

UDC 004.912

AUTOMATED DETECTION AND ANALYSIS OF CUSTOMER COMPLAINTS IN ONLINE STORE REVIEWS USING NLP

Aziyeva A.

Master student, Computing and Data Science Department, Applied Data Analytics, Astana IT University, Astana, Kazakhstan

This study applies Natural Language Processing (NLP) techniques to detect and analyze customer complaints in online store reviews using the Amazon US Customer Reviews Dataset. The methodology includes data preprocessing, feature extraction, and classification using Naïve Bayes, BERT, and a hybrid model combining rule-based filtering with deep learning. Results show that BERT outperforms traditional models. Analysis reveals that verified purchase reviews are more likely to contain genuine complaints. Despite improvements, challenges remain in classifying neutral reviews and handling domain-specific terminology. Future work will focus on multilingual expansion, real-time detection, and Explainable AI integration to enhance model transparency and business applicability.

Keywords: Natural language processing, sentiment analysis, complaint detection, bert, hybrid model.

Introduction. In the era of digital commerce, online reviews are essential in influencing purchasing choices and guiding business strategies. Customers often share their experiences—both praises and complaints—through reviews on e-commerce platforms. Accurately identifying and analyzing these complaints is crucial for businesses to enhance their services, boost customer satisfaction, and stay competitive. However, manual analysis of large-scale customer reviews is time-consuming and inefficient. Natural Language Processing (NLP) techniques offer a scalable and automated solution for extracting valuable insights from textual feedback.

This study proposes an automated approach to detecting and analyzing customer complaints in online store reviews using NLP. Using sentiment analysis, aspect-based sentiment analysis (ABSA), and machine learning models, we strive to detect complaint-related reviews, categorize them appropriately, and pinpoint key factors that drive negative customer experiences. Our approach builds upon existing research in sentiment analysis and complaint detection, extending it with modern deep learning techniques and domain-specific text processing methods.

Previous studies have explored various aspects of sentiment analysis and complaint detection in e-commerce reviews. Traditional sentiment analysis categorizes reviews as positive, negative, or neutral but struggles to differentiate between general dissatisfaction and specific complaints. In contrast, aspect-based sentiment analysis provides a more detailed understanding by linking sentiments to particular aspects of a product or service. Recent advancements in transformer-based language models, like Bidirectional Encoder Representations from Transformers (BERT) and Generative Pre-trained Transformer (GPT), have greatly enhanced text classification and information extraction, opening up new possibilities for more precise complaint detection.

The novelty of this study lies in its practical and simplified approach to complaint detection. Instead of developing complex hybrid models, we focus on fine-tuning transformer-based models specifically for e-commerce complaints and leveraging simple rule-based techniques to enhance interpretability. By concentrating on key complaint categories (e.g., product quality, shipping issues, customer service), our method provides a more streamlined yet effective way to identify and classify customer grievances.

This paper aims to address the gap in automated complaint detection by developing a model that combines NLP techniques with machine learning classifiers to enhance the accuracy and interpretability of results. The proposed approach is evaluated on publicly available datasets of online store reviews to demonstrate its effectiveness in real-world scenarios. The findings of this study can be beneficial for businesses seeking to optimize their customer support processes and for researchers interested in advancing the field of NLP-based consumer feedback analysis.

Literature review. The automated identification and analysis of customer complaints in online store reviews using Natural Language Processing (NLP) has been extensively researched, leading to the development of various methods aimed at improving classification accuracy, processing efficiency, and real-time analytics. This section reviews key research contributions in this field, focusing on different methodological approaches and their impact.

Machine learning and deep learning models have been extensively used for classifying customer complaints. Recurrent Neural Networks (RNNs), when combined with sentiment analysis and text preprocessing, have demonstrated high classification accuracy. For instance, Taneja et al. (2019) achieved an 82% accuracy using an RNN-based model, while Khedkar and Shinde (2018) reported 79.83% accuracy with an ensemble classifier. Expanding on these techniques, Fong et al. (2021) integrated RNNs with Latent Dirichlet Allocation (LDA) topic modeling, attaining high ROC AUC scores of 0.930 and 0.894. Similarly, Tasmia and Quraishi (2022) explored the application of a Fast Text Deep Neural Network (DNN) for Bengali-language complaint detection, achieving a 74% accuracy rate.

Beyond supervised learning models, researchers have also explored lexicon-based and topic modeling methods for complaint detection. Adams et al. (2017) utilized specialized "smoke word" dictionaries to identify safety and efficacy concerns in pharmaceutical product reviews, demonstrating superior performance over traditional sentiment analysis techniques. Meanwhile, Omurca et al. (2021) and Zhan et al. (2009) applied topic modeling techniques such as Latent Dirichlet Allocation (LDA) and Gibbs Sampling for Dirichlet Multinomial Mixtures (GSDMM) to structure and categorize complaints effectively, thereby providing a more organized view of customer feedback.

