

Enhanced Sentiment Interpretation via a Lexicon-Fuzzy-Transformer Framework

Shayan Rokhva, Mousa Alizadeh*, Maryam Abdollahi Shamami

Department of Information Technology, Faculty of Industrial & Systems Engineering, Tarbiat Modares University, Tehran, Iran

School of Engineering, Royal Melbourne Institute of Technology (RMIT University), Melbourne, VIC 3001, Australia

Department of Information Technology, Faculty of Industrial & Systems Engineering, Tarbiat Modares University, Tehran, Iran

Abstract

Accurately detecting sentiment polarity and intensity in product reviews and social media posts remains challenging due to informal and domain-specific language. To address this, we propose a novel hybrid lexicon-fuzzy-transformer framework that combines rule-based heuristics, contextual deep learning, and fuzzy logic to generate continuous sentiment scores reflecting both polarity and strength. The pipeline begins with VADER-based initial sentiment estimations, which are refined through a two-stage adjustment process. This involves leveraging confidence scores from DistilBERT—a lightweight transformer—and applying fuzzy logic principles to mitigate excessive neutrality bias and enhance granularity. A custom fuzzy inference system then maps the refined scores onto a 0–1 continuum, producing expert-like judgments. The framework is rigorously evaluated on four domain-specific datasets: food delivery, e-commerce, tourism, and fashion. Results show improved alignment with user ratings, better identification of sentiment extremes, and reduced misclassifications. Both quantitative metrics (distributional alignment, confusion matrices) and qualitative insights (case studies, runtime analysis) affirm the model’s robustness and efficiency. This work demonstrates the value of integrating symbolic reasoning with neural models for interpretable, fine-grained sentiment analysis in linguistically dynamic domains.

Keywords:

Sentiment Analysis, Fuzzy Logic, Fuzzy Inference System, Transformers, Product Reviews, Deep Learning

1. Introduction

In the age of pervasive digital communication, sentiment analysis (SA) has become a critical tool for interpreting public opinion embedded in product reviews, social media content, and online discourse. By computationally extracting and evaluating emotional cues from textual data, SA provides actionable insights into consumer preferences, product performance, and brand perception [1],[2],[3]. The recent evolution of SA has been significantly driven by two complementary paradigms: deep learning (DL) and fuzzy logic. DL models, particularly those based on neural networks, excel at uncovering complex linguistic structures within large-scale datasets, enabling robust sentiment interpretation across diverse domains. In contrast, fuzzy logic provides a mathematical framework for modeling the ambiguity and imprecision intrinsic to natural language, thus supporting nuanced, human-like reasoning [4],[5],[6]. The integration of these approaches has led to the development of intelligent and adaptive systems capable of both high accuracy and interpretability, reflecting a broader trend toward hybrid models that balance data-driven performance with symbolic reasoning[7].

The utility of fuzzy logic in SA stems from its ability to represent sentiment gradation through membership functions[8]. For example,[9] mapped Support Vector Machine (SVM) output to seven triangular fuzzy sets and employed a fuzzy inference system (FIS) to detect nuanced sentiment intensity, outperforming traditional classifiers on Twitter data. Similarly,[10, 11] applied Gaussian fuzzy sets within a Bidirectional Long Short-Term Memory (BiLSTM) framework to process Saudi dialect tweets and achieved 86% accuracy in binary sentiment classification for e-learning feedback. In decision-making contexts, [12, 13] modeled buyer preferences under

uncertainty by transforming review tags into picture fuzzy values and computing entropy-based weights for ranking new energy vehicles. Other efforts, such as that by [14], combined semantic word pairs with SentiWordNet (SWN) polarity scores and used fuzzy rules for sentiment extraction in transportation data. [15] developed an unsupervised Mamdani FIS utilizing polarity lexicons to classify tweets without requiring labeled data. In the financial domain [4]

used GloVe embeddings and TF-IDF-based polarity seeding to create a five-level fuzzy lexicon for nuanced analysis of economic news.

