvind Dagur, Sohit Agarwal, Dhirendra Kumar Shukla, Shabir Ali, Sandhya Sharm	<1%
17 Publication	
Cristina M. Quintella, Pamela D. Rodrigues, Jorge L. Nicoleti, Samira A. Hanna. "Co	<1%
18 Internet	
americaspg.com	<1%
19 Student papers	
Glasgow Caledonian University	<1%
20 Student papers	
University of Hertfordshire	<1%
21 Student papers	
University of Wales Institute, Cardiff	<1%
22 Internet	
egusphere.copernicus.org	<1%
23 Publication	
Manish Kumar Chandan, Shrabanti Mandal. "A comprehensive survey on sentime	<1%
24 Internet	
library.nih.go.kr	<1%





25 Internet	
www.arxiv-vanity.com	<1%
26 Internet	
hal.science	<1%
nai.science	
27 Internet	
saucis.sakarya.edu.tr	<1%
28 Internet	
www.ijert.org	<1%
29 Internet	, , ,
www.pure.ed.ac.uk	<1%
30 Internet	
ijsrem.com	<1%
31 Internet	
jebas.org	<1%
	-1//
32 Internet	
arxiv.org	<1%
33 Internet	
core.ac.uk	<1%
34 Internet	
export.arxiv.org	<1%
35 Internet	
icon2021.nits.ac.in	<1%
36 Internet	
ijnrd.org	<1%
37 Internet	
ir.juit.ac.in:8080	<1%
70 V	
su-plus.strathmore.edu	<1%
su-pius.stratiiiilore.edu	<19 ⁴





39 Internet	
www.researchgate.net	<1%
40 Publication	
"AI and Digital Transformation: Opportunities, Challenges, and Emerging Threats	<1%
41 Publication	
A.M. Rajeswari, M. Mahalakshmi, R. Nithyashree, G. Nalini. "Sentiment Analysis fo	<1%
42 Publication	
Narong Pleerux, Attawut Nardkulpat. "Sentiment analysis of restaurant custome	<1%
43 Publication	
Viviane Ito. "Natural Language Processing for Understanding Chronic Illness Pati	<1%
44 Internet	
annals-csis.org	<1%
45 Internet	
api-depositonce.tu-berlin.de	<1%
46 Internet	
jurnal.iaii.or.id	<1%
47 Internet	
repository.biust.ac.bw	<1%
48 Publication	
Muhammad Umer, Saima Sadiq, Hanen karamti, Ala' Abdulmajid Eshmawi, Michel	<1%
49 Publication	
Arvind Dagur, Karan Singh, Pawan Singh Mehra, Dhirendra Kumar Shukla. "Intelli	<1%
50 Publication	
R. N. V. Jagan Mohan, B. H. V. S. Rama Krishnam Raju, V. Chandra Sekhar, T. V. K. P	<1%





2

4

5

6

7

8

9

10

12

14

16

18

19

20

21

22

23

10 24

4 11

9 13

12 15

16 17

Comparing Sentiment Analysis Methods: Flipkart Reviews

3 Abstract

- Sentiment analysis is widely applied to examine opinions, evaluations, attitudes, judgments, and emotions toward a product. In this study, online reviews of the e-commerce portal Flipkart have been analyzed. The dataset contains 189874 rows and 5 columns of product information such as product name, product price, rate, and review and summary of 104 different types of products. Natural Language Processing has been used to carry out the sentiment analysis of online reviews. These techniques are used to conclude the amount of positive and negative reviews received by a product and further identify the opinions of consumers towards the product. In this research, three approaches of sentiment analysis machine learning, unsupervised lexicon-based analysis and large language models have been used. Machine learning classifiers like logistic regression, support vector machine, decision trees, and eXtreme Gradient Boosting(XGBoost); lexicon-based approaches -Valence Aware Dictionary and sEntiment Reasoner(VADER) lexicon and SentiWordNet; and large language models(LLMs)-Generative pre-trained Transformer(GPT) and Bidirectional Encoder Representations from Transformers(BERT) are used. These approaches have been compared based on accuracy, F1-score, recall, precision, and kappa, to determine their effectiveness in sentiment analysis, revealing that machine learning and lexicon-based approaches provide robust performance while the large language models are computationally intensive, timeconsuming, and show comparatively lower accuracy. The identification of context-specific limitations of LLMs in sentiment classification is a significant finding of this work. This comprehensive evaluation of different approaches can help us select the most suitable model for sentiment analysis.
- Keywords: Online reviews, sentiment analysis, machine learning models, lexicon-based
 models, large language models
- 14 27 Introduction
 - Sentiment analysis, often referred to as opinion mining, is a technique of Natural Language
 - Processing (NLP). It helps to identify the human sentiment- positive, negative, or neutral,
- expressed in the textual data [1]. In the era of big data, this technique is vital, as it enables the



36

38

40

42

43

45

52

56

57

59

60

22 61

10 58

1 46

31 understanding of customer opinions shared across various e-commerce sites, social media

32 platforms, and review forums.