Sentiment analysis has played a fundamental role in complaint detection, with various hybrid approaches emerging to enhance its effectiveness. For example, VADER-based sentiment analysis, as implemented by Taneja et al. (2019), demonstrated improved preprocessing when combined with RNNs. Additionally, hybrid models that integrate machine learning with rule-based techniques have been explored. Fong et al. (2021) combined RNNs with LDA to improve early defect detection, while Oelke et al. (2009) introduced a visual analytics approach to opinion mining, incorporating discrimination-based techniques for better interpretability and user insights.

The practical impact of automated NLP techniques on complaint analysis has been significant. Studies highlight the substantial improvements in efficiency and scalability that these methods provide. Taneja et al. (2019) reported reducing processing time from 10 days to under 30 seconds per review, showcasing the power of automation. Omurca et al. (2021) further contributed by developing a mobile application that visualizes topic distributions in customer complaints, enhancing accessibility for business users. Additionally, Singh et al. (2020) explored Neural Machine Translation (NMT) to process complaints in low-resource languages, underscoring the potential for multilingual applications in e-commerce settings.

Other researchers have also made significant contributions. Sulova (2016) proposed an approach for the automatic analysis of online store product and service reviews, integrating machine learning and NLP techniques. S. İ. Omurca et al. (2021) focused on detecting topics in customer complaints using artificial intelligence techniques. Their study, along with Zhan et al. (2009), emphasized topic extraction and summarization methods for better complaint classification. Furthermore, David Z. Adams et al. (2017) explored risk mitigation through automated detection of safety concerns, demonstrating the impact of NLP in sectors beyond e-commerce.

Many current models require significant computational power and large labeled datasets for training, which restricts their scalability. Furthermore, domain adaptation remains a significant hurdle, as models trained on one dataset often struggle to generalize effectively across different e-commerce platforms. Future studies should investigate methods like transfer learning, semi-supervised approaches, and fine-

tuning transformer-based models to improve adaptability and classification accuracy across different retail settings.

Existing research demonstrates that automated NLP techniques can successfully identify and analyze customer complaints, providing businesses with valuable insights. However, further advancements are needed to enhance model interpretability, scalability, and adaptability across different domains. This study helps address these challenges by utilizing transformer-based models and improving complaint classification techniques, aiming to further the use of NLP in e-commerce complaint analysis.

Methods. This study employs NLP techniques to automatically detect and analyze customer complaints in online store reviews. The dataset used is the Amazon US Customer Reviews Dataset, which includes millions of reviews across various product categories. The dataset includes key fields like review text, star rating, helpful votes, and verification status, providing a valuable resource for sentiment analysis and complaint classification. The research follows a systematic approach, encompassing data preprocessing, feature extraction, model development, evaluation, and interpretation.

The dataset undergoes preprocessing to enhance text quality and model performance. First, missing values, duplicates, and short reviews (fewer than three words) are removed. HyperText Markup Language (HTML) tags, emojis, special characters, and Uniform Resource Locators (URL) are stripped from the text. Tokenization and lemmatization are applied using spaCy, while common stopwords from Natural Language Toolkit (NLTK) are removed to focus on meaningful words. Reviews are categorized by sentiment, where 1-2 star reviews are labeled as negative (potential complaints), 3-star reviews as neutral, and 4-5 star reviews as positive. Given the imbalance in complaint distribution, the Synthetic Minority Over-sampling Technique (SMOTE) is applied to balance the dataset.

Feature engineering plays a crucial role in enhancing classification accuracy. Aspect-Based Sentiment Analysis (ABSA) is implemented using Named Entity Recognition (NER) and dependency parsing to extract complaint-related aspects such as "product quality" or "delivery issues." Sentences with keywords like "arrived late," "damaged," and "poor quality" are flagged. Term Frequency–Inverse Document Frequency (TF-IDF) is applied to traditional models, while Word2Vec and fastText embeddings are used for deep learning models. Bigrams and trigrams are extracted to capture sequential complaint patterns, and Latent Dirichlet Allocation (LDA) is used to identify key complaint topics.

Three models are developed for classification: a Naïve Bayes classifier using TF-IDF features as a baseline, a BERT-based classifier fine-tuned on the dataset, and a hybrid model combining rule-based filtering with deep learning. The BERT model is trained using Transformers from Hugging Face with the AdamW optimizer and a learning rate of $2e-5$, running for five epochs. The dataset is split into 80% training,

10% validation, and 10% testing. GridSearchCV optimizes the Naïve Bayes classifier, while Bayesian optimization fine-tunes BERT's batch size, dropout rate, and attention heads.

