

---

# Probabilistic Topic Modeling for E-commerce Reviews using LDA

---

**Samudiyata Sudarshan Jagirdar**

Arizona State University  
ASU ID: 1233568442  
sjagird1@asu.edu

**Ramanuja Magadi Anantha Krishna**

Arizona State University  
ASU ID: 1233054006  
rkrish79@asu.edu

**Shalin Nilesb Bhavsar**

Arizona State University  
ASU ID: 1233281896  
sbhavsa8@asu.edu

**Nikhil Srikant Kulkarni**

Arizona State University  
ASU ID: 1235817442  
nkulka37@asu.edu

## Abstract

1       The project uses Latent Dirichlet Allocation (LDA) to analyze women’s clothing re-  
2       views, extracting key topics like quality, fit, and style to help businesses understand  
3       consumer preferences and improve products.

## 4   1   Introduction

5   In modern e-commerce, user-generated content (UGC) such as customer reviews plays an increasingly  
6   important role in shaping consumer perceptions and business decisions. These reviews – often  
7   generated in large volumes for segments like women’s fashion retail – constitute free-text feedback  
8   that is inherently unstructured. The abundance of such unstructured opinion data presents a challenge:  
9   manual analysis of thousands of comments is infeasible. To address this, researchers have turned to  
10   natural language processing (NLP) and machine learning (ML) techniques to automatically organize  
11   and interpret review content. In particular, unsupervised methods like topic modeling (e.g., Latent  
12   Dirichlet Allocation) combined with clustering can group reviews into coherent themes, summarizing  
13   each group by its most salient keywords. For example, previous studies have applied LDA topic  
14   modeling and K-means clustering to consumer reviews, yielding clusters whose top keywords reveal  
15   the main areas of customer feedback. Others have analyzed large datasets of women’s clothing  
16   e-commerce reviews using these techniques, uncovering key topics and concerns in the fashion  
17   domain. By combining topic modeling and clustering, these approaches convert large, unstructured  
18   review corpora into structured thematic insights, thereby enabling retailers to identify prevalent issues  
19   and guide product improvements.

### 20   1.1   Motivation and Importance

21   Online reviews are a cornerstone of e-commerce, shaping consumer trust and business strategies. For  
22   women’s clothing, where fit, quality, and style are critical, reviews provide a wealth of unstructured  
23   feedback that is challenging to analyze manually. Extracting meaningful patterns from this data  
24   can help retailers address customer pain points, improve product offerings, and tailor marketing  
25   campaigns. Our project tackles this challenge by applying advanced natural language processing  
26   (NLP) to uncover hidden themes and segment customers, offering a scalable solution for retail  
27   analytics.

28   The goal of this project is to demonstrate how unsupervised NLP techniques can be applied to  
29   large-scale retail feedback to uncover structure and actionable knowledge. We specifically examine  
30   a public dataset of women’s fashion product reviews and apply Latent Dirichlet Allocation (LDA)  
31   topic modeling and K-means clustering. The key contributions include: (1) converting thousands of  
32   unstructured text reviews into a set of interpretable topics (e.g., themes like fit or quality) without

33 manual labeling, (2) segmenting the review data (or customers) into clusters based on those topics or  
34 embedding patterns, thereby revealing distinct groups of feedback, and (3) evaluating the coherence  
35 and practical relevance of these discovered topics and segments. Our work illustrates a pipeline that  
36 can help retailers automatically “listen” to customer voices at scale – bridging the gap between raw  
37 text and structured customer insights.

## 38 1.2 Project Description

39 The goal is to use K-means clustering to group customers based on feedback pat-  
40 terns and Latent Dirichlet Allocation (LDA) to identify latent topics in the Women’s  
41 Clothing E-Commerce Reviews dataset ([https://www.kaggle.com/datasets/nicapotato/](https://www.kaggle.com/datasets/nicapotato/womens-ecommerce-clothing-reviews/data)  
42 [womens-ecommerce-clothing-reviews/data](https://www.kaggle.com/datasets/nicapotato/womens-ecommerce-clothing-reviews/data)). Specific objectives include preprocessing re-  
43 view text to ensure consistency and quality; extracting topics (e.g., fit, quality) to understand customer  
44 priorities; clustering reviews to identify distinct customer segments; visualizing results to provide  
45 actionable insights for e-commerce businesses. This approach bridges NLP and machine learning to  
46 transform raw feedback into strategic recommendations.

