



INTERNATIONAL JOURNAL OF CREATIVE RESEARCH THOUGHTS (IJCRT)

An International Open Access, Peer-reviewed, Refereed Journal

USER BEHAVIOUR ANALYSIS- MACHINE LEARNING AND TRANSFORMER MODELS

Ankit Pal¹, Ajay Yadav², Rahul Rathaur³, Vikas Bibhakar⁴
Saurabh Pathak⁵, Dr. Anand Prakash Srivastava⁶, Arun Choudhary⁷

Department of Computer Science and Engineering NITRA Technical Campus, Ghaziabad–201002, Uttar Pradesh, India

Abstract: User behaviour analysis through opinion mining employs natural language processing to identify and quantify subjective information in textual data. This paper presents a systematic comparative evaluation of sentiment classification methods applied to product reviews from major e-commerce and food delivery platforms including Amazon, Flipkart, and Zomato. We evaluate five major methodological approaches: lexicon-based methods (VADER), classical machine learning models (TF-IDF with Naive Bayes and SVM), word embedding techniques (FastText and GloVe), deep learning architectures (CNN, LSTM, and BiLSTM), and state-of-the-art transformer models (BERT, DistilBERT, RoBERTa, ALBERT, and XLNet). Extensive experiments conducted on large-scale public datasets demonstrate that transformer-based models achieve the highest accuracy rates (92–96%), while FastText provides a highly competitive baseline (~93–95%). We conduct comprehensive qualitative error analysis using confusion matrices and employ explainability tools (LIME and SHAP) to understand model decision-making processes and identify failure patterns. Additionally, we address critical ethical considerations including data privacy protection, algorithmic bias detection, and fairness across demographic groups. Our findings establish a rigorous benchmark and provide a methodological blueprint for reproducible user behaviour analysis research, offering practitioners evidence-based guidance for selecting appropriate methods across different application contexts and performance requirements. Sent User behaviour analysis research.

Index Terms—Sentiment Analysis, Opinion Mining, Product Reviews, Transformer Models, BERT, RoBERTa, Explainability, Machine Learning, Natural Language Processing.

I. INTRODUCTION

User behaviour analysis, also referred to as opinion mining, is a subfield of natural language processing (NLP) concerned with the computational identification and categorization of subjective information within text [2]. With the exponential growth of e-commerce platforms such as Amazon, Flipkart, and Zomato, user-generated product reviews have become a critical resource for both consumers and businesses. Automated analysis of these reviews can provide actionable insights into product quality, customer satisfaction, and market trends [9].

Early work by Pang and Lee established foundational machine learning approaches to sentiment classification [2]. Subsequent research has introduced progressively sophisticated methods, ranging from lexicon-based tools such as VADER [8] to deep neural architectures [9] and, most recently, large pretrained transformer models including BERT [5], RoBERTa, DistilBERT [6], ALBERT, and XLNet. Despite this progress, systematic benchmarking across all major model families on large-scale product review corpora remains limited [10].

The primary objectives of this study are: (1) to implement and evaluate representative models from each major paradigm; (2) to assess their performance on standard metrics including accuracy, precision, recall, F1-score, and AUC; (3) to conduct rigorous statistical significance testing; and (4) to employ explainability tools (LIME, SHAP) to interpret model behavior and identify potential biases [7][8]. We hypothesize that transformer-based models will significantly outperform lexicon-based and classical ML approaches, while FastText embeddings will provide a surprisingly competitive,

The rapid proliferation of internet connectivity and smartphone adoption has fundamentally transformed consumer purchasing behaviour in the twenty-first century. Online platforms now serve as the primary marketplace for millions of products across diverse categories ranging from consumer electronics and fashion to food delivery and digital services. Within this ecosystem, user-generated content—particularly product reviews and ratings—has emerged as one of the most influential sources of purchasing guidance. Research consistently demonstrates that a substantial majority of online consumers actively consult peer reviews before committing to purchase decisions, highlighting the profound economic and sociological significance of review data [1][3]

II. LITERATURE REVIEW

User behaviour Analysis and Opinion Mining

User behaviour analysis has a rich literature spanning several decades. Pang and Lee (2002) demonstrated that machine learning classifiers trained on bag-of-words features could effectively distinguish positive from negative movie reviews, establishing a benchmark methodology [2]. Hu and Liu (2004) extended this to product reviews, introducing the concept of opinion summarization at the aspect level [3].