DL has also played a transformative role in SA by capturing rich semantic and contextual patterns from raw text. [16] Fine-tuned DistilBERT for analyzing long corporate reviews while using VADER for shorter posts, fusing both outputs to build a real-time Corporate Sentiment Index. [17] demonstrated that GPT-4o, applied in a zero-shot setting with aspect-specific prompts, achieved near-human performance on hotel reviews. [18] collected 9000 TikTok reviews, labeled them using VADER, and trained a one-vs-one SVM model that achieved an F1 score of 0.80, without requiring manual annotation. Hybrid approaches have also gained traction; for example, [19] integrated fuzzy alpha-cuts with longitudinal trend analysis for brand monitoring. [20, 21] benchmarked multiple SA tools across heterogeneous datasets, revealing inconsistencies in classification accuracy between lexicon-based, machine learning, and hybrid models.

Despite these advances, both fuzzy logic and DL-based SA methods exhibit limitations. Fuzzy systems, while interpretable, often simplify sentiment into binary or ternary categories, losing the richness of sentiment expression [10],[15]. In many cases, intensity cues are static or rule-based, rather than dynamically learned or context-aware [14],[12].

Lexicon-based systems like VADER and SWN, though fast and transparent, are prone to misclassifying informal or domain-specific expressions, such as slang or sarcasm, that frequently appear in user-generated content [22],[23],[18]. Moreover, the constantly evolving informal language prevalent among teenagers and youngsters makes fine-tuning lexicons even more challenging [24][25][26]. Conversely, DL models, including DistilBERT and GPT-based approaches, often produce rigid or overly polarized sentiment classifications when confronted with ambiguous or nuanced inputs [17],[27],[16].

Moreover, many existing pipelines fail to address computational complexity or runtime performance [28, 29],[12], and a few explicitly incorporate multi-stage correction mechanisms or cognitively inspired reasoning to align computed sentiment with human interpretation.

These shortcomings underscore the need for a unified, interpretable, and intensity-aware SA framework that bridges the gap between symbolic reasoning and contextual depth while maintaining efficiency and adaptability to informal language. To support this analysis, Table 1 presents a structured review of recent studies, covering their domain applications, methodological approaches, and observed strengths and limitations.

This study addresses these gaps by proposing a hybrid end-to-end SA framework that integrates the efficiency of rule-based scoring (VADER), the contextual learning capabilities of a pre-trained light-weight transformer (DistilBERT), and the interpretability of a fuzzy logic-based inference system. The pipeline operates in three stages: (1) initial sentiment scoring using VADER, (2) refinement using Transformer-based insights and fuzzy logic perspective, and (3) aggregation through an FIS that outputs continuous sentiment scores, capturing both polarity and intensity, mean- while mimicking an expert-like decision-making process.

The main contributions of this study are as follows:

- Introducing an integrated SA framework that merges rapid rule-based VADER scoring, contextual understanding from a pre-trained light Trans- former model (DistilBERT), and fuzzy logic-based refinements to improve SA within informal textual content commonly found in product reviews and digital media.
- Development of a tailored FIS capable of generating continuous sentiment scores, improving expressiveness and interpretability compared to discrete classifiers.
- Conducting empirical evaluations on four real-world benchmark datasets, showcasing a reduction in neutrality bias and superior cor- correspondence with human-annotated sentiments, all achieved without dependence on supervised learning processes.

The remainder of the study is structured as follows: Section 2 details the proposed methodology, in both segmented

and integrated formats, Section 3 evaluates the methodology on four diverse datasets, proposing experimental results and discussion; and Section 4 concludes the paper and outlines potential directions for future work.

2. Proposed Method

This section presents the end-to-end sentiment scoring framework proposed under three subsections, and they are delineated below. Section 2.1 formalizes the problem and delineates each framework component, specifying its role within the pipeline. Section 2.2 details the datasets utilized for model evaluation across diverse domains. Section 2.3 describes the integration strategy, illustrating how the components interact through a two-stage refinement process and an FIS to enhance sentiment estimation, mitigate neutrality bias, and improve the alignment between computed sentiment scores and user perceptions.

2.1. Problem Formulation

2.1.1. Objective

This study aims to mitigate the neutrality bias inherent in traditional rule-based SA tools, thereby enhancing their ability to detect both sentiment polarity and intensity, particularly in domain-specific, highly informal textual data.