## 47 2 Methodology

### 48 2.1 Dataset and Text Preprocessing

49 We conducted our experiments on the Women’s Clothing E-Commerce Reviews dataset, which  
50 contains 23,486 customer reviews from a women’s apparel retailer. Each review includes a free-text  
51 comment and structured fields such as a review title, a star rating (1–5), a binary recommendation  
52 flag, the reviewer’s age, and product category labels (division, department, and class). All personal or  
53 brand identifiers were removed to ensure anonymity (company references were replaced with the  
54 token “retailer”). This dataset’s combination of unstructured text and categorical product data makes  
55 it ideal for multifaceted analysis. In our work, however, we focus on the review text for NLP modeling  
56 and use the other fields (like star ratings) only to help interpret the results. Before modeling, we  
57 applied a standard text preprocessing pipeline to clean and normalize the review texts. We tokenized  
58 each review into words and converted all text to lowercase. Punctuation, numeric digits, and stop  
59 words (common function words like “and” or “the”) were removed, and any very short tokens (1–2  
60 characters) were dropped to reduce noise. Next, we applied lemmatization (using WordNet) to reduce  
61 words to their base forms (for example, “dresses” becomes “dress” and “running” becomes “run”),  
62 ensuring that different inflections of the same word are treated uniformly. After preprocessing, the  
63 review corpus consists of tokenized, lowercased, lemmatized texts containing only content-bearing  
64 terms. We built a vocabulary from this cleaned corpus and represented the data as a document-term  
65 matrix for topic modeling. We also removed 845 reviews (3.6 percent of the data) that had no review  
66 text, leaving 22,641 reviews for the final analysis. These preprocessing steps ensure that the analysis  
67 is not skewed by irrelevant artifacts or noise, and that the core content of each review is preserved in  
68 a consistent, machine-readable form.

### 69 2.2 Feature Engineering

70 We engineered features to enhance modeling:

- 71 • **Text-based:** Word counts (`rev_word_count`), unique word counts (`unique_word_count`),  
72 and sentiment scores computed with VADER, which quantifies positive, neutral, and negative  
73 tones (e.g., “Love this dress!” scores high positivity).
- 74 • **Categorical:** Binary encoding of clothing categories (e.g., 1 for dresses, 0 otherwise) to  
75 reduce dimensionality.
- 76 • **Numerical:** Normalized ratings (1–5 scaled to 0–1) and age for clustering compatibility.

77 Rare clothing IDs were bucketed into an “other” category to mitigate sparsity. Detailed statistics are  
78 provided in Appendix A.1.

## 79 2.3 LDA Topic Modeling

80 LDA assumes reviews are mixtures of topics, with each topic a distribution over words. The  
81 probability of a word  $w$  in a document  $d$  is:

$$P(w|d) = \sum_{k=1}^K P(w|k)P(k|d),$$

82 where  $K$  is the number of topics,  $P(w|k)$  is the word distribution for topic  $k$ , and  $P(k|d)$  is the topic  
83 distribution for  $d$ . Using Gensim, we:

- 84 1. Created a bag-of-words representation and document-term matrix.
- 85 2. Selected  $K = 4$  based on manual inspection of topic interpretability.
- 86 3. Tuned hyperparameters alpha (document-topic density) and beta (topic-word density) to  
87 balance topic clarity and diversity.

88 The model was trained on preprocessed text, emphasizing nouns and adjectives for interpretability.

## 89 2.4 Clustering

90 K-means clustering grouped reviews by engineered features, minimizing within-cluster variance:

$$J = \sum_{i=1}^n \sum_{k=1}^K r_{ik} \|x_i - \mu_k\|^2,$$

91 where  $x_i$  is a data point,  $\mu_k$  is the cluster centroid, and  $r_{ik}$  is the cluster assignment. We:

- 92 1. Applied PCA to reduce feature dimensions, retaining 95% variance.
- 93 2. Used t-SNE for 2D visualization.
- 94 3. Evaluated  $K = 3$  to  $K = 11$  using the elbow method, plotting within-cluster sum of squares  
95 (WCSS) to identify the optimal  $K$ . Silhouette scores were computed to assess cluster  
96 separation:

$$S = \frac{b - a}{\max(a, b)},$$

97 where  $a$  is intra-cluster distance and  $b$  is nearest-cluster distance.