Evaluation metrics include accuracy, precision, recall, and F1-score, with ROC-AUC measuring classification performance. A confusion matrix visualizes false positives and false negatives, while SHAP values and attention heatmaps from BERT aid in model interpretability.

Results. The BERT-based classifier significantly outperforms the Naïve Bayes baseline. The Naïve Bayes model achieves 74.2% accuracy with an F1-score of 69.8%, while the BERT classifier improves accuracy to 85.6% and F1-score to 82.3%. The hybrid model further enhances performance, reaching 87.1% accuracy and an F1-score of 85.2%. These results demonstrate the advantages of deep learning models in accurately identifying complaints.

Table 1: Classification Performance Results

Model	Model Evaluation Metrics			
	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1-score</i>
Naïve Bayes	74.2%	71.5%	68.2%	69.8%
BERT Classifier	85.6%	84.1%	81.2%	82.3%
Hybrid Model	87.1%	86.5%	84.0%	85.2%

The Aspect-Based Sentiment Analysis (ABSA) successfully categorizes complaints, revealing that 42% of negative reviews pertain to product quality issues, 35% to delivery problems, and 23% to customer service dissatisfaction. Further analysis indicates that verified purchase reviews are 68% more likely to contain genuine complaints compared to unverified reviews, highlighting the importance of filtering for reliability.

The confusion matrix below shows the classification performance of the BERT model, where most misclassifications occur in neutral (3-star) reviews due to their ambiguous sentiment.

Table 2: Confusion Matrix

	Predicted Complaint	Predicted Non-Complaint
Actual Complaint	4325	610
Actual Non-Complaint	780	12540

The Receiver Operating Characteristic (ROC) curve shows the trade-off between true positive and false positive rates for different classification thresholds. The BERT model achieves an AUC of 0.92, indicating strong discrimination between complaints and non-complaints.

A detailed error analysis reveals several challenges. Neutral reviews (3 stars) are difficult to classify due to mixed sentiments. Short complaints (e.g., "Terrible!" or "Not good.") lack contextual detail, reducing model confidence. Domain-specific language introduces classification errors, particularly for niche product categories. Addressing these issues requires more refined feature engineering and specialized fine-tuning.

To further improve accuracy and applicability, several areas are identified for future research. Expanding dataset coverage by incorporating multilingual complaints will allow cross-regional comparison. Explainable AI (XAI) techniques such as Local Interpretable Model-agnostic Explanations (LIME) and Shapley Additive Explanations (SHAP) can improve model transparency, making decisions more interpretable. A real-time complaint detection system using Flask or FastAPI would enable e-commerce platforms to process complaints dynamically. Graph-based aspect refinement using Graph Neural Networks (GNNs) could enhance aspect extraction, reducing misclassification. Finally, adaptive complaint classification via transfer learning will improve model performance across different e-commerce domains.

This study highlights the potential of NLP-driven complaint detection to improve customer experience analysis. By integrating deep learning, aspect-based sentiment analysis, and real-time processing, businesses can better address customer concerns, ultimately enhancing service quality and brand reputation.

Discussions and Conclusions. This study demonstrates the effectiveness of Natural Language Processing techniques for automating the detection and classification of customer complaints in online store reviews. The results indicate that deep learning models, particularly fine-tuned transformer-based architectures like BERT, significantly outperform traditional approaches such as Naïve Bayes in detecting complaints with high accuracy and reliability. The implementation of Aspect-Based Sentiment Analysis (ABSA) further enhances interpretability by identifying key complaint categories such as product quality, delivery issues, and customer service dissatisfaction. The hybrid model, combining rule-based filtering with deep learning, provides the most robust performance, minimizing false positives and improving overall classification precision. However, challenges remain in accurately classifying neutral reviews and handling domain-specific terminology. The findings suggest that verified purchase reviews are a more reliable source of complaint data, which can be leveraged for improving automated complaint detection systems. Future work should explore multilingual adaptation, real-time processing, and the integration of Explainable AI (XAI) techniques to enhance model transparency and usability in business applications. By advancing NLP-based complaint detection, this research contributes to improving customer service analytics and optimizing response strategies for e-commerce platforms.

