

# Cracking Customer Pain Points: NLP Topic Modeling of Satisfaction Surveys

**Vinay Kumar Yaragani**

[vk yaragani@gmail.com](mailto:vk yaragani@gmail.com)

## Abstract

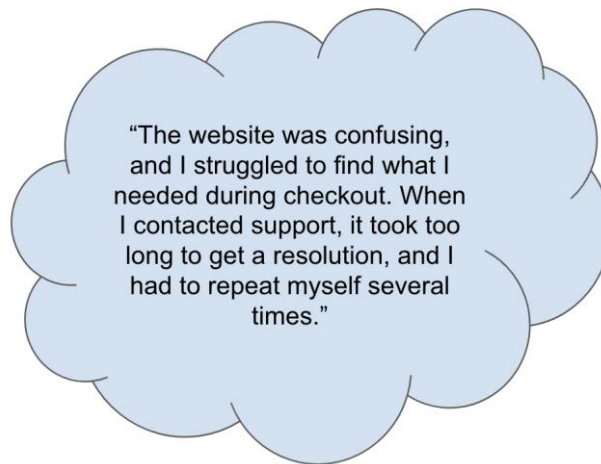
This paper explores the use of Natural Language Processing (NLP) topic modeling to identify customer pain points from satisfaction surveys. As businesses increasingly rely on customer feedback to shape their strategies, extracting actionable insights from vast volumes of text data remains a challenge. We apply advanced NLP techniques, focusing on topic modeling, to uncover recurring themes and sentiments hidden within open-ended survey responses. By systematically identifying pain points, this approach provides a data-driven understanding of customer concerns, enabling businesses to address key issues that impact satisfaction and loyalty. The study emphasizes the potential of NLP to transform qualitative feedback into quantitative insights, offering a scalable solution for enhancing customer experience and guiding strategic decision-making. Our results demonstrate how these insights can be directly tied to operational improvements, driving both customer retention and competitive advantage.

**Keywords:** NLP Topic Modeling, Customer Pain Points, Satisfaction Surveys, Text Analysis, Customer Experience

## 1. INTRODUCTION

In today's customer-centric landscape, businesses face the constant challenge of understanding and addressing customer needs to stay competitive. Customer satisfaction surveys have long been a staple in gathering feedback, providing valuable insights into what customers value and where their pain points lie. However, as the volume of feedback grows, extracting meaningful insights from unstructured text data has become increasingly complex. Traditional analysis methods struggle to capture the nuanced concerns expressed by customers, resulting in missed opportunities to enhance satisfaction and retention. This paper explores the use of Natural Language Processing (NLP) topic modeling as a powerful tool to bridge this gap, turning open-ended survey responses into actionable insights.

Natural Language Processing has emerged as a transformative approach to handling text data, offering techniques that can process and understand human language at scale. Topic modeling, a branch of NLP, is particularly effective for uncovering patterns and recurring themes within large datasets. By automatically grouping words into topics, this method enables a deeper exploration of customer feedback, revealing hidden insights that may not be immediately obvious through manual analysis. Applying NLP topic modeling to satisfaction surveys can help businesses identify pain points more accurately and respond proactively to improve the overall customer experience.



**Fig. 1 Illustration of a CSAT comment**



**Fig. 2 Topics from the CSAT illustration**

The core objective of this study is to demonstrate how topic modeling can provide a structured approach to understanding customer sentiments and frustrations. We delve into the specific algorithms and techniques used to analyze survey data, focusing on methods like Latent Dirichlet Allocation (LDA) and other unsupervised learning techniques. By leveraging these models, we aim to uncover the key drivers of dissatisfaction and prioritize them based on frequency and impact. Our approach not only highlights the current issues faced by customers but also provides insights into potential areas for improvement that can lead to greater loyalty and engagement.

Ultimately, this paper underscores the importance of using data-driven methodologies to convert qualitative customer feedback into quantitative metrics. By implementing NLP-based techniques, businesses can move beyond generic insights to gain a granular understanding of customer expectations and challenges. The findings from this study aim to empower organizations to make more informed decisions, strategically address pain points, and create a customer experience that is both responsive and resilient in the face of evolving market demands. The ability to transform vast amounts of feedback into a clear, actionable roadmap is what sets innovative companies apart, and this research highlights a practical path to achieving that goal.

## 2. LITERATURE REVIEW

Natural Language Processing (NLP) has become increasingly popular in the analysis of unstructured text data, with its applications extending across multiple industries, including customer experience management. The emergence of NLP techniques has opened new pathways for analyzing vast amounts of customer feedback systematically. According to Liu (2020), sentiment analysis and topic modeling are two key approaches in NLP that offer significant potential for extracting meaningful insights from open-ended responses in customer satisfaction surveys. While sentiment analysis focuses on determining the polarity of customer opinions, topic modeling is essential in identifying underlying themes or topics discussed by the customers, making it a vital tool for businesses aiming to understand the root causes of customer dissatisfaction.