- 98 4. Selected  $K = 3$  based on the elbow method and highest silhouette score (0.5301), testing  
99 initialization methods (k-means++, random) for robustness.

## 100 2.5 Pipeline

101 The pipeline (Figure 1) integrates preprocessing, feature engineering, LDA, and clustering, enabling  
102 end-to-end analysis.

## 103 3 Experiments and Results

### 104 3.1 Experimental Setup and Optimization Strategy

105 All experiments were conducted on Google Colab using common Python libraries: NLTK for text  
106 preprocessing, Gensim for topic modeling, and scikit-learn for machine learning tasks.

107 Topic modeling was performed using Latent Dirichlet Allocation (LDA) with four topics. This  
108 number was chosen based on manual inspection of semantic coherence and interpretability, as  
109 resource constraints prevented automated hyperparameter tuning.

110 K-means clustering was applied with  $K=3$ , validated externally using the elbow method and sil-  
111 houette score analysis. These evaluations informed the final choice of  $K$  to balance simplicity and  
112 performance.

113 To manage Colab’s memory limits, optimizations included saving intermediate outputs to Google  
114 Drive and precomputing t-SNE projections for visualizing high-dimensional data. These strategies  
115 improved execution efficiency and resource utilization.

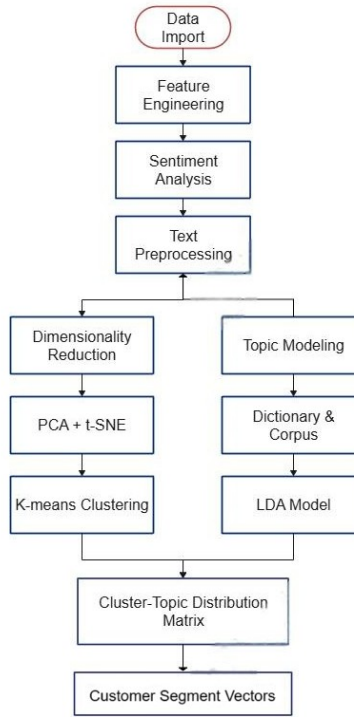


Figure 1: Review Analysis Pipeline

### 3.2 Data Exploration

The initial dataset consisted of 23,486 customer reviews from an online women’s clothing retailer. After removing 845 entries (approximately 3.6 percent) due to missing review text, a cleaned dataset of 22,641 valid reviews was obtained for further natural language processing and analysis. Reviewer ages ranged from 18 to 99 years, with a mean age of 43.28 and a median of 41. This broad age distribution indicates a diverse customer base, which enhances the generalizability and relevance of the insights derived from the analysis.

In addition to review text and reviewer age, the dataset included numerical ratings on a scale of 1 to 5, representing customer satisfaction levels. It also categorized reviewed products into various clothing types such as dresses, tops, bottoms, and other apparel categories. These features—age, rating, and product category—were incorporated as contextual metadata to support downstream tasks such as clustering and topic modeling, providing meaningful dimensions for segmenting the data and extracting interpretable themes from the reviews. Additional visualizations (e.g., age and rating distributions) are included in Appendix A.2.

### 3.3 Feature Engineering Results

Feature engineering yielded

- **Word Counts:** Average review length was 59.35 words (rev\_word\_count, min 2, max 115), with 33.17 unique words (unique\_word\_count, min 2, max 62).
- **Sentiment Scores:** VADER compound scores averaged 0.673 (min -0.9735, max 0.9923), reflecting generally positive sentiment. Positive, negative, and neutral scores were also computed.

Detailed statistics are provided in Appendix A.1.

### 3.4 Clustering Results

K-means clustered 22,641 reviews into three balanced segments (Table 1), visualized via t-SNE (Figure 3):

- **Cluster 0:** 33.06% (7,484 reviews), mixed sentiment, focused on fit issues (e.g., “Dress too tight”).
- **Cluster 1:** 34.29% (7,763 reviews), positive sentiment, style-driven (e.g., “Stylish and comfortable”).
- **Cluster 2:** 32.66% (7,394 reviews), neutral sentiment, quality concerns (e.g., “Material faded quickly”).

The elbow method confirmed  $K = 3$  as optimal, with the WCSS plot (Figure 3) showing a clear elbow at  $K = 3$ . Silhouette scores (Table 2) indicated good cluster separation for  $K = 3$  (0.5301), with higher  $K$  values yielding lower scores, suggesting reduced cohesion. Additional cluster analysis visualizations are provided in Appendix A.2.