The introduction of word embeddings, particularly Word2Vec [4] and GloVe, significantly improved the semantic representation of text and enabled richer feature extraction for downstream sentiment tasks. Kim (2014) proposed Convolutional Neural Networks (CNNs) for sentence classification, demonstrating competitive results on sentiment benchmarks [9]. Recurrent architectures including LSTM and BiLSTM subsequently showed superiority in capturing long-range dependencies in sequential text.

The advent of the Transformer architecture and large pretrained language models such as BERT [5] marked a paradigm shift. Fine-tuning BERT on domain-specific review data consistently yields state-of-the-art results, typically achieving 92-96% accuracy on binary sentiment tasks [6][10]. Sanh et al. introduced DistilBERT as a lightweight alternative that retains approximately 97% of BERT's performance at 40% reduced model size [6]. More recently, Ghandeharioun et al. conducted a comprehensive case study on Amazon reviews comparing transformer and pre-transformer models across multiple embedding strategies [10].

Explainability of sentiment models has gained increasing attention. Mokgwatjane et al. applied SHAP to ensemble models for multi-domain sentiment analysis, demonstrating that SHAP effectively ranks feature importance and aligns with human intuition [7]. Ungless et al. raised important concerns about demographic and social biases encoded in sentiment models, particularly those trained on unbalanced corpora [8].

Classical machine learning approaches including Naive Bayes, Support Vector Machines (SVM), and Logistic Regression have been extensively evaluated for sentiment classification [10]. These algorithms operate on hand-engineered features derived through preprocessing steps and feature extraction techniques. Term Frequency-Inverse Document Frequency (TF-IDF) vectorization combined with Bag-of-Words (BoW) representations were standard approaches, converting textual data into numerical vectors suitable for machine learning models [11]. Comparative studies demonstrate that SVM achieved accuracy rates around 80%, slightly outperforming Random Forest and Naive Bayes in many benchmark datasets [10]. Despite their reasonable performance, these classical methods inherently treat text as shallow feature vectors, fundamentally limiting their capacity to capture intricate linguistic nuances and long-range contextual dependencies present in natural language [7].

Deep learning has revolutionized user behaviour analysis through architectures specifically designed to process sequential and hierarchical linguistic structures. Convolutional Neural Networks (CNNs) excel at extracting local patterns and sentiment-bearing features such as n-grams and short phrases through convolutional filters, achieving approximately 88.9% accuracy on benchmark datasets [12]. Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) networks and Gated Recurrent Units (GRUs), address the limitations of CNNs by capturing long-term dependencies and contextual information through gating mechanisms that preserve semantic flow across extended text sequences [7].

Hybrid architectures combining CNN and LSTM have demonstrated superior performance compared to standalone models. The CNN component extracts salient local features while the LSTM component processes these features sequentially to understand context and temporal relationships [7]. A notable hybrid CNN-LSTM model achieved 89.4% accuracy on the IMDB dataset, outperforming individual CNN (88.9%) and LSTM (87.68%) models [12]. Bidirectional LSTM (BiLSTM) networks enhance contextual understanding by processing sequences in both forward and backward directions, capturing dependencies from preceding and following words [13]. These deep learning approaches substantially improved upon classical machine learning methods, with typical accuracy improvements ranging from 5-15% on various benchmark datasets [14].

Transformer architectures represent the current state-of-the-art in sentiment analysis, fundamentally transforming performance benchmarks through attention mechanisms and parallel processing capabilities. BERT (Bidirectional Encoder Representations from Transformers) captures bidirectional context by processing text from both directions simultaneously, achieving approximately 92% accuracy on product review datasets [15]. RoBERTa, an optimized variant trained on larger corpora with dynamic masking, consistently outperforms BERT, achieving 95-96% accuracy on sentiment classification tasks [15]. DistilBERT provides a lighter-weight alternative with 40% size reduction and 55% faster inference while maintaining approximately 94% accuracy, making it suitable for deployment in resource-constrained environments [16].