2.1.2. VADER

The VADER is a widely used rule-based SA tool, valued for its speed and interpretability. However, its effectiveness diminishes when applied to industry-specific informal language, where nuanced or implicit sentiments are common. In this framework, VADER is used to generate initial, normalized sentiment scores, positive, negative, and neutral, which are then refined to enhance accuracy.

Table 1: A Structured Review of Recent Studies

Paper	Objective	Domain	Proposed Method	Pros & Cons
[28]	Transfer concept polarity	Cross-domain	Triangular fuzzy functions; graph propagation; trapezoidal refinement; fuzzy-mean aggregation	+ Intensity retained; Domain adaptation - Large semantic graph; Unreported runtime
[30]	Generalize SA to new domains	Amazon reviews	Fuzzy feature extraction, trapezoidal membership refinement, domain-overlap sentiment inference	+ Fuzzy intensity; Cross-domain transfer - Five-hour modeling
[12]	Rank NEVs by sentiment	Auto reviews	Picture fuzzy conversion; Entropy-based attribute weights; Regret-theory NEV ranking model	+ Uncertainty modelling - Slow processing; Non-interactive PFN; Intensity ignored
[15]	Model tweet sentiment intensity	Twitter	Extract features; Fuzzy profiles; SVM/MaxEnt Classification; fuzzy inference	+ Granular scoring; ML-fuzzy combo - Small corpus; Moderate performance; Timing absent
[14]	Infer transport sentiment	Smart city	Feature extraction via fuzzy ontology; SWRL rules infer five-level polarity	+ 5-level output; 96% accuracy - Reasoning cost unknown; Domain specific
[19]	Track brand sentiment trends	Brand data	Fuzzy α -cuts quantifying sentiment; longitudinal aggregation	+ Adjustable granularity; Longitudinal insight - API reliance; No benchmarks
[10]	Process dialect e-learning tweets	E-learning	Gaussian fuzzification; BiLSTM processing; defuzzification; sigmoid classification	+ Dialect robustness; Internal intensity - Only binary output; Small dataset
[17]	Aspect-based GPT sentiment	Hotels	GPT-4.0 w/ prompts + similarity validation	+ 95% human match; Zero-shot - No intensity scale; API cost
[31]	Extract film review sentiment	Movies	GAN-filtered Gaussian SVD + fuzzy rules	+ Intensity levels; Spam robust - No runtime; Small dataset
[9]	Unsupervised fuzzy tweet SA	Twitter	Lexicon scoring + Mamdani rules on 3-level scores	+ Label-free; Intensity encoded - Limited gradation; Slow on

[32]	Explore feedback	tourism	TripAdvisor	VADER + LDA w/ year-wise trends	SWN + 18-year span; Label-free; VADER-LDA synergy - Outputs' intensity flattened; No runtime
[16]	Real-time sentiment index	corp.	Corporate reputation	VADER + DistilBERT fusion; LDA scaling	+ Rule-transformer fusion; Continuous index - No label; Fuzzy logic absent; Empirical weights; No timing
[20]	Benchmark toolsets	SA	SA tools	Compare 13 packages on 7 datasets	+ Captures intensity-polarity; Timing reported - Deep Learning excluded; Some packages are slow
[33]	Baboon sentiment classifier	Baboon corpus	Baboon corpus	Unsupervised lexicon + 6 polarity-shifting rules (negations, intensifiers, etc.)	+ Interpretable rules; Pioneering - Handcrafted lexicon; Language-specific rules; Limited corpus
[4]	Fuzzy lexicon for stock news	Finance		TF-IDF polarities diffused via QLogic five-level fuzzy lexicon	+ Context-specific intensity; No labels - Heavy compute; Bigger corpus needed
[18]	Classify sentiment	TikTok	TikTok	VADER auto-labels; one-vs-one SVM classifies TikTok sentiment	+ Good performance - Fuzzy absent; No intensity; No timing
[34]	Mitigate sparsity via hybrid SA-CB-CF	Groceries	Ho-tels	GloVe-BiGRU SA, BERT-AE-MLP CB, SVD CF, RF fusion	+ BERT deeper comprehension-RMSE reduction - Intensity ignored, NO run-