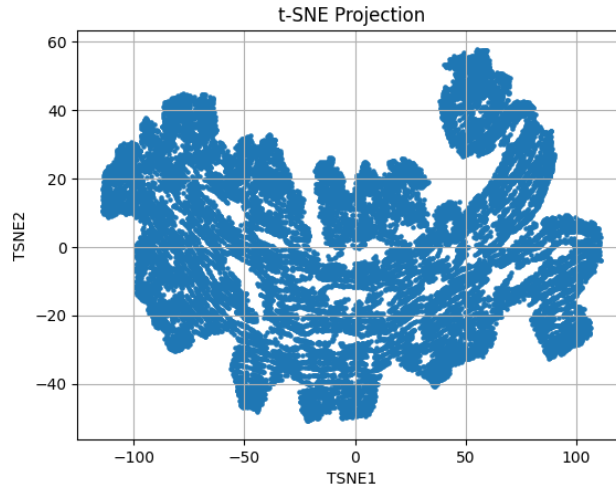


Figure 2: Unlabeled t-SNE

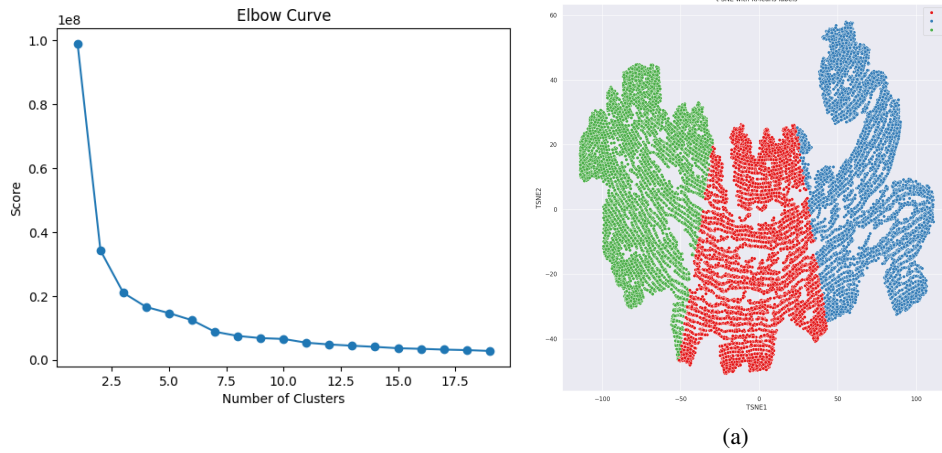


Figure 3: (a) Elbow Method Plot for K-means Clustering (b) t-SNE Visualization with K-means Clusters

### 151 3.5 Topic Modeling Results

152 The LDA model identified four interpretable topics, derived from the review text using Gensim:

- 153 1. **Dresses and Skirts:** Words like “dress,” “size,” “skirt,” “fabric,” “petite,” “waist,” “beautiful”  
154 (e.g., “The dress fits perfectly but is a bit short for petites”).
- 155 2. **Jeans and Pants:** Words like “jean,” “pant,” “size,” “great,” “jacket,” “comfortable,” “pair”  
156 (e.g., “These jeans are super comfortable and fit well”).
- 157 3. **Sweaters and Colors:** Words like “color,” “great,” “sweater,” “soft,” “love,” “nice,” “black”  
158 (e.g., “Love the soft sweater, especially in black”).
- 159 4. **Sizing and Shirts:** Words like “size,” “small,” “fabric,” “large,” “medium,” “shirt,” “didn’t,”  
160 “cute” (e.g., “The shirt runs small, but the fabric is cute”).

### 161 3.6 Evaluation

162 LDA performance was assessed with:

- 163 • **Perplexity:** 69.56, indicating reasonable model fit to the data, though lower values would  
164 suggest better generalization.
- 165 • **Jensen-Shannon Divergence (JSD):** Pairwise JSD between topics ranged from 0.2284 to  
166 0.3116, showing moderate topic distinctiveness, computed as:

$$JSD(P||Q) = \frac{1}{2}(KL(P||M) + KL(Q||M)),$$

167 where  $M = \frac{P+Q}{2}$ . Specific values include:

- 168 – Topic 0 vs. 1: 0.2692
- 169 – Topic 0 vs. 2: 0.2747
- 170 – Topic 0 vs. 3: 0.2284
- 171 – Topic 1 vs. 2: 0.2814
- 172 – Topic 1 vs. 3: 0.3116
- 173 – Topic 2 vs. 3: 0.2974

174 Clustering performance was evaluated with:

- 175 • **Silhouette Score:** 0.5301 for  $K = 3$ , indicating good cluster separation, with scores for  
176 higher  $K$  values (Table 2) showing reduced cohesion.