33 Customer feedback is important for every business. Customer reviews are vital for 34

understanding purchasing decisions, brand perception, and product development. The reviews

help the business understand the product strengths, flaws and preferences. The feedback of

customers in the form of reviews is instrumental in helping businesses decide their future

directions. Sentiment analysis efficiently processes large volumes of customer feedback in the 37

form of unstructured data [2]. It helps businesses optimize their marketing strategies,

promotion strategies, improve customer support, and product design. This brings a competitive 39

advantage to the organization by predicting market trends and improving customer experience.

Sentiment analysis uses a blend of traditional linguistic techniques and modern machine 41

learning methods [2]. The machine learning models use labelled datasets to train algorithms

that classify text [3] as positive, negative, or neutral. Lexicon-based approaches [4][5] make

20 44 use of a predefined corpus of sentiment-bearing words. Recent advancements in deep learning

and transformer-based architectures[6], such as Bidirectional Encoder Representations from

Transformers (BERT)[7] have significantly enhanced the accuracy and applicability of

47 sentiment analysis.

The comparative analysis of sentiment analysis approaches—machine learning, lexicon-based, 23 48

49 and large language models (LLMs)—reveals varied performance depending on the contexts.

50 Traditional machine learning models are simple, robust and offer foundational accuracy.

51 Transformer models, like Generative pre-trained Transformer (GPT) and BERT, give higher

performance in handling complex linguistics, but they are computationally intensive. Lexicon-

53 based approaches rrequire predefined lexicons, a collection of terms that belong to a particular

54 subject or language. They often lack the precision of machine learning and LLMs.

28 55 Earlier studies focused on comparing three machine learning algorithms—logistic regression,

random forest, and naive bayes for sentiment analysis, without addressing lexicon-based or

large language models, thus limiting its scope to these traditional machine learning approaches

[8]. A comparative study of various machine learning models for sentiment analysis in

financial services, highlighted that neural networks perform well. The importance of sentiment

words and expressions, suggesting that texts rich in these elements yield more reliable

sentiment evaluations was emphasized [9]. Later, machine learning approaches, specifically

Support Vector Machine (SVM) and Long Short-Term Memory (LSTM), were compared with 62



- BERT, and it was highlighted that the BERT model achieved the highest performance metrics 49 63 on a balanced dataset[10]. In one research, LLMs like GPT-4, ChatGPT, and Llama-2-chat 15 64 15 65 were compared against rule-based for sentiment analysis of app reviews. The findings 26 66 suggested LLMs are more promising for structured text and not reviews [11]. Another study evaluated the performance of machine learning and deep learning models in handling 67 customer feedback data, emphasizing the advantages of these models in improving sentiment 68 69 analysis performance[12]. Another approach hasperformed a similar study, showing that SVM outperforms deep learning models like LSTM and BERT[13]. The effectiveness of sentiment 70 71 analysis techniques using Naive Bayes (machine learning), TextBlob (lexicon-based), and LSTM (deep learning) was evaluated and compared on a multi-source dataset from various 72 73 social media platforms [14].
- In the present study, sentiment analysis has been done using machine learning models, 36 74 lexicon-based models and large language models. Machine learning models like Logistic 18 75 Regression(LR), Decision Tree(DT), Random Forest(RF), SVM and XGBoost; lexicon-based 29 76 77 models Valence Aware Dictionary and sEntiment Reasoner(VADER) and SentiWordNet and large language models BERT and GPT2 have been deployed for sentiment analysis of 78 8 79 customer reviews. These models were evaluated based on various performance metrics to 80 determine the most effective approach for extracting insights from customer reviews. The 81 present work employs all 3 above approaches and works on a dataset, which is not balanced. 82 Also, the customer reviews are unstructured text; hence, the text doesn't provide context. This 83 comparative analysis helps in choosing the appropriate according to the specific needs of 84 businesses and data characteristics.

Material & Methodology

- 86 In this research, a dataset of product reviews of the e-commerce portal Flipkart available at the
- 87 Kaggle repository[15] has been analysed using sentiment analysis. This dataset contains
- 88 189874 rows and 5 columns of product information, such as product name, product price,
- rating(on a scale of 1 to 5), review and summary of 104 different types of products. The
- 90 "summary" column has been considered for sentiment analysis.