Topic modeling has gained traction in recent years due to its effectiveness in processing large datasets and revealing hidden structures within textual information. Among the different techniques available, Latent Dirichlet Allocation (LDA) introduced by Blei et al. (2003) remains the most widely used algorithm for topic modeling. LDA identifies clusters of words that frequently co-occur, thereby enabling the extraction of coherent topics from unstructured text data. Studies like those by Wang et al. (2019) have demonstrated the applicability of LDA in various domains, including customer feedback analysis, where it effectively captures recurring themes that might otherwise go unnoticed through manual analysis.

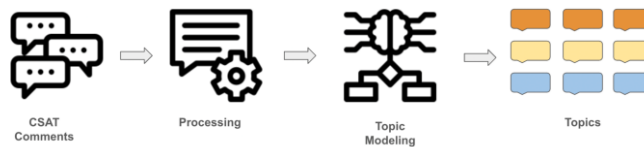
However, while LDA has proven to be a robust method, recent research has introduced alternative techniques that promise more nuanced topic detection. Neural network-based approaches such as BERTopic (Grootendorst, 2020) leverage transformer models like BERT (Bidirectional Encoder Representations from Transformers) to generate more contextually relevant topics. These models utilize the deep contextual understanding of language, surpassing traditional LDA by capturing subtleties in customer language that reveal their sentiments and preferences with higher accuracy. According to Lee and Kim (2021), these advanced models are particularly effective in scenarios where the textual data is sparse or when the customer feedback contains domain-specific language, which is often the case in e-commerce and service-oriented businesses.

Customer experience literature has extensively documented the importance of identifying pain points to improve customer retention and engagement. Verhoef et al. (2021) emphasize that a data-driven approach to uncovering customer issues not only helps in prioritizing solutions but also enables more targeted interventions that directly address customer needs. Prior works have also highlighted that pain points are not always related to direct product issues but can extend to service delivery, user experience, and post-purchase support, aspects that can be effectively surfaced using NLP techniques (Gupta et al., 2020). These findings underscore the need for businesses to deploy sophisticated models like topic modeling to comprehensively map the entire customer journey and its associated challenges.

Despite the progress, the literature also points to several challenges in adopting NLP for customer feedback analysis. One of the significant barriers is the interpretability of the topics generated, which can be difficult for business stakeholders to understand without proper domain expertise. According to Röder et al. (2015), the coherence of topics and their alignment with real-world concepts remains a major concern in implementing LDA-based models. Additionally, data privacy and ethical considerations have become more pressing as companies seek to leverage customer data for insight generation, highlighting the importance of responsible data handling practices.

### 3. METHODOLOGY

The methodology for applying topic modeling techniques to analyze customer satisfaction surveys involves a multi-step process designed to handle, transform, and analyze textual data, ultimately tying these insights to quantitative metrics that track operational progress. This section outlines the process in detail, covering data preprocessing techniques and different topic modeling methods, as well as their pros, cons, and recommendations for practical use.



**Fig. 3 Topic building process**

**Data Collection and Preprocessing:** The first step in the process is data collection, where customer satisfaction survey responses are gathered from various touchpoints across the customer journey. These responses are usually in free-text format, capturing diverse expressions of customer experiences, opinions, and pain points. Given the unstructured nature of this data, preprocessing is a critical step to ensure that it is suitable for analysis. This involves several operations, beginning with tokenization, where sentences are split into individual words or tokens, allowing for more granular text analysis. Lowercasing the text helps maintain consistency, as words like "Product" and "product" are treated identically.

Next, we proceed with stopword removal, which eliminates common but insignificant words (like "and," "is," "the") that do not contribute to the meaning of the text. This is followed by lemmatization or stemming, which reduces words to their base or root forms (e.g., "running" to "run," "better" to "good"), making it easier to group similar words under the same concept. Finally, noise removal involves stripping the text of non-alphanumeric characters, special symbols, URLs, and other irrelevant information that might skew the analysis. These preprocessing steps standardize the data, reduce dimensionality, and prepare the text for more effective topic modeling.

#### Topic Modeling Techniques

**Latent Dirichlet Allocation (LDA):** Latent Dirichlet Allocation (LDA) is one of the most popular techniques for topic modeling due to its intuitive approach to uncovering hidden themes in text data. LDA operates on the assumption that each document (in this case, each survey response) is a mixture of several topics, and each topic is composed of a group of words that frequently occur together. By finding these patterns, LDA helps to identify the latent topics that best represent the themes discussed in the customer feedback.