177 Qualitative analysis via manual review sampling and pyLDAvis visualizations confirmed topic  
178 relevance (e.g., “Dresses and Skirts” topic aligned with petite sizing complaints). No coherence  
179 metrics (e.g.,  $C_v$ ) were computed due to computational constraints, relying on manual inspection for  
180 topic selection.

### 181 3.7 Analysis

182 The clustering analysis identified three distinct customer segments: Cluster 0 (fit-focused) emphasized  
183 garment sizing and comfort, Cluster 1 (style-driven) focused on design and trends, and Cluster 2  
184 (quality-concerned) highlighted fabric durability and construction. These segments inform targeted  
185 marketing strategies, such as promoting fit accuracy to Cluster 0 or emphasizing aesthetics to Cluster  
186 1.

187 Topic modeling with LDA uncovered key themes: sizing issues in dresses and skirts (especially for  
188 petite customers), comfort concerns for jeans and pants, visual and tactile preferences for sweaters,  
189 and inconsistent sizing for shirts. These insights guide product improvements, like refining size  
190 guides and material descriptions.

191 Sentiment analysis using VADER revealed a positive overall sentiment (mean score: 0.673), reflecting  
192 general customer satisfaction.

193 Analytical challenges included noisy text, which was resolved through NLTK-based preprocessing  
194 and stopword removal. Topic clarity was enhanced by focusing on nouns and adjectives during tok-  
195 enization. Computational constraints were mitigated using checkpointing and sentence embeddings  
196 to improve efficiency.

## 197 4 Conclusion

198 The project applied Latent Dirichlet Allocation (LDA) topic modeling and K-means clustering to  
199 22,641 women’s clothing reviews to extract meaningful insights. LDA uncovered four dominant  
200 themes aligned with product categories and customer concerns: dresses and skirts; jeans and pants;  
201 sweaters and color preferences; and sizing issues with shirts. Meanwhile, K-means grouped reviews  
202 into three distinct customer segments by primary concern: fit-focused, style-driven, and quality-  
203 conscious shoppers. These complementary results linked what customers discuss with who is  
204 discussing it, offering a comprehensive view of customer feedback. The analysis surpassed its initial  
205 objectives by delivering robust findings and visualizations beyond preliminary goals. The approach is  
206 also modular and scalable, enabling updates with new data and adjustments to the number of topics  
207 or clusters as needed. Overall, the study demonstrates that advanced text mining techniques like  
208 LDA and K-means yield actionable insights from unstructured reviews, providing a roadmap for  
209 data-driven improvements in product design, sizing policies, and targeted marketing in the women’s  
210 apparel industry, ultimately enhancing customer satisfaction and business outcomes.