91 Dataset Preprocessing

- 92 The dataset consisted of unstructured and unformatted data which had to be converted into a
- 93 structured format. Data without a proper framework and structure is difficult to work with and



95

96

97

98

43 99

100

101

102

103

1 04

105

106

107

108

45 09

21 10

111

48 12

39 13

114

115

causes unnecessary errors while running programmes. Structured data draws attention to the characteristics in reviews so that tokenization becomes easier for the algorithm. Tokenization is when text is broken down into smaller units to make it meaningful to the machine without losing the text's initial essence[16]. The data set was cleaned by converting the text into lower case to maintain uniformity and punctuation marks were replaced with spaces to ease tokenization. The missing values,Not a Number (NaN) were replaced with spaces and non-alphanumeric characters were removed. The common stop-words were filtered out using Python's Natural Language Toolkit(NLTK) English stop words list[17]. Tokenization was performed which divides text into smaller pieces, such as phrases or words. The Python function WhitespaceTokenizer()was used for tokenization. This was followed by lemmatization, which returns the words to their original form. The NLTK library was utilized to complete this procedure[18]. The Python function WordNetLemmatizer() was used for lemmatization.

After preprocessing, feature extraction was performed to convert unprocessed raw text data into numerical features suitable for processing while retaining the data from the original dataset [19]. To extract features from the corpus, i.e. the "cleaned" text data was transformed into a sparse matrix. For sentiment analysis, the machine learning model has to know the sentiment score of every unique word in the text data, and its frequency of occurrence. The features and their target values were specified to train the machine learning models. The features are the transformed text data using Term Frequency(TF) and Inverse Document Frequency(IDF) vectorizer [20][21].

Machine Learning Models

The study uses different types of machine learning models, Decision Trees, Random Forest, 4 16 117 Logistic Regression, XGBoost and SVM to predict sentiment labels. The sentiment labels were calculated using TextBlob[22]. The models are trained using TF-IDF vectors and 8 18 evaluated based on precision, accuracy, F1-score, and recall [23]. Their performances are 119 compared.Logistic Regression is used to predict results based on probability [24]. It is a 46 20 121 machine learning model that is quite popular and is simple to understand and use for binary 122 classification tasks such as positive/negative reviews. In our work, Word-Level Logistic Regression has been employed. This technique uses the individual words (unigrams) as 123 features. Words are converted into numerical features using methods like TF-IDF. A decision 13 24 tree is used for tasks concerning classification and is a non-parametric supervised learning 125 37 26 algorithm. It operates by splitting the data into various subsets based on the most important



128

129

6 30

131

132

133

134

13 35

136

137

31 38

139

33 40 141

142

143

144

145

146

147

41 48

4 49

150

151

152

153

5 54

5 55

<u>19</u> 56

157

158

159

160

features. It creates a structure similar to that of a tree made of decisions. Each node in the tree denotes a feature [25]. XGBoost is a popular classification algorithm, an implementation of gradient-boosting decision trees, that is suitable when data training is involved. It is an implementation of gradient boosting machines (GBM) which is known as one of the bestperforming algorithms utilised for supervised learning [26]. It is designed for speed, convenience and performance on large datasets. SVM seeks to find a hyperplane which best divides data into different classes. In a sentiment analysis setting, SVM tries to find a boundary between negative and positive reviews that are as far from each other as possible. SVM works well in high-dimensional spaces and is, therefore, appropriate for text data that can be represented by thousands of features, such as TF-IDF scores of words or n-grams. It is also known for efficiently working with unbalanced datasets and preventing overfitting, which explains its quite good accuracy in your results[27].Random Forest classifier can be described as a collection of tree-structured classifiers. Each tree classifies the input into a class based on its features. The classification with the highest votes is selected by the forest (over all the trees in the forest). The random forest is a classification technique made up of several decision trees[28].

Lexicon-based Approaches

The lexicon-based sentiment analysis or unsupervised sentiment analysis leverages predefined dictionaries of words (lexicons) to assess the sentiment of text by evaluating the polarity of individual words. Polarity scores are numerical scores ranging from 1 (most positive) to -1 (most negative) that indicate the overall sentiment and tone of a phrase or word[29]. A column mentioning the polarity score was added to the dataset, which categorised each review as positive or negative based on the score. In our work, VADER (Valence Aware Dictionary and Sentiment Reasoner) and SentiWordNet libraries have been used. The VADER lexicon [30] is a curated vocabulary of words and phrases for sentiment analysis and their magnitude of the polarity. VADER's sentiment scoring mechanism involves splitting the text into individual tokens like words, phrases, emoticons and retrieving its sentiment intensity score. Sentiment scores are adjusted based on contextual elements such as punctuation, capitalization, degree modifiers, negations, and conjunctions[31]. The adjusted scores are aggregated to compute an overall sentiment score for the text and lastly normalized to produce a compound score ranging from -1 (most negative) to +1 (most positive). The SentiWordNet approach [32] uses SentiWordNet, a lexical resource built on the WordNet database, to assign sentiment scores to words and phrases. It is widely used for analyzing text to determine its sentiment polarity (positive, negative, or neutral). TF-IDF has been used for feature engineering [33]. Sentiment



44 64

165

32 66

169

170

171

172

173

1 74

1 75

1 77

3 83

27 84

185

38 86

188

189

190

161 Analysis with SentiWordNet also tokenizes the input text into words and performs part-of-162

speech (POS) tagging to identify whether a word is verb, noun, adjective, etc.[34]. For each

word, its synsets are identified in WordNet based on the POS tag. The sentiment score from

SentiWordNet is retrieved and the aggregate score is calculated by combining the score of all

words; weighting methods are used to give importance to adjectives and adverbs.