2.1.3. Predefined Transformers

To address VADER's limitations in capturing nuanced sentiment polarity, a pre-trained Transformer model, *DistilBERT*-base-uncased, is integrated into the workflow. *DistilBERT*, a compressed variant of BERT, maintains approximately 97% of BERT's performance while reducing model size and inference time by approximately 40%. Fine-tuned on the Stanford Sentiment Treebank v2 (SST-2) dataset, the Transformer outputs probability-based confidence scores for positive and negative sentiment classifications [35, 36].

These confidence scores are used to refine VADER's initial outputs by targeting the reallocation of the neutral component. Specifically, when the Transformer's confidence in a positive (*po*) or negative (*ne*) label exceeds 90%, 30% of the neutral score is shifted toward the corresponding sentiment. For confidence levels between 80%–90% and 70%–80%, the reallocation percentages are adjusted to 20% and 10%, respectively. This recalibration process is formalized in Eq.1 and Eq.2, where (P_{o1} , Ne_1 , Nu_1) represent the original VADER scores, and (P_{o2} , Ne_2 , Nu_2) denote the refined sentiment scores.

If $Nu_1 = 1$, $Ne_1 = 0$, $Po_1 = 0$:

$$\begin{cases} Nu_2 = 0.7, Ne_2 = 0, Po_2 = 0.3 & \text{if } 0.9 \leq C_{po} < 1.0 \\ Nu_2 = 0.8, Ne_2 = 0, Po_2 = 0.2 & \text{if } 0.8 \leq C_{po} < 0.9 \\ Nu_2 = 0.9, Ne_2 = 0, Po_2 = 0.1 & \text{if } 0.7 \leq C_{po} < 0.8 \end{cases} \quad (1)$$

If $Nu_1 = 1$, $Ne_1 = 0$, $Po_1 = 0$:

$$\begin{cases} Nu_2 = 0.7, Ne_2 = 0.3, Po_2 = 0 & \text{if } 0.9 \leq C_{pe} < 1.0 \\ Nu_2 = 0.8, Ne_2 = 0.2, Po_2 = 0 & \text{if } 0.8 \leq C_{pe} < 0.9 \\ Nu_2 = 0.9, Ne_2 = 0.1, Po_2 = 0 & \text{if } 0.7 \leq C_{pe} < 0.8 \end{cases} \quad (2)$$

2.1.4. Fuzzy Logic Applications

Fuzzy logic, which is inherently adept at modeling gradations of meaning, is employed in the second refinement

phase to amplify the sentiment intensity that remains underrepresented after transformer-based corrections. While the initial refinement phase addresses cases of complete misclassification into the neutral category, fuzzy logic deals with subtle sentiment that is present yet insufficiently expressed. Specifically, when the Transformer's confidence score for either the positive or negative category exceeds 0.95, the corresponding sentiment score is intensified through a square-root transformation. This operation aligns with the fuzzy logic concept of applying linguistic hedges, such as elevating "happy" to "very happy", to enhance membership strength. Following the transformation, all sentiment scores are re-normalized to ensure they collectively sum to one.

This adjustment process is mathematically defined in Eqs 3 and 4, where Po_3 , Ne_3 , and Nu_3 represent the final sentiment scores after this phase.

If $C_{po} \geq 0.95$, then:

$$\begin{aligned} Po_3 &= \frac{\sqrt{Po_2}}{\sqrt{Po_2} + Ne_2 + Nu_2} \\ Ne_3 &= \frac{Ne_2}{\sqrt{Po_2} + Ne_2 + Nu_2} \\ Nu_3 &= \frac{Nu_2}{\sqrt{Po_2} + Ne_2 + Nu_2} \end{aligned} \quad (3)$$

If $C_{ne} \geq 0.95$, then:

$$\begin{aligned} Ne_3 &= \frac{\sqrt{Ne_2}}{\sqrt{Po_2} + Ne_2 + Nu_2} \\ Po_3 &= \frac{Po_2}{\sqrt{Po_2} + Ne_2 + Nu_2} \\ Nu_3 &= \frac{Nu_2}{\sqrt{Po_2} + Ne_2 + Nu_2} \end{aligned} \quad (4)$$