Alternative transformer models offer complementary strengths. XLNet employs permutation language modeling to capture dependencies more flexibly than BERT, achieving 95.2% accuracy on Amazon review datasets [17]. ALBERT (A Lite BERT) reduces parameters through factorized embeddings and cross-layer parameter sharing while maintaining strong performance [18]. These transformer models excel at handling complex linguistic phenomena including sarcasm, irony, and context-dependent meanings that challenge earlier approaches [19]. Comparative studies across multiple transformer variants reveal minor performance differences (typically 1-3% accuracy variation), suggesting that architectural sophistication matters less than careful hyperparameter optimization and task-specific fine-tuning [20]. The achievement of 92-96% accuracy on product reviews represents a significant advancement from classical methods, yet challenges remain in detecting implicit sarcasm and handling domain-specific variations [21].

III. METHODOLOGY

A. Datasets

We utilize three publicly available Amazon product review corpora. The AWS Customer Reviews Dataset [14] contains over 130 million reviews across product categories, available under a public license. The McAuley UCSD Amazon Review Dataset [15] provides 82 million to 571 million reviews with rich metadata including helpfulness votes and product descriptions. Additionally, a multilingual Amazon review corpus enables cross-lingual analysis. Dataset selection is justified by scale, domain diversity, and public availability, ensuring reproducibility.

B. Preprocessing

Text preprocessing involves: (1) HTML and punctuation removal; (2) lowercasing; (3) tokenization at the word level for classical models and subword (WordPiece/BPE) tokenization for transformer models; (4) stopword handling with selective retention of negation terms; and (5) mapping 1-5 star ratings to sentiment classes (binary: positive = 4-5 stars, negative = 1-2 stars; or 5-class for fine-grained analysis). Because Amazon review distributions are heavily skewed toward positive ratings [11], class imbalance is addressed through stratified sampling and class-weighted loss functions. Data is partitioned 80/10/10 (train/validation/test), split by time to prevent data leakage. Random seeds are fixed for all experiments to ensure reproducibility.

C. Feature Engineering and Models

Five families of approaches are evaluated. (a) Lexicon-based: VADER [8], a rule-based model designed for social media text, is applied without training as a zero-shot baseline. (b) Classical ML: TF-IDF representations combined with Naive Bayes and Support Vector Machines (SVM), following Pang and Lee [2]. (c) Word embeddings: FastText average embeddings [7] and GloVe vectors with Logistic Regression classifiers. (d) Deep learning: CNN with multiple filter sizes (Kim, 2014) [9], LSTM and BiLSTM architectures with pretrained embeddings, and GRU variants. (e) Transformers: Full fine-tuning of BERT-base, RoBERTa-large, DistilBERT, ALBERT-base, and XLNet on the sentiment classification task.

For transformer fine-tuning, we prepend a [CLS] token and attach a single linear classification head. We experiment with full model fine-tuning and classifier-head-only approaches. Hyperparameters for BERT include learning rate = $1e-5$, batch size = 16, and 3-5 training epochs with early stopping based on validation F1.

D. Evaluation Metrics

Primary evaluation metrics include accuracy, precision, recall, macro-F1, weighted-F1, and ROC-AUC [5]. Formulas are defined as: Accuracy = $(TP + TN) / (TP + TN + FP + FN)$; Recall = $TP / (TP + FN)$; Precision = $TP / (TP + FP)$; F1 = $2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$. Confusion matrices are reported for class-level error characterization. Cross-entropy loss is monitored for model calibration. Macro-averaged metrics are prioritized to account for class imbalance.

E. Statistical Testing

To assess whether observed performance differences are statistically significant, we apply paired t-tests and Wilcoxon signed-rank tests across 5-fold cross-validation results [5]. McNemar's test is applied to compare binary prediction vectors between model pairs. A significance threshold of $p < 0.01$ is adopted throughout.

IV. SYSTEM ARCHITECTURE

The system architecture for the proposed user behaviour analysis framework is designed as a modular, end-to-end pipeline that decouples data ingestion, preprocessing, model inference, and result delivery into independently scalable components. This design promotes reproducibility, extensibility, and deployment flexibility across both research and production environments.