Large Learning Models(LLMs)

LLMs like OpenAI's GPT [35] and Google's BERT [36] are being used for sentiment analysis. 167

These models are capable of understanding context, tone, and even humor, detect sarcasm and 168

even identifying sentiment shifts within a single document as they have been trained on

extensive datasets. In the research, pre-trained LLMs GPT2 [37] and BERT have been

used. In addition to the data preprocessing mentioned above, some additional steps were

performed. The BERT and GPT2 tokenizer functions were employed in this study's tokenizing

procedure. The corpus utilized is "microsoft/DialogRPT-updown" for GPT2 and "bert-base-

uncased" for BERT. The "microsoft/DialogRPT-updown" corpus has 50,257 words and 1,024

token classes, while the "bert-base-uncased" corpus has 30,522 words and 768 token classes.

The quality of the generated tokens may be impacted by the differing tokenization techniques 176

used by BERT and GPT2. This was followed by the encoding function from BERT and GPT2.

Word tokens are transformed into tensor-formatted vectors by this technique. After that, 178

padding is done to make the vector length for all of the data 32. The dataset was divided into 3 79

180 80% and 20% for training and testing respectively. GPT2 model's learning process employed

the AdamW() optimization algorithm, utilizing a 2e-5 learning rate, 1e-8 as eps(Adams 181

182 epsilon) and processing data in groups of 75930 training batches and 18983 testing batches. In

this work, BERT pre-trained model 'bert-base-uncased' was employed, the dataset was

divided into 70%, 15% and 15% for training, validation and testing respectively. The

AdamW() optimizer was used. The number of epochs used for both is 10.

The comparison between the classifiers was done on the basis of various performance metrics

187 namely accuracy, precision, F1-score, confusion matrix, recall and kappa statistics [38].

Results and Discussion

191 The dataset has the "summary" column which has the customer review, it has been analysed

192 for sentiment analysis. The TextBlob python library has been used to calculate the polarity and

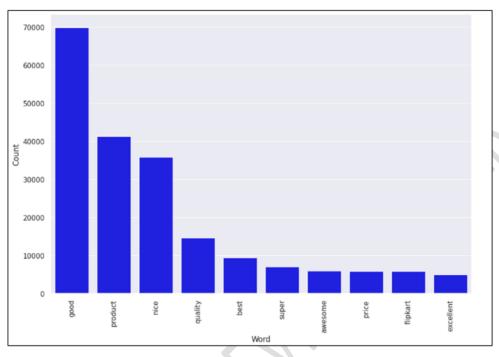
193 subjectivity of the "summary" column, then the column "sentiment label" is calculated as



positive or negative. The neutral sentiment has not been considered in this study. The top 10 words in positive and negative reviews in the dataset are shown in Figure 1 and 2.

196

195



197 198

Figure 1: Words used most frequently in positive reviews

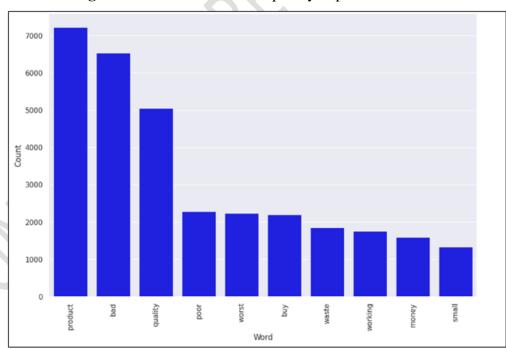


Figure 2: Words used most frequently in negative reviews



202

203

204

The frequently occurring positive and negative words in customer reviews were identified and plotted as word clouds. The word clouds [39] for positive and negative words are shown in Figure 3 and 4.





207

17 08 209

210

211

34 12

213

2 14

215

216

2 17

2 18219

17 06



Figure 3: Word clouds of frequent words in positive reviews



Figure4: Word clouds of frequent words in negative reviews

In this research, traditional machine learning techniques, lexicon-based techniques and large language models were used. The models were run on the dataset and their performance was compared.

Traditional Machine Learning Algorithms

The traditional machine learning algorithms employed TF-IDF for feature extraction. The classification algorithms employed —logistic regression (LR), decision tree (DT), SVM and two ensemble methods—random forest (RF) and XGBoost through TF-IDF features. The function *tfidVectorizer* from Python library *sklearn* was used to generate the TF-IDF feature set. The dataset was split into training and testing sets. The results of sentiment analysis using various traditional classifiers are given in Table 1.