## 211 5 References

- 212 1. Dataset: [https://www.kaggle.com/datasets/nicapotato/](https://www.kaggle.com/datasets/nicapotato/womens-ecommerce-clothing-reviews/data)  
213 [womens-ecommerce-clothing-reviews/data](https://www.kaggle.com/datasets/nicapotato/womens-ecommerce-clothing-reviews/data)
- 214 2. LDA for Sentiment Analysis: [https://datascience.stackexchange.com/](https://datascience.stackexchange.com/questions/53271/lda-for-sentiment-analysis)  
215 [questions/53271/lda-for-sentiment-analysis](https://datascience.stackexchange.com/questions/53271/lda-for-sentiment-analysis)
- 216 3. LDA in NLP: [https://mohamedbakrey094.medium.com/](https://mohamedbakrey094.medium.com/all-about-latent-dirichlet-allocation-lda-in-nlp-6cfa7825034e)  
217 [all-about-latent-dirichlet-allocation-lda-in-nlp-6cfa7825034e](https://mohamedbakrey094.medium.com/all-about-latent-dirichlet-allocation-lda-in-nlp-6cfa7825034e)
- 218 4. Understanding LDA: [https://medium.com/@pinakdatta/](https://medium.com/@pinakdatta/understanding-lda-unveiling-hidden-topics-in-text-data-9bbbd25ae162)  
219 [understanding-lda-unveiling-hidden-topics-in-text-data-9bbbd25ae162](https://medium.com/@pinakdatta/understanding-lda-unveiling-hidden-topics-in-text-data-9bbbd25ae162)
- 220 5. Visualize Document Clusters: [https://www.mathworks.com/help/textanalytics/](https://www.mathworks.com/help/textanalytics/ug/visualize-document-clusters-using-lda-model.html)  
221 [ug/visualize-document-clusters-using-lda-model.html](https://www.mathworks.com/help/textanalytics/ug/visualize-document-clusters-using-lda-model.html)
- 222 6. Clustering with LDA: [https://stats.stackexchange.com/questions/292281/](https://stats.stackexchange.com/questions/292281/clustering-with-latent-dirichlet-allocation-lda-distance-measure)  
223 [clustering-with-latent-dirichlet-allocation-lda-distance-measure](https://stats.stackexchange.com/questions/292281/clustering-with-latent-dirichlet-allocation-lda-distance-measure)
- 224 7. Clustering and Topic Modeling Evaluation: [https://pmc.ncbi.nlm.nih.gov/](https://pmc.ncbi.nlm.nih.gov/articles/PMC9040385/)  
225 [articles/PMC9040385/](https://pmc.ncbi.nlm.nih.gov/articles/PMC9040385/)
- 226 8. Age-wise Sentiment Analysis: [https://www.kaggle.com/code/adhok93/](https://www.kaggle.com/code/adhok93/understanding-age-wise-sentiments-using-k-means)  
227 [understanding-age-wise-sentiments-using-k-means](https://www.kaggle.com/code/adhok93/understanding-age-wise-sentiments-using-k-means)
- 228 9. COVID-19 Literature Clustering: [https://www.kaggle.com/code/maksimeren/](https://www.kaggle.com/code/maksimeren/covid-19-literature-clustering/notebook)  
229 [covid-19-literature-clustering/notebook](https://www.kaggle.com/code/maksimeren/covid-19-literature-clustering/notebook)
- 230 10. Topic Modeling Tutorial: [https://github.com/adashofdata/](https://github.com/adashofdata/nlp-in-python-tutorial/blob/master/4-Topic-Modeling.ipynb)  
231 [nlp-in-python-tutorial/blob/master/4-Topic-Modeling.ipynb](https://github.com/adashofdata/nlp-in-python-tutorial/blob/master/4-Topic-Modeling.ipynb)
- 232 11. t-SNE Guide: <https://distill.pub/2016/misread-tsne/>
- 233 12. Hutto, C.J. & Gilbert, E.E. (2014). VADER: A Parsimonious Rule-based Model for Senti-  
234 ment Analysis. ICWSM-14, Ann Arbor, MI.
- 235 13. Lossio-Ventura, J.A., et al. (2021). Evaluation of Clustering and Topic Modeling Methods.  
236 Artificial Intelligence in Medicine, 117:102096. [https://doi.org/10.1016/j.artmed.](https://doi.org/10.1016/j.artmed.2021.102096)  
237 [2021.102096](https://doi.org/10.1016/j.artmed.2021.102096)
- 238 14. ChatGPT Prompts:
  - 239 • Create LaTeX table template for dataset summary.
  - 240 • Provide Python code for LDA topic modeling.
  - 241 • <https://chatgpt.com/share/67da5ad2-b7f4-800d-8751-bd44a5c92b87>

- <https://chatgpt.com/share/6803110c-a440-800d-8f3d-d7249d6d0ed5>
- <https://chatgpt.com/share/6803112a-edb0-800d-be61-91f218d71557>
- Explain LDA and topic modeling in detail.
- List potential challenges in NLP projects.
- Debug Python code for LDA implementation.

## A Supplementary Material

This appendix provides additional figures and feature engineering details to support the main report.

### A.1 Feature Engineering Details

Table 3 summarizes the statistics for engineered features, including word counts and VADER sentiment scores, computed for the 22,641 reviews.

### A.2 Additional Figures

This section includes supplementary visualizations for data exploration, clustering, and topic modeling.