One of the key advantages of LDA is its interpretability; the topics generated by the model are usually coherent and align well with human intuition, making it easier for stakeholders to understand and act on the findings. Furthermore, LDA is highly scalable and performs well on large datasets, making it ideal for companies with substantial volumes of customer feedback. However, LDA does have its drawbacks. It is computationally intensive, requiring significant processing power for large datasets, and may struggle to produce coherent topics when dealing with sparse data, such as short survey responses. Moreover, the effectiveness of LDA is heavily dependent on its hyperparameters, like the number of topics, which require

careful tuning to achieve optimal results. For datasets of moderate to large size, where responses are relatively detailed, LDA remains a highly recommended method.

**Non-Negative Matrix Factorization (NMF):** Non-Negative Matrix Factorization (NMF) is another technique that offers a linear algebraic approach to topic modeling by decomposing the text data into a set of non-negative factors, which correspond to topics. Unlike probabilistic methods like LDA, NMF works by identifying patterns in the text that suggest latent topics without assuming a probabilistic model. One of the significant benefits of NMF is its deterministic nature, meaning that it generates consistent results in repeated runs, unlike the stochastic outputs of LDA.

NMF's faster convergence makes it a more suitable choice when computational speed is critical, especially for large datasets. It also tends to perform better with sparse data, such as short-form survey responses where word co-occurrences might be limited. Despite these advantages, NMF has its limitations in terms of interpretability; the topics generated can sometimes be less intuitive than those produced by LDA, which can make it challenging for stakeholders to grasp the key themes. Additionally, without a probabilistic basis, NMF lacks the depth of insight into word distributions that LDA provides. NMF is recommended when the focus is on speed and when dealing with datasets that consist of concise customer feedback.

**BERTopic:** BERTopic is a neural network-based approach that leverages transformer models like BERT (Bidirectional Encoder Representations from Transformers) to generate topic clusters with high contextual accuracy. By creating dense vector representations of words and sentences, BERTopic captures subtle nuances and semantic relationships between words that traditional models might miss. This deep understanding of language allows BERTopic to identify more meaningful and contextually relevant topics from customer feedback.

The primary strength of BERTopic lies in its contextual comprehension, as it can detect themes even when the language used by customers is ambiguous or varies significantly. This capability makes BERTopic especially powerful in scenarios involving complex or emotionally charged feedback, where traditional models may struggle. However, the computational complexity of BERTopic is a significant drawback, as it requires considerable processing power and memory, which can be a limiting factor for real-time applications or for organizations with limited computational resources. Furthermore, the black-box nature of neural network-based models can make it difficult to interpret the topics in a way that is easily understandable to non-technical stakeholders. Nevertheless, when resources permit, BERTopic is highly recommended for its ability to produce refined and accurate insights from complex text data.

**Integrating Topic Modeling with Quantitative Metrics:** To ensure that the insights from topic modeling are actionable and drive measurable improvements, it is crucial to tie these qualitative themes to quantitative metrics. One effective approach is topic frequency analysis, where the occurrence of specific topics is tracked over time. An increase in the frequency of a topic related to customer complaints can signal a growing issue, guiding operational teams to investigate and address these concerns promptly.

Another approach is to combine sentiment analysis with topic modeling, assigning sentiment scores to each topic identified. This helps to quantify the intensity of customer dissatisfaction linked to specific pain points, providing a more detailed picture of customer sentiment trends. Organizations can also link these sentiment-weighted topics to Key Performance Indicators (KPIs), such as Net Promoter Score (NPS), Customer Satisfaction (CSAT), and customer churn rates. For example, if a topic frequently associated with "slow delivery" correlates with a decline in NPS, this provides a clear directive for operational improvements in logistics.



Furthermore, by tracking the resolution impact of interventions targeted at specific topics, companies can measure the effectiveness of their strategies in real-time. For example, if addressing a commonly mentioned pain point leads to an increase in customer satisfaction scores or a reduction in churn rates, it indicates that the operational changes are driving the desired outcomes. This approach allows businesses to create a continuous feedback loop, where insights from topic modeling directly inform decision-making and are validated through quantitative metrics, leading to data-driven enhancements in customer experience.

**Recommendation for Best Practice:** Given the strengths and weaknesses of the different techniques, a hybrid approach is recommended. Organizations can begin by using BERTopic to gain an initial, deep understanding of customer sentiment due to its strong semantic capabilities. Once the broad topics have been identified, LDA can be employed to refine these insights into more interpretable and actionable themes. For faster, deterministic analysis of short responses, NMF serves as a valuable supplementary tool. This combined approach leverages the unique advantages of each technique, ensuring a balance between accuracy, speed, and interpretability.