📆 turnitin



221

222

223

224

225

226

227

30 28

229

230

231

232

233234

235

236

Table 1:Performance of machine learning classifiers for sentiment analysis

Machine learning <mark>models</mark>	Avg. Accuracy (%)	Avg. Precision (%)	Avg. Recall (%)	Avg. F1- score (%)	Kappa (%)
LR	96	87	93.5	90	80
DT	98	93.5	94	94	87
RF	98	92.5	97	94.5	89
SVM	97	90.5	93.5	92	84
VCRoost	08	92.5	96.5	04.5	- 80

It can be observed that the accuracy given by classifier is 96% (LR), 98% (DT), 98% (RF), 97% (SVM) and 98% (XGBoost), also the performance of the ensemble methods - RF and XGBoost is slightly higher, XGBoost has higher accuracy (98%) as compared to other classifiers, taking other performance metrics into consideration.

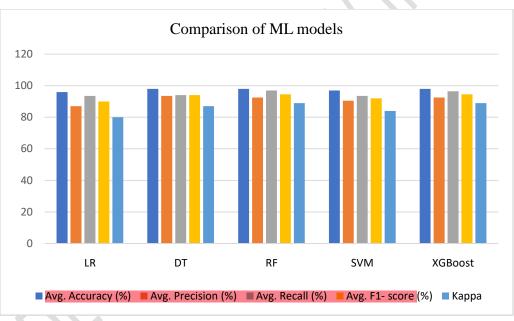


Figure5:Performance of MLmodels

RF and XGBoost showed the best overall performance with high accuracy, precision, recall, and F1-score as is clear from Figure 5. The ensemble models outperform as they combine multiple base models to provide better results and reduce overfitting. These methods can capture the variations and ambiguity in language more effectively than other traditional machine learning models.

Lexicon-based Approaches

For lexicon-based sentiment analysis, VADER and SentiWordNet classifiers were used. TF-IDF has been used for feature extraction. The results of the lexicon-based approach using VADER and SentiWordNet are summarized in Table 2 and 3.



238

Table 2:Performance of classifiers using VADER lexicon-based approach

1 7

Machine learning	Avg. Accuracy	Avg. Precision	Avg. Recall	Avg. F1- score	Kappa (%)
models	(%)	(%)	(%)	(%)	(,,,)
LR	96	94.5	94.5	94.5	89
DT	97	96	95.5	95.5	90
RF	97	95	96	95.5	91
SVM	97	95.5	94.5	95	90
XGBoost	97	96.5	95	95.5	91

239

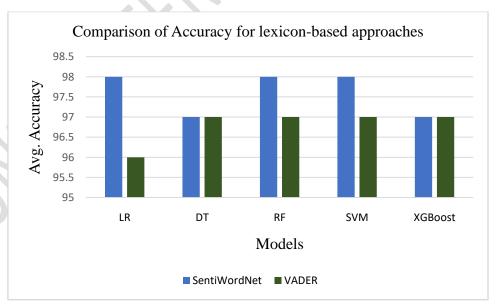
240

Table 3:Performance of classifiers using SentiWordNet lexicon-based approach

Machine learning models	Avg. Accuracy (%)	Avg. Precision (%)	Avg. Recall (%)	Avg. F1- score (%)	Kappa (%)
LR	98	95	92	93.5	87
DT	97	92.5	93	93	85
RF	98	95.5	92	94	88
SVM	98	95	95	95	90
XGBoost	97	95	90	92.5	85

242

243



244245

Figure6: Comparison of accuracy of VADER and SentiWordNet lexicon-based approaches

246 247



250

251

252

253

254

255

256

257

25 58

259

260

8 61

24 62

263

264

265

266

267

268

269

It can be observed from Figure 6 that lexicon-based technique using SentiWordNet gives better accuracy for the classifiers; LR, RF and SVM all of which give 98% accuracy and for DT and XGBoost, the accuracy remains unchanged(97%). SentiWordNet can assign partial sentiment scores to different parts of a sentence, offering more granular analysis. It performs better than VADER because of its lexical richness, contextual sensitivity and adaptability to specific domains. VADER, while good with slang and social media language, often struggles with such context-dependent meanings.

LLMs

For executing the LLM models, the GPU available in Google Colab was deployed. The LLM model BERT used the model 'bert-base-uncased', the batch size was 16, learning rate 1e-5 and number of epochs 10, gave an accuracy of 93%. The training loss and validation loss were 0.274 and 0.282 respectively. The model used for GPT was GPT2, the learning rate was 2e-5, epsilon parameter was 1e-8 and epochs 10, it resulted in accuracy of 95%. The training loss and validation loss were 0.257 and 0.23 respectively. The performance of both these models is summarized in Table 4.