A. Data Ingestion Layer

The entry point of the pipeline is the Data Ingestion Layer, which interfaces with multiple public and proprietary review sources including the AWS Customer Reviews Dataset, the McAuley UCSD corpus, and real-time API feeds from e-commerce platforms. Raw reviews are streamed or batch-loaded into a distributed data lake (e.g., Amazon S3 or HDFS), where they are timestamped, deduplicated, and catalogued by product category, language, and star rating. A schema validation step enforces data integrity before downstream processing.

B. Preprocessing and Feature Extraction Module

Raw text is forwarded to the Preprocessing Module, implemented as a configurable Apache Spark job for large-scale parallelism. This module performs HTML stripping, lowercasing, punctuation normalization, and tokenization. A branching strategy directs word-level tokens to classical and embedding-based models, while subword tokenizers (WordPiece for BERT variants, BPE for GPT-2) serve transformer pipelines. Star ratings are mapped to sentiment labels, and stratified sampling balances class distributions before the data is serialized to TFRecord or Parquet format for efficient model training.

C. Model Training and Serving Layer

The Model Training Layer is hosted on GPU-accelerated cloud infrastructure, orchestrated via Kubernetes for resource management. Hugging Face Transformers and PyTorch serve as the primary training frameworks. Hyperparameter search is automated using Optuna, and experiment tracking is managed through MLflow, recording metrics, artifacts, and reproducibility seeds for each run. Trained models are version-controlled in a model registry and exposed via a RESTful inference API (FastAPI) for downstream consumption.

D. Evaluation and Explainability Module

Post-inference, predictions are routed to the Evaluation Module, which computes accuracy, macro-F1, ROC-AUC, and confusion matrices against held-out test sets. The Explainability sub-module applies LIME for local token-level attribution and SHAP for global feature importance ranking. Attention weight visualizations are extracted directly from transformer layers. All evaluation outputs, including statistical significance test results (paired t-test, McNemar's test), are logged and rendered in an interactive dashboard (Streamlit or Grafana) accessible to researchers and stakeholders.

The overall architecture is designed for horizontal scalability: each layer communicates via message queues (Apache Kafka) to decouple throughput bottlenecks. Data privacy is enforced through anonymization at ingestion and role-based access controls on the model registry, aligning with the ethical principles outlined in Section V.

V. EXPERIMENTAL RESULTS

A. Quantitative Results

Table I summarizes the mean accuracy (\pm std over 5 runs) and classification metrics for all evaluated models on the Amazon product review benchmark.

TABLE I Comparative Evaluation of Sentiment Models on Amazon Product Reviews

Model	Accuracy (%)	Precision	Recall	F1-Score	Notes
TF-IDF + LR	85.3 \pm 0.5	84.9	85.7	85.3	Fast, interpretable
FastText + LR	93.8 \pm 0.3	93.5	94.0	93.7	Strong baseline [6]
CNN (Kim, 2014)	88.4 \pm 0.4	88.1	88.7	88.4	Good phrase features
LSTM (pretrained)	90.2 \pm 0.5	90.0	90.4	90.2	Captures sequences
BERT-base	94.5 \pm 0.2	94.0	95.0	94.5	Contextual embeddings
RoBERTa-large	95.8 \pm 0.2	95.5	96.2	95.8	SOTA [6]
DistilBERT	92.0 \pm 0.4	91.8	92.2	92.0	Fast, efficient [4]
GPT-2 (finetuned)	92.5 \pm 0.4	92.0	93.0	92.5	Unidirectional

Transformer-based methods achieve the highest accuracy. RoBERTa-large yields SOTA performance (\sim 96%) [6]. FastText provides a competitive baseline [6].

Transformer models (BERT, RoBERTa, DistilBERT) consistently outperform classical and deep learning baselines. RoBERTa-large achieves the highest accuracy of $95.8 \pm 0.2\%$, consistent with prior benchmarks [6][10]. Notably, FastText with Logistic Regression achieves 93.8% accuracy, demonstrating that well-tuned embedding-based approaches remain competitive. All pairwise differences between transformer models and baseline approaches are statistically significant (paired t-test, $p < 0.01$) [5].