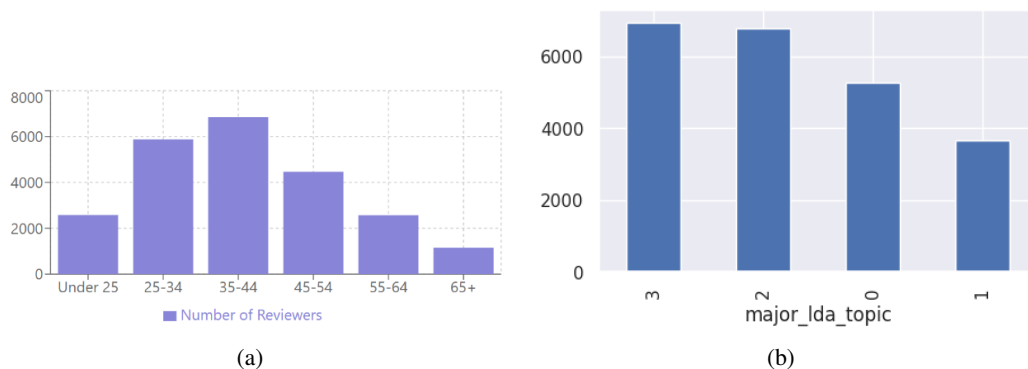


Figure 4: (a) Age Distribution of Reviewers (b) LDA Topic Value Counts

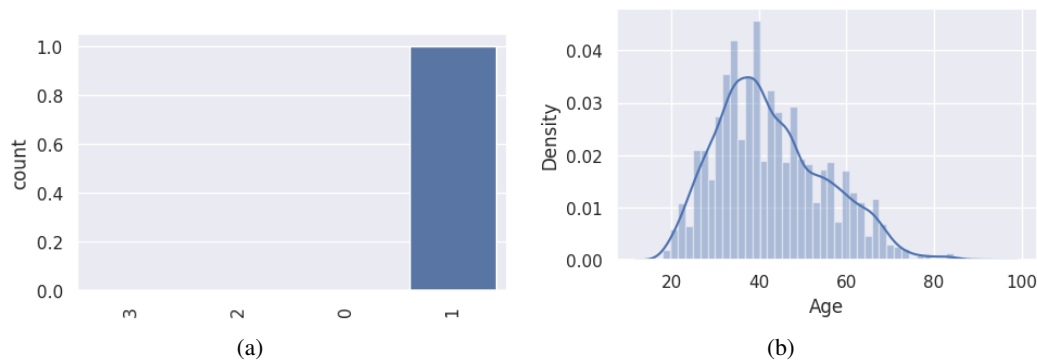


Figure 5: (a) Categorical Feature Distribution for Cluster 0 (b) Numerical Feature Distribution for Cluster 0



Table 1: Cluster Distribution Analysis ( $K = 3$ )

Cluster	Number of Points	Percentage (%)
Cluster 0	7,484	33.06
Cluster 1	7,763	34.29
Cluster 2	7,394	32.66

Total reviews analyzed: 22,641

Table 2: Silhouette Scores for Different  $K$  Values

$K$	Silhouette Score
3	0.5301
4	0.4545
5	0.3890
6	0.4019
7	0.4120
8	0.4084
9	0.4236
10	0.4241
11	0.4089

Table 3: Feature Engineering Statistics

Feature	Mean	Std	Min	Median	Max
rev_word_count	59.35	28.28	2	57	115
unique_word_count	33.17	14.08	2	32	62
Compound Sentiment	0.673	0.390	-0.9735	0.8481	0.9923
Positive Sentiment	0.279	0.165	0.0	0.260	1.0
Negative Sentiment	0.048	0.089	0.0	0.0	0.737
Neutral Sentiment	0.673	0.180	0.0	0.677	1.0