Table 4: Performance of classifiers using LLMs

11	

LLMs	Avg. Accuracy (%)	Avg. Precision (%)	Avg. Recall (%)	Avg. F1- score (%)	Kappa (%)
GPT2	95	88.5	94.5	91	82
BERT	93	86.5	91.5	88	77



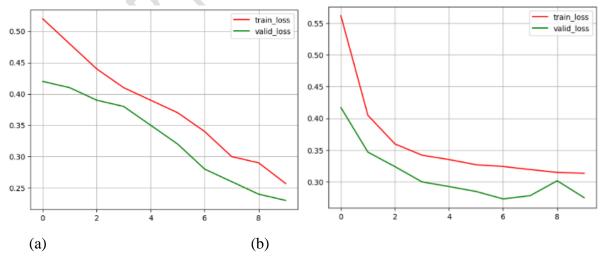


Figure 7: Training and validation loss of GPT and BERT in 10 epochs



It can be observed that the GPT2 model is performing better than the BERT model, although 35 70 271 the execution time taken by GPT2 is longer as compared to the BERT model. The graphs given in Figure 7(a) & (b) indicate that the training and validation loss for both models is 50 72 273 decreasing and hence there is no overfitting. GPT2 was trained using the 274 "microsoft/dialogrpt-updown" corpus, which is tailored for opinion ranking and response 275 modeling, aligning more closely with review sentiment tasks, whereas BERT's "bert-base-276 uncased" lacks domain-specific fine-tuning. It can be concluded that fine-tuned GPT2 model 277 can surprisingly handle classification tasks well at the sentence-level, especially when large 278 datasets are involved.

Conclusions

280

2 81

282

283

16 84

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

This study presented a comparative evaluation of sentiment analysis on thecustomer review dataset using traditional machine learning classifiers, lexicon-based approach and deep learning models. The performance of traditional models, LR, DT, SVM, RF and XGBoost, lexicon-based approaches VADER and SentiWordNet and deep learning models GPT2 and BERT has been compared based on accuracy, F1-score, recall and kappa statistics. The studydemonstrated that the ensemble classifiers such as RF and XGBoost have outperformed other traditional ML models. Lexicon-based models were slightly less robust than ML models, but SentiWordNet demonstrated stronger contextual sensitivity than VADER. On the contrary, large language models like GPT2 and BERT underperformed in comparison to traditional models in this specific application. The lack of domain-specific fine-tuning, high computation requirement and direct binary sentiment classification task, could be the reasons for their reduced performance. Customer reviews often contain complex expressions of sentiment, such as sarcasm, mixed feelings or indirect sentiment. Despite their strength in understanding context, LLMs may still miss these subtleties or generate inaccurate sentiment predictions.

The identification of context-specific limitations of LLMs in sentiment classification is a significant finding of this work. Our findings show that traditional models perform better in domain-specific applications like sentiment analysis of customer review as compared to LLMswhich have shown good performance in several research. The work can be extended by including the neutral sentiments and aspect-based sentiment analysis, which could offer more granular insights.

References

303 [1] B. Liu, Sentiment Analysis and Opinion Mining. Morgan & Claypool Publishers, 2012.