2.1.5. Fuzzy Inference System

The FIS is employed to interpret the sentiment scores refined through the two-step adjustment process. Its objective is to emulate expert-like reasoning by producing a continuous sentiment score between 0 and 1, thereby capturing both sentiment polarity and intensity with greater fidelity. The process begins with fuzzification, wherein each refined sentiment score, positive, negative, or neutral, is mapped onto corresponding linguistic intensity levels: "low," "medium," or "high." The membership functions and associated boundaries governing this transformation are formally defined in Eqs. 5–7. Through this structure, the FIS enables the modeling of nuanced sentiment expressions beyond simple categorical classification.

A comprehensive and customized knowledge-based system was developed, with its rules detailed in Table 2. The system is carefully designed to ensure continuity, consistency, and completeness across all possible input scenarios. Beyond its standardized formulation, the knowledge-based system includes deliberate adjustments aimed at mitigating the prevalent neutrality. Specifically, rows 8, 9, 16, 19, 25, and 26 address cases with dominant neutrality, but the presence of a medium- or high-intensity secondary sentiment prompts the system to favor positive or negative polarity instead. The membership function constructed for the final sentiment score—the framework's output variable—directly mirrors the intensity functions defined in Eqs. 5–7.

Specifically, the "Negative," "Neutral," and "Positive" sets employ mathematical forms and parameters identical to those used for the "Low," "Medium," and "High" input categories. The inference mechanism adopts the minimum operator as the t-norm for rule activation and the maximum operator as the s-norm for rule aggregation, as specified in Eq. 8. The final step involves defuzzification, where a continuous scalar value is derived using the centroid method, effectively capturing both the sentiment polarity and its intensity.

$$\mu_{\text{Low}}(x) = \begin{cases} 1 & x \leq 0.3 \\ \frac{0.5-x}{0.2} & 0.3 < x \leq 0.5 \\ 0 & x > 0.5 \end{cases} \quad (5)$$

$$\mu_{\text{Medium}}(x) = \begin{cases} 0 & x \leq 0.3 \\ \frac{x-0.3}{0.2} & 0.3 < x < 0.5 \\ \frac{0.7-x}{0.2} & 0.5 < x < 0.7 \\ 0 & x \geq 0.7 \end{cases} \quad (6)$$

$$\mu_{\text{High}}(x) = \begin{cases} 0 & x < 0.5 \\ \frac{x-0.5}{0.2} & 0.5 \leq x < 0.7 \\ 1 & x \geq 0.7 \end{cases} \quad (7)$$

$$\mu_{\text{output}}(z) = \max_{i=1 \dots 27} [\min(\mu_{\text{Neg}}(N_{\text{e}3}), \mu_{\text{Neu}}(N_{\text{u}3}), \mu_{\text{Pos}}(P_{\text{o}3}))] \quad (8)$$

$$\text{Final Sentiment Score} = \frac{\int z \mu_{\text{output}}(z) dz}{\int \mu_{\text{output}}(z) dz} \quad (9)$$



Figure 1: User Rating Distributions

Table 2 – Constructed Knowledge-based System

Rule	Negative	Neutral	Positive	Outcome
1	low	low	low	neutral
2	low	low	medium	positive
3	low	low	high	positive
4	low	medium	low	neutral
5	low	medium	medium	positive
6	low	medium	high	positive
7	low	high	low	neutral
8	low	high	medium	positive
9	low	high	high	positive
10	medium	low	low	negative
11	medium	low	medium	neutral
12	medium	low	high	positive
13	medium	medium	low	negative
14	medium	medium	medium	neutral
15	medium	medium	high	positive
16	medium	high	low	negative
17	medium	high	medium	neutral
18	medium	high	high	positive
19	high	low	low	negative
20	high	low	medium	negative
21	high	low	high	neutral
22	high	medium	low	negative
23	high	medium	medium	negative
24	high	medium	high	neutral
25	high	high	low	negative
26	high	high	medium	negative