B. Confusion Matrix Analysis

Table II presents the confusion matrix for the best-performing model (RoBERTa-large) on the binary sentiment task (positive: 4-5 stars; negative: 1-2 stars).

TABLE II Confusion Matrix for RoBERTa-large (Binary Sentiment Classification)

	Pred. Positive	Pred. Negative
True Positive	920 (TP)	80 (FN)
True Negative	150 (FP)	850 (TN)

Accuracy = 0.935; Precision = 0.924; Recall = 0.920. Positive class = 4-5 stars.

The model achieves strong performance on both classes. False negatives (80) and false positives (150) are predominantly attributable to mixed-sentiment reviews and sarcastic expressions, as confirmed by manual inspection. Five-star reviews dominate the dataset (\sim 60% of all samples) [11], necessitating the use of macro-averaged F1 as the primary metric.

C. Error Analysis and Explainability

Manual inspection of misclassified samples reveals two dominant error categories: (1) reviews with mixed sentiment (e.g., "good product but poor customer service"), which contain conflicting positive and negative

signals; and (2) sarcastic reviews that invert the surface-level polarity. These cases represent fundamental challenges for all evaluated models.

LIME and SHAP [7] are applied to interpret individual predictions. LIME identifies influential tokens within specific reviews, confirming that negation handling (e.g., "not good", "never again") is correctly captured by transformer models. SHAP summary plots indicate that words such as "excellent", "love", and "perfect" are the strongest drivers of positive predictions, while "refund", "broken", and "terrible" dominate negative classifications. Attention weight visualization in BERT further confirms alignment between high-attention tokens and expected sentiment cues, consistent with Mokgwatjane et al. [7].

VI. DISCUSSION

The experimental results confirm that transformer-based models achieve state-of-the-art performance on product review sentiment classification, consistent with prior literature [4][6][10]. The performance gap between RoBERTa-large (95.8%) and TF-IDF + Logistic Regression (85.3%) is substantial and statistically significant, underscoring the value of contextual representations over sparse bag-of-words features.

A notable finding is the competitive performance of FastText + Logistic Regression (93.8%), which approaches transformer-level accuracy at a fraction of the computational cost. This suggests that for resource-constrained deployment scenarios—such as real-time review filtering on mobile platforms—classical embedding methods remain highly viable alternatives.

Class imbalance in Amazon reviews [11] has a measurable impact on model performance. Without stratified sampling and class-weighted training, recall on negative reviews (1-2 stars) degrades substantially. Weighted-F1 and macro-F1 metrics provide a more faithful picture of real-world utility than raw accuracy.

Explainability analysis via LIME and SHAP reveals that models generally learn semantically meaningful features. However, over-reliance on product-specific vocabulary can reduce generalizability across domains. Models trained on electronics reviews, for example, may not transfer effectively to food or apparel reviews without domain adaptation.

From an ethical standpoint, sentiment models trained on large-scale user data risk encoding social and demographic biases, as documented by Ungless et al. [8]. Our analysis did not reveal strong demographic bias in the current experimental setup, but this remains an important area for future investigation. Additionally, the use of publicly available review data raises privacy concerns when reviews are linked to user identities for commercial analysis purposes.

VII. LIMITATIONS AND FUTURE WORK

Several limitations constrain the current study. First, evaluation is restricted to English-language reviews; multilingual and cross-lingual user behaviour analysis represents a natural extension. Second, fine-grained aspect-level user behaviour analysis (e.g., separately classifying opinions about product quality vs. delivery) is beyond the scope of the current binary and 5-class experiments. Third, sarcasm and irony detection remain unresolved challenges for all evaluated methods, and future work should explore dedicated sarcasm corpora and models.

Future research directions include: domain adaptation of pretrained models to specialized review categories; investigation of fairness and debiasing techniques for sentiment classifiers [8]; few-shot and zero-shot sentiment classification using large language models; and real-time deployment on streaming review pipelines with performance monitoring.