- 304 [2] E. Sullivan, "Understanding from Machine Learning Models," *Br. J. Philos. Sci.*, vol. 305 73, no. 1, pp. 109–133, Mar. 2022, doi: 10.1093/bjps/axz035.
- 306 [3] E. Cambria, B. Schuller, Y. Xia, and C. Havasi, "New Avenues in Opinion Mining and Sentiment Analysis," *IEEE Intell. Syst.*, vol. 28, no. 2, pp. 15–21, Mar. 2013, doi: 10.1109/MIS.2013.30.
- V. Bonta, N. Kumaresh, and N. Janardhan, "A Comprehensive Study on Lexicon Based
 Approaches for Sentiment Analysis," *Asian J. Comput. Sci. Technol.*, vol. 8, no. S2, pp.
 1–6, Mar. 2019, doi: 10.51983/ajcst-2019.8.S2.2037.
- V. Singh, G. Singh, P. Rastogi, and D. Deswal, "Sentiment Analysis Using Lexicon Based Approach," in 2018 Fifth International Conference on Parallel, Distributed and Grid Computing (PDGC), Solan Himachal Pradesh, India: IEEE, Dec. 2018, pp. 13–18. doi: 10.1109/PDGC.2018.8745971.
- 316 [6] A. Gillioz, J. Casas, E. Mugellini, and O. A. Khaled, "Overview of the Transformer-317 based Models for NLP Tasks," presented at the 2020 Federated Conference on 318 Computer Science and Information Systems, Sept. 2020, pp. 179–183. doi: 319 10.15439/2020F20.
- 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*, Minneapolis, Minnesota: Association for Computational Linguistics, 2019, pp. 4171–4186. doi: 10.18653/v1/N19-1423.
- V. M. Abia and E. H. Johnson, "Sentiment Analysis Techniques: A Comparative Study of Logistic Regression, Random Forest, and Naive Bayes on General English and Nigerian Texts," *J. Eng. Res. Rep.*, vol. 26, no. 9, pp. 123–135, Sept. 2024, doi: 10.9734/jerr/2024/v26i91268.
- [9] A. Garg, S. Dhanasekaran, R. Gupta, and J. Logeshwaran, "Exploring Machine 328 329 Learning Models for Intelligent Sentiment Analysis in Financial Services," in 2024 IEEE International Conference on Contemporary Computing and Communications 330 (InC4),Bangalore, India: IEEE, Mar. 2024, doi: 331 pp. 1–6. 332 10.1109/InC460750.2024.10649215.
- 333 [10] A. M. ElMassry, A. Alshamsi, A. F. Abdulhameed, N. Zaki, and A. N. Belkacem, 334 "Machine Learning Approaches for Sentiment Analysis on Balanced and Unbalanced 335 Datasets," in 2024 IEEE 14th International Conference on Control System, Computing 336 and Engineering (ICCSCE), Penang, Malaysia: IEEE, Aug. 2024, pp. 18–23. doi: 337 10.1109/ICCSCE61582.2024.10695972.
- 138 [11] F. A. Shah, A. Sabir, R. Sharma, and D. Pfahl, "How Effectively Do LLMs Extract Feature-Sentiment Pairs from App Reviews?," 2024, arXiv. doi: 10.48550/ARXIV.2409.07162.
- [12] T. Ananth Kumar, J. Zaafira, P. Kanimozhi, R. Rajmohan, A. Christo, and A. A. Sunday,
 "Machine Learning and Sentiment Analysis: Analysing Customer Feedback," in
 Advances in Marketing, Customer Relationship Management, and E-Services, R.
 Masengu, O. T. Chiwaridzo, M. Dube, and B. Ruzive, Eds., IGI Global, 2024, pp. 245–262. doi: 10.4018/979-8-3693-2165-2.ch014.
- 346 [13] N. Jing Xiang, K. M. Lim, C. P. Lee, Q. Z. Lim, E. K. H. Ooi, and N. K. N. Loh, "Sentiment Analysis Using Learning-based Approaches: A Comparative Study," in 2023 347 348 11th International Conference on Information and Communication Technology 349 (ICoICT), Melaka, Malaysia: IEEE, 2023, 469-474. doi: Aug. pp. 10.1109/ICoICT58202.2023.10262604. 350
- 351 [14] K. Kapur and R. Harikrishnan, "Comparative Study of Sentiment Analysis for Multi-352 Sourced Social Media Platforms," 2022, *arXiv*. doi: 10.48550/ARXIV.2212.04688.



- [15] H. Sharma and V. Kakran, "Sentiment Analysis: Analyzing Flipkart Product Reviews using NLP and Machine Learning," in 2024 International Conference on Computing, Sciences and Communications (ICCSC), Ghaziabad, India: IEEE, Oct. 2024, pp. 1–7. doi: 10.1109/ICCSC62048.2024.10830421.
- Tokenization Approaches in Sentiment Classification," *IEEE Access*, vol. 11, pp. 134951–134968, 2023, doi: 10.1109/ACCESS.2023.3337354.
- [17] D. Yogish, T. N. Manjunath, and R. S. Hegadi, "Review on Natural Language Processing Trends and Techniques Using NLTK," in *Recent Trends in Image Processing and Pattern Recognition*, vol. 1037, K. C. Santosh and R. S. Hegadi, Eds., in Communications in Computer and Information Science, vol. 1037., Singapore: Springer Singapore, 2019, pp. 589–606. doi: 10.1007/978-981-13-9187-3_53.
- [18] N. Hardeniya, J. Perkins, D. Chopra, N. Joshi, and i. Mathur, *Natural Language Processing: Python and NLTK*. Packt Publishing, 2016.
- 19] R. Ahuja, A. Chug, S. Kohli, S. Gupta, and P. Ahuja, "The Impact of Features Extraction on the Sentiment Analysis," *Procedia Comput. Sci.*, vol. 152, pp. 341–348, 2019, doi: 10.1016/j.procs.2019.05.008.
- [20] A. Aizawa, "An information-theoretic perspective of tf-idf measures," *Inf. Process.* Manag., vol. 39, no. 1, pp. 45–65, Jan. 2003, doi: 10.1016/S0306-4573(02)00021-3.
- 372 [21] S. Qaiser and R. Ali, "Text Mining: Use of TF-IDF to Examine the Relevance of Words to Documents," *Int. J. Comput. Appl.*, vol. 181, no. 1, pp. 25–29, July 2018, doi: 10.5120/ijca2018917395.
- 375 [22] S. Loria, "Textblob Documentation. Release 0.15," 2018.
- [23] N. W. S. Wardhani, M. Y. Rochayani, A. Iriany, A. D. Sulistyono, and P. Lestantyo,
 "Cross-validation Metrics for Evaluating Classification Performance on Imbalanced
 Data," in 2019 International Conference on Computer, Control, Informatics and its
 Applications (IC3INA), Tangerang, Indonesia: IEEE, Oct. 2019, pp. 14–18. doi:
 10.1109/IC3INA48034.2019.8949568.
- [24] A. Poornima and K. S. Priya, "A Comparative Sentiment Analysis Of Sentence
 Embedding Using Machine Learning Techniques," in 2020 6th International Conference
 on Advanced Computing and Communication Systems (ICACCS), Coimbatore, India:
 IEEE, Mar. 2020, pp. 493–496. doi: 10.1109/ICACCS48705.2020.9074312.
- [25] L. Rokach and O. Maimon, "Decision Trees," in *Data Mining and Knowledge Discovery Handbook*, O. Maimon and L. Rokach, Eds., New York: Springer-Verlag, 2005, pp. 165–192. doi: 10.1007/0-387-25465-X_9.
- T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," *Knowl. Discov. Data Min.*, 2016.
- [27] T. Evgeniou and M. Pontil, "Support Vector Machines: Theory and Applications," in Machine Learning and Its Applications, vol. 2049, G. Paliouras, V. Karkaletsis, and C. D. Spyropoulos, Eds., in Lecture Notes in Computer Science, vol. 2049. , Berlin, Heidelberg: Springer Berlin Heidelberg, 2001, pp. 249–257. doi: 10.1007/3-540-44673-7_12.
- 395 [28] L. Breiman, "Random Forests. Machine Learning," *Mach. Learn.*, vol. 45, no. 1, pp. 5–396 32, 2001, doi: 10.1023/A:1010933404324.
- [29] D. H. Abd, A. R. Abbas, and A. T. Sadiq, "Analyzing sentiment system to specify polarity by lexicon-based," *Bull. Electr. Eng. Inform.*, vol. 10, no. 1, pp. 283–289, Feb. 2021, doi: 10.11591/eei.v10i1.2471.
- 400 [30] D. C. Youvan, "Understanding Sentiment Analysis with VADER: A Comprehensive Overview and Application," 2024, doi: 10.13140/RG.2.2.33567.98726.