2.2. Input Data

Four datasets from diverse application domains were utilized to evaluate the effectiveness and generalizability of the proposed framework. The **Food** dataset, compiled by the authors, comprises over 1,250 restaurant reviews collected in Tehran. The **TripAdvisor** dataset contains nearly 20,500 accommodation-related reviews in the tourism sector. The **Flipkart** dataset includes around 10,000 user reviews focused on electronic products, while the **Women’s E-Commerce** dataset encompasses approximately 23,500 comments pertaining to clothing and fashion items. Each dataset contains two essential features: (1) user-generated review texts, which are processed through the proposed SA pipeline, and (2) associated user ratings, serving as a reference for sentiment polarity and intensity. While the three external datasets use a conventional 1–5 star rating scale, the **Food** dataset follows a 2–5 star scale, with a portion of missing ratings primarily corresponding to reviews that express negative sentiment.

2.3. Strategy

The methodological flow illustrated in Figure 2 begins with VADER, which provides preliminary sentiment scores. However, due to its rule-based nature, these scores often exhibit neutrality bias and lack expressive granularity. To benchmark performance, these raw scores are first passed through the Fuzzy Inference System (FIS) to establish a comparative baseline. Concurrently, review texts are evaluated using the proposed Transformer model operating in inference mode, which outputs confidence probabilities for each sentiment class. When confidence in positive or negative sentiment exceeds predefined thresholds, a proportion of the neutral score is redistributed accordingly, as detailed in Section 2.1.3. This constitutes the first-stage correction, aimed at

mitigating erroneous neutrality.

The second-stage refinement, outlined in Section 2.1.4, targets sentiment under-expression. Using fuzzy logic and linguistic hedging, the framework intensifies sentiment magnitudes for high-confidence Transformer outputs, ensuring better representation of implicit emotional cues. The resulting refined scores are then processed by the custom FIS, which produces a continuous sentiment value in the $[0,1]$ range. This final output is compared to the unrefined baseline to assess the cumulative effect of both refinement stages. Although probability (from Transformers) and possibility (from fuzzy logic) are theoretically distinct, leveraging Transformer-derived confidences to guide fuzzy-based amplification is empirically sound. Prior Natural Language Processing (NLP) evidence confirms that texts rich in sentiment cues yield high confidence under both frameworks, justifying this hybrid integration despite the lack of a formal probabilistic–possibilistic mapping.