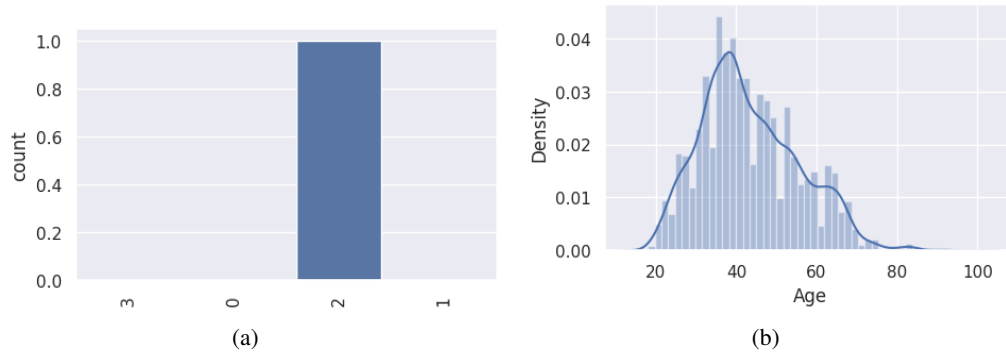


Figure 6: (a) Categorical Feature Distribution for Cluster 1 (b) Numerical Feature Distribution for Cluster 1

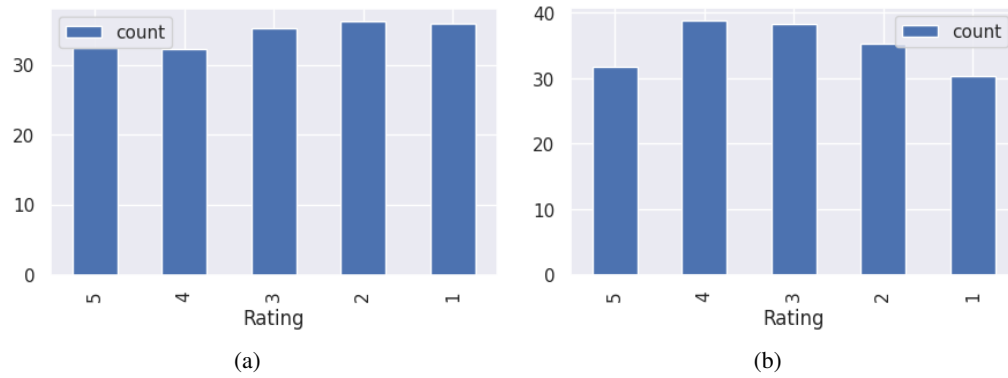


Figure 7: (a) Rating Distribution for Cluster 0 (b) Rating Distribution for Cluster 1

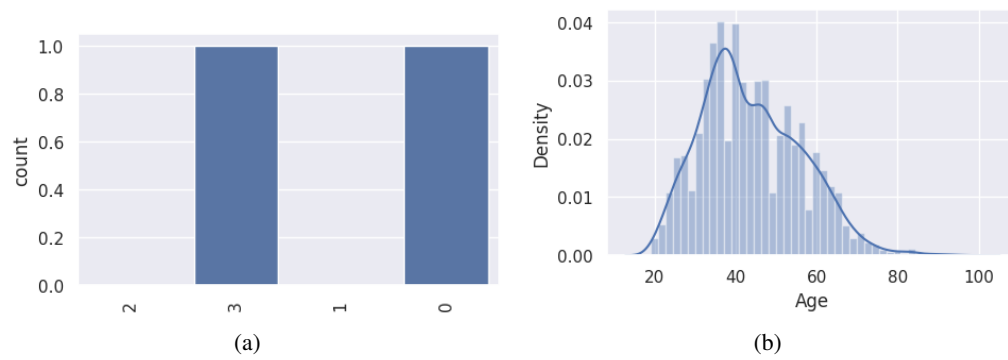


Figure 8: (a) Categorical Feature Distribution for Cluster 2 (b) Numerical Feature Distribution for Cluster 2

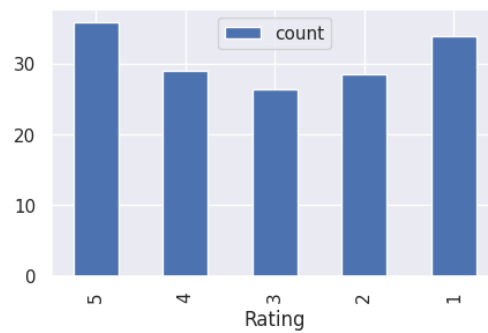


Figure 9: Rating Distribution for Cluster 2

Topic	Keyword	Probability
0	dress	0.1312
	size	0.0353
	true	0.0060
1	jean	0.0397
	pant	0.0359
	work	0.0077
2	color	0.0458
	great	0.0301
	little	0.0076
3	size	0.0390
	small	0.0333
	side	0.0067

Table 4: Top 3 Keywords per Topic with Probabilities