406 407

408

- 402 [31] C. Hutto and E. Gilbert, "VADER: A Parsimonious Rule-Based Model for Sentiment 403 Analysis of Social Media Text," *Proc. Int. AAAI Conf. Web Soc. Media*, vol. 8, no. 1, pp. 404 216–225, May 2014, doi: 10.1609/icwsm.v8i1.14550.
 - [32] A. Cernian, V. Sgarciu, and B. Martin, "Sentiment analysis from product reviews using SentiWordNet as lexical resource," in 2015 7th International Conference on Electronics, Computers and Artificial Intelligence (ECAI), Bucharest, Romania: IEEE, June 2015, p. WE-15-WE-18. doi: 10.1109/ECAI.2015.7301224.
- 409 [33] S. Akuma, T. Lubem, and I. T. Adom, "Comparing Bag of Words and TF-IDF with different models for hate speech detection from live tweets," *Int. J. Inf. Technol.*, vol. 14, no. 7, pp. 3629–3635, Dec. 2022, doi: 10.1007/s41870-022-01096-4.
- 412 [34] S. Baccianella, A. Esuli, and F. Sebastiani, "SentiWordNet 3.0: An Enhanced Lexical 413 Resource for Sentiment Analysis and Opinion Mining," presented at the Proceedings of 414 the International Conference on Language Resources and Evaluation, LREC 2010, 17-415 23, Valletta, Malta, May 2010.
- 416 [35] K. I. Roumeliotis and N. D. Tselikas, "ChatGPT and Open-AI Models: A Preliminary Review," *Future Internet*, vol. 15, no. 6, p. 192, May 2023, doi: 10.3390/fi15060192.
- 418 [36] R. Sudharsan, Getting Started with Google BERT: Build and Train State-Of-the-art 419 Natural Language Processing Models Using BERT. Packt Publishing, 2021.
- 420 [37] M. Hanna, O. Liu, and A. Variengien, "How does GPT-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model," 2023, *arXiv*. doi: 10.48550/ARXIV.2305.00586.
- 423 [38] M. Aldwairi and A. Alwahedi, "Detecting Fake News in Social Media Networks," 424 *Procedia Comput. Sci.*, vol. 141, pp. 215–222, 2018, doi: 10.1016/j.procs.2018.10.171.
- [39] F. Heimerl, S. Lohmann, S. Lange, and T. Ertl, "Word Cloud Explorer: Text Analytics 425 426 Based on Word Clouds," in 2014 47th Hawaii International Conference on System 427 Sciences. Waikoloa, HI: IEEE, Jan. 2014, 1833-1842. doi: pp. 10.1109/HICSS.2014.231. 428

429

430

431

432

433

434 435

436 437