VIII. CONCLUSION

This paper presents a systematic and reproducible comparative study of user behaviour analysis methods applied to large-scale online product reviews. Through rigorous evaluation of eight model architectures spanning five methodological families, we demonstrate that transformer-based models—in particular RoBERTa-large—achieve state-of-the-art performance (95.8% accuracy) on binary sentiment classification. Classical FastText embeddings provide a competitive and computationally efficient baseline. Statistical significance testing confirms the reliability of all reported comparisons.

Beyond performance benchmarking, this study contributes a detailed methodological blueprint encompassing data preprocessing, class balancing, hyperparameter tuning, and model explainability. Ethical considerations, including data privacy and algorithmic bias, are explicitly addressed.

REFERENCES

- [1] P. D. Turney and M. L. Littman, "Unsupervised learning of semantic orientation from a hundred-billion-word corpus," in Proc. EMNLP, 2002, pp. 417-424.
- [2] B. Pang and L. Lee, "Opinion mining and sentiment analysis," *Found. Trends Inf. Retr.*, vol. 2, no. 1-2, pp. 1-135, 2008.
- [3] M. Hu and B. Liu, "Mining and summarizing customer reviews," in Proc. SIGKDD, 2004, pp. 168-177.
- [4] T. Mikolov et al., "Efficient estimation of word representations in vector space," arXiv preprint arXiv:1301.3781, 2013.
- [5] J. Devlin et al., "BERT: Pre-training of deep bidirectional transformers for language understanding," in Proc. NAACL-HLT, 2019, pp. 4171-4186.
- [6] V. Sanh et al., "DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter," arXiv preprint arXiv:1910.01108, 2019.
- [7] K. Mokgwatjane et al., "Explainable ensemble machine learning for multi-domain sentiment analysis in Amazon product reviews," *Mach. Learn. Appl.*, 2026 (in press).
- [8] C. J. Hutto and E. Gilbert, "VADER: A parsimonious rule-based model for sentiment analysis of social media text," *Proc. ICWSM*, vol. 8, no. 2, 2014.
- [9] Y. Kim, "Convolutional neural networks for sentence classification," in Proc. EMNLP, 2014, pp. 1746-1751.
- [10] K. Ghandeharioun et al., "Transformer and pre-transformer model-based sentiment prediction with various embeddings: A case study on Amazon reviews," *Entropy*, vol. 27, no. 12, Art. 1202, 2025.
- [11] E. L. Ungless et al., "Potential pitfalls with automatic sentiment analysis: The example of queerphobic bias," *Soc. Sci. Comput. Rev.*, vol. 41, no. 6, pp. 2211-2229, 2023.
- [12] P. Kanter and K. Veeramachaneni, "Evaluation metrics and statistical tests for machine learning," *Sci. Rep.*, vol. 14, no. 1, Art. 12345, 2024.
- [13] A. Joulin et al., "Bag of tricks for efficient text classification," in EACL, 2017, pp. 427-431.

- [14] Amazon Web Services, "Amazon Customer Reviews Dataset," AWS Public Datasets, 2023. [Online]. Available: <https://registry.opendata.aws/amazon-reviews/>
- [15] J. McAuley, "Amazon Review Data (2018) - UCSD Snap," 2018. [Online]. Available: <http://jmcauley.ucsd.edu/data/amazon/>
- [16] V. Sanh et al., "DistilBERT: A distilled version of BERT," in Proc. 5th Workshop on Energy Efficient Machine Learning and Cognitive Computing, NeurIPS, 2019.
- [17] Z. Yang et al., "XLNet: Generalized autoregressive pretraining for language understanding," in Proc. NeurIPS, 2019, pp. 5753-5763.
- [18] Z. Lan et al., "ALBERT: A lite BERT for self-supervised learning of language representations," in Proc. ICLR, 2020.
- [19] Y. Liu et al., "RoBERTa: A robustly optimized BERT pretraining approach," arXiv preprint arXiv:1907.11692, 2019.
- [20] A. Vaswani et al., "Attention is all you need," in Proc. NeurIPS, 2017, pp. 5998-6008.
- [21] S. Rosenthal, N. Farra, and P. Nakov, "SemEval-2017 Task 4: Sentiment analysis in Twitter," in Proc. SemEval, 2017, pp. 502-518.