REFERENCES

1. Adams, D. Z., Gruss, R., & Abrahams, A. (2017). Automated discovery of safety and efficacy concerns for joint & muscle pain relief treatments from online reviews. *International Journal of Medical Informatics*, 103, 1–12.
2. Fong, T. H. Y., Sarkani, S., & Fossaceca, J. (2021). Auto defect detection using customer reviews for product recall insurance analysis. *Frontiers in Applied Mathematics and Statistics*, 7, 632847.
3. Khedkar, S. A., & Shinde, S. K. (2018). Customer review analytics for business intelligence. *Proceedings of the International Conference on Computational Intelligence and Computing Research (ICCIC)*, 1–6.
4. Oelke, D., Hao, M., Rohrdantz, C., Keim, D., Dayal, U., Haug, L., & Janetzko, H. (2009). Visual opinion analysis of customer feedback data. *2009 IEEE Symposium on Visual Analytics Science and Technology*, 187–194.
5. Omurca, S. İ., Ekinci, E., Yakupoğlu, E., Arslan, E., & Çapar, B. (2021). Automatic detection of the topics in customer complaints with artificial intelligence. *Balkan Journal of Electrical and Computer Engineering*, 9(2), 123–132.
6. Singh, R., Haque, R., & Hasanuzzaman, M. (2020). Identifying complaints from product reviews in low-resource scenarios via neural machine translation. *Proceedings of the ACL Workshop on NLP for Low-Resource Languages*, 45–53.
7. Sulova, S. (2016). An approach for automatic analysis of online store product and services reviews. *Proceedings of the International Conference on Advances in Computing, Communications and Informatics (ICACCI)*, 89–94.
8. Taneja, S., Jha, K., Lakhotia, N., Kapoor, V., & Swarnalatha, P. (2019). Customer feedback analyzer. *International Journal of Engineering and Advanced Technology*, 9(2), 123–130.
9. Tasmia, F. S., Quraishi, A. H. (2022). Customer complaint detection in e-commerce platforms using NLP. *International Journal of Computer Applications*, 182(3), 21–27.
10. Zhan, J., Loh, H., & Liu, Y. (2009). Gather customer concerns from online product reviews - A text summarization approach. *Expert Systems with Applications*, 36(2), 2107–2115.

АВТОМАТИЗИРОВАННОЕ ОБНАРУЖЕНИЕ И АНАЛИЗ ЖАЛОБ КЛИЕНТОВ В ОТЗЫВАХ ИНТЕРНЕТ-МАГАЗИНОВ С ИСПОЛЬЗОВАНИЕМ NLP

Азиева А.

В данном исследовании применяются методы обработки естественного языка (Natural Language Processing, NLP) для выявления и анализа жалоб клиентов в отзывах интернет-магазинов на основе датасета Amazon US Customer Reviews. Методология включает предобработку данных, извлечение признаков и классификацию с использованием Наивного Байеса, модели BERT и гибридной модели, сочетающей фильтрацию на основе правил с методами

глубокого обучения. Результаты показали, что модель BERT превосходит традиционные алгоритмы. Анализ выявил, что отзывы от проверенных покупателей чаще содержат подлинные жалобы. Несмотря на достигнутые успехи, остаются проблемы с классификацией нейтральных отзывов и обработкой специализированной терминологии. В дальнейшем планируется расширение на мультязычные данные, реализация обработки в реальном времени и интеграция объяснимого ИИ (Explainable AI) для повышения прозрачности моделей и их бизнес-применимости.

Ключевые слова: Обработка естественного языка, сентимент-анализ, выявление жалоб, BERT, гибридная модель.

АВТОМАТТАНДЫРЫЛҒАН ШАҒЫМДАРДЫ АНЫҚТАУ ЖӘНЕ ТАЛДАУ: ONLINE ДҮКЕН ПІКІРЛЕРІН NLP АРҚЫЛЫ ЗЕРТТЕУ

Азиева А.

Бұл зерттеу Amazon АҚШ тұтынушы пікірлерінің деректер жиынтығын пайдаланып, онлайн дүкен пікірлеріндегі тұтынушы шағымдарын анықтау және талдау үшін табиғи тілді өңдеу (NLP) әдістерін қолданады. Әдістеме мәліметтерді алдын ала өңдеуді, белгілерді (фичаларды) шығаруды және Наив Байес, BERT және ережеге негізделген сүзгілеуді терең оқытумен біріктіретін гибридік модельді қолданатын классификацияны қамтиды. Нәтижелерге сәйкес, дәстүрлі модельдермен салыстырғанда BERT анағұрлым тиімді нәтиже көрсетеді. Талдау расталған сатып алуларға қатысты пікірлерде нақты шағымдардың жиі кездесетінін көрсетті. Алайда, бейтарап пікірлерді жіктеу және салалық терминологияны дұрыс түсіну секілді қиындықтар әлі де сақталып отыр. Болашақта зерттеу бағыттары көптілді кеңейтуге, нақты уақыттағы анықтауға және түсінікті Жасанды Интеллект (Explainable AI) элементтерін енгізуге бағытталады, бұл модельдің ашықтығы мен бизнестік қолданылуын арттыруға сеп болады.

Кілт сөздер: Тілді өңдеу, сентимент-анализ, шағымдарды анықтау, BERT, гибридік модель.