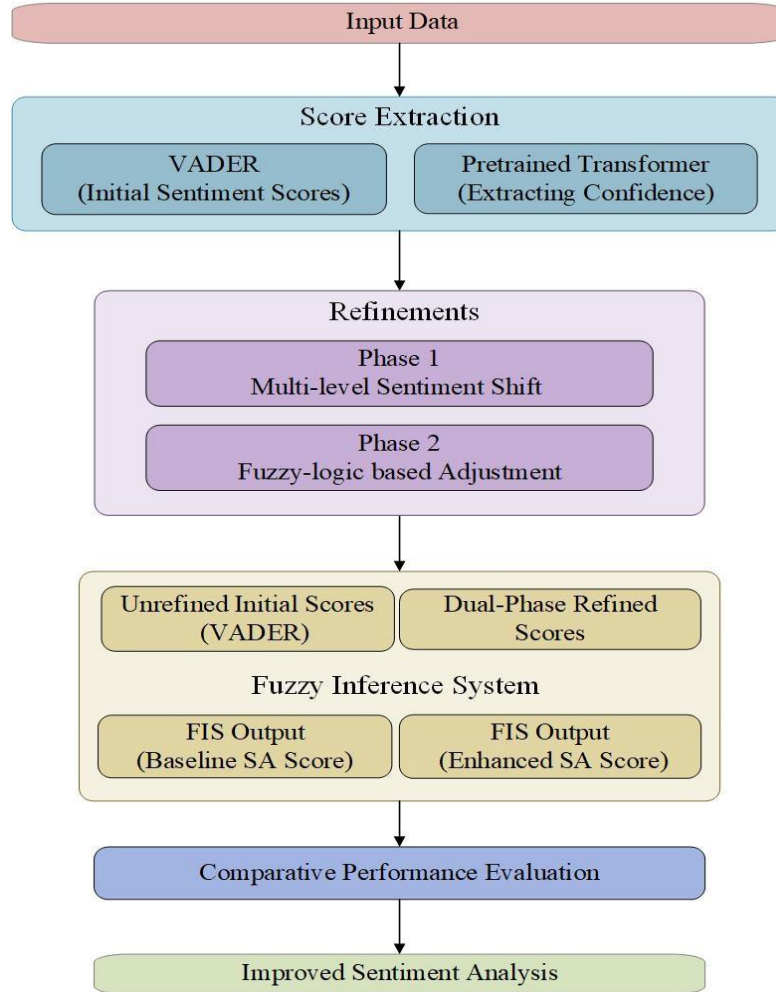


Figure 2: System workflow of the proposed SA framework.

3. Results & Discussion

The results are organized into a structured, multi-layered analysis. First, the framework’s impact on sentiment score distribution is visualized to reveal global improvements. This is followed by a statistical evaluation against user-assigned ratings to validate alignment. Next, a categorical partitioning of sentiment bands enables a finer-grained examination of the model’s interpretability across intensity levels. Confusion matrices are then utilized to identify desirable and undesirable sentiment transitions during refinement, offering insight into error dynamics. A qualitative review complements the quantitative analysis, showcasing how refined scores align with review content. Finally, performance metrics such as runtime, along with a critical reflection on the framework’s strengths, limitations, and future research opportunities, are presented to contextualize the system’s practicality and potential.

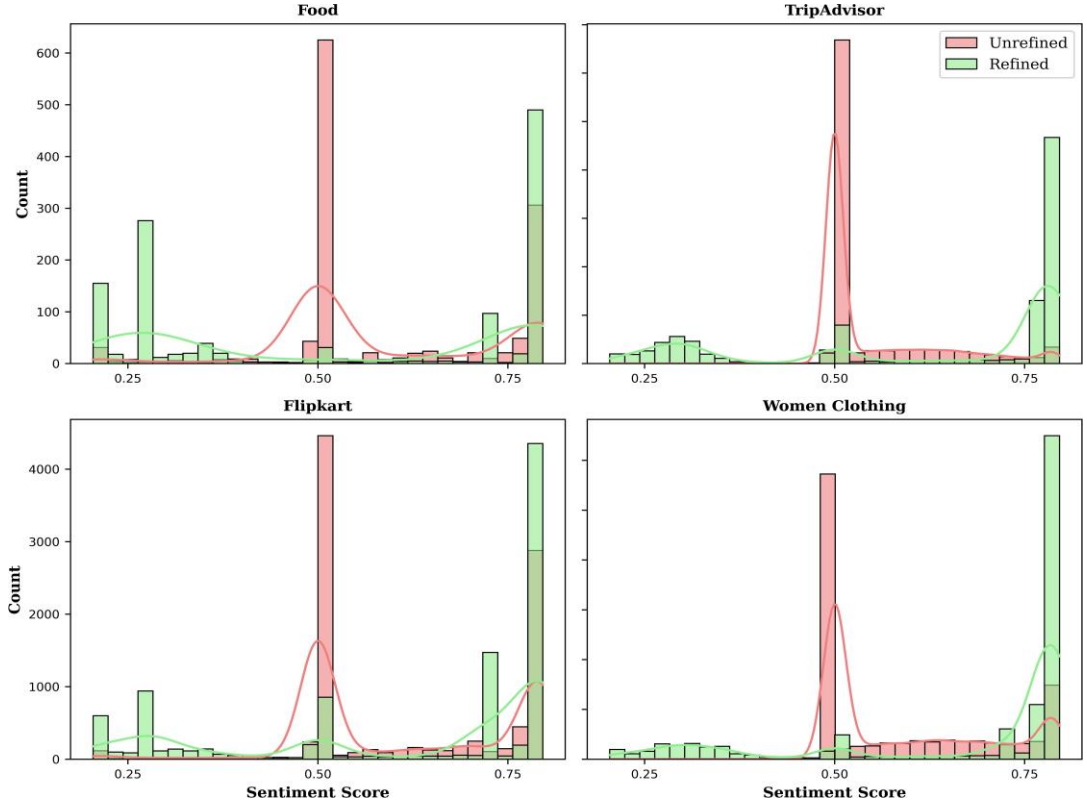


Figure 3: Distribution of sentiment scores before and after refinement across four datasets: Food, TripAdvisor, Flipkart, and Women’s Clothing. Based on the proposed framework, each subplot compares unrefined (red) vs. refined (green) sentiment outputs.