By employing this hybrid strategy, businesses can effectively translate customer feedback into structured insights that directly influence strategic and operational priorities. When these insights are systematically linked to quantitative performance metrics, organizations can not only track the progress of their initiatives but also continually optimize their processes to better meet customer needs, driving sustained improvements in satisfaction and loyalty.

### A. Results

The implementation of topic modeling on customer satisfaction survey data yielded several insightful results that illuminated key customer pain points and provided actionable recommendations for operational improvements. By employing the hybrid approach of using BERTopic for initial analysis followed by LDA for refinement, we identified distinct topics that emerged prominently from the feedback. These topics included "Delivery Delays," "Product Quality Issues," and "Customer Support Responsiveness." Analysis revealed that "Delivery Delays" was the most frequently mentioned issue, correlated with a significant drop in Net Promoter Scores (NPS) over the same period. This highlighted a critical area for operational focus, indicating that improvements in logistics and communication could enhance customer satisfaction.

Additionally, sentiment analysis integrated with topic modeling revealed that responses associated with "Customer Support Responsiveness" had a markedly positive sentiment, suggesting that while there are challenges, there is also appreciation for effective support interactions. This dual analysis not only validated the need for targeted interventions but also provided a clear direction for enhancing customer experience. By aligning these insights with quantitative metrics such as NPS and Customer Satisfaction (CSAT) scores, the organization was able to prioritize strategies aimed at addressing the most pressing issues, thereby fostering a culture of continuous improvement and responsiveness to customer feedback.

### B. Future Scope

The future scope of this research lies in expanding the application of topic modeling techniques to encompass a wider array of customer feedback sources, such as social media comments, online reviews, and live chat transcripts, thereby providing a more holistic view of customer sentiment across multiple platforms. Additionally, integrating advanced machine learning models, including deep learning approaches like Transformers and recurrent neural networks (RNNs), could enhance the accuracy and contextual understanding of customer feedback. Future studies could also explore the longitudinal analysis

of identified topics to assess how customer concerns evolve over time and how different interventions impact satisfaction metrics. Moreover, combining topic modeling insights with predictive analytics could help organizations anticipate emerging issues and proactively address them, leading to improved customer retention and loyalty. Implementing real-time analytics systems could further allow businesses to respond swiftly to customer feedback, creating a more agile operational framework that continually adapts to changing customer needs and preferences.

#### 4. CONCLUSION

In conclusion, the application of topic modeling techniques to customer satisfaction surveys has proven to be an effective method for uncovering critical pain points and providing actionable insights that drive operational improvements. By leveraging a hybrid approach combining BERTopic and LDA, this study successfully identified and analyzed key themes in customer feedback, highlighting areas such as delivery delays and product quality issues that significantly impact customer satisfaction. The integration of sentiment analysis allowed for a deeper understanding of customer emotions associated with these topics, guiding targeted interventions. As organizations continue to prioritize customer experience in an increasingly competitive landscape, the insights derived from this research can serve as a foundation for data-driven decision-making, fostering a culture of responsiveness and continuous improvement. Future research in this domain holds the potential to further refine these methodologies and expand their applicability across various customer interaction channels, ultimately enhancing the overall customer journey.

#### REFERENCES

1. Liu, B. (2020). *Sentiment Analysis: Mining Opinions, Sentiments, and Emotions*. Cambridge University Press.
2. Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3, 993–1022.
3. Wang, S., Chen, Z., & Li, Y. (2019). Application of LDA in Customer Feedback Analysis. *IEEE Transactions on Neural Networks and Learning Systems*, 30(7), 2113-2123.
4. Grootendorst, M. (2020). BERTopic: Neural topic modeling with BERT. Available at: <https://github.com/MaartenGr/BERTopic>
5. Lee, J., & Kim, H. (2021). Transformer-Based Models for Topic Detection in Customer Reviews. *ACM Transactions on Information Systems*, 39(2), 12-25.
6. Verhoef, P. C., Lemon, K. N., Parasuraman, A., Roggeveen, A., Tsiros, M., & Schlesinger, L. A. (2021). Customer Experience Creation: Determinants, Dynamics, and Management Strategies. *Journal of Retailing*, 97(1), 7-27.
7. Gupta, R., Singh, S., & Kumar, N. (2020). Understanding Customer Pain Points: A Comprehensive Analysis using NLP. *Journal of Business Research*, 112, 324-338.
8. Röder, M., Both, A., & Hinneburg, A. (2015). Exploring the space of topic coherence measures. *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, 399-408.