3.1. Global Distribution Analysis

As shown in **Figure 3**, the unrefined sentiment scores (pink) are heavily concentrated around the **midpoint**, revealing a pronounced neutrality bias and an inability to capture true sentiment polarity. This indicates that **VADER**, even when processed through the **FIS**, fails to accurately reflect **user sentiment**. In contrast, the **refined scores** (green) exhibit a well-balanced redistribution, concentrated predominantly in the **positive region**, followed by the **negative**, and only minimally around the neutral center. This distribution pattern aligns closely with the user-assigned ratings illustrated in Figure 1, confirming the effectiveness of the proposed refinement pipeline.

3.2. Rating-Stratified Distributional Evaluation

To assess the framework’s alignment with user perception, the final sentiment scores are evaluated against user-assigned star ratings. As illustrated in Figure 4, the unrefined scores display a strong central bias around 0.5, particularly in low-rating groups (1 and 2), where mean values seldom drop below the midpoint—signaling neutrality bias. After refinement, the distributions shift meaningfully: ratings 1 and 5 center around 0.30 and 0.70, respectively, accurately reflecting their sentiment polarity. Intermediate ratings (2 and 4) demonstrate consistent mean shifts, while the three-star group remains appropriately near the neutral zone.

Refinement also enhances the statistical profile. In the post-refinement state, box plots for extreme ratings fall entirely below or above the midpoint, with improved gradation observed in mid-tier groups. The reduction in outliers further indicates greater consistency and stronger alignment with the underlying rating structure. While Section 3.1 highlighted global improvements, this analysis demonstrates that the framework achieves finer-grained enhancements across rating-specific strata.

3.3. Band-wise Sentiment Intensity Alignment

Building on Section 3.2, this section provides a finer-grained analysis by categorizing continuous sentiment scores into five bands—Negative, Slightly Negative, Neutral, Slightly Positive, and Positive—as defined in Figure 5. This partitioning, grounded in the overlapping fuzzy regions (see Eqs. 5–7) and the five-level user rating system, is used solely for evaluation and does not influence the model’s internal computations. While the intermediate bands occur less frequently due to their narrow widths, this represents expected statistical behavior rather than a model weakness. Figure 6 illustrates sentiment category distributions across rating levels. The lower bars (unrefined) display a high concentration of gray, indicating neutrality bias, particularly within lower-rated groups. After refinement (upper bars), the distribution realigns with user sentiment: 5-star reviews are dominated by Positive tones (dark green), while 1-star and 2-star groups prominently feature Negative tones (red/orange), which were previously underrepresented. The 3-star group shows a balanced composition, consistent with its intermediate nature, and the 4-star group demonstrates a slight reduction in positivity, reflecting mild satisfaction.

A notable exception emerges in the Food dataset. Unlike other domains, the 3-star and 4-star reviews here exhibit a more critical tone—the 4-star distributions resemble 3-star patterns from other datasets, while 3-star reviews lean toward negativity, similar to 2-star ratings. Manual inspection confirmed that these reviews often emphasize flaws despite higher ratings, a nuance effectively captured by the model. To further evaluate this observation, the unrated portion (19%) of the Food dataset was analyzed separately. As discussed in Section 2.2, most of these reviews (75–80%) conveyed negative sentiment. Figure 7 reveals that while raw VADER outputs failed to detect this trend, the refined framework produced a negative sentiment distribution with minimal neutral presence, closely aligning with both content and visual inspection.

3.4. Error-Transition Profiling via Confusion Matrices

Figure 8 presents confusion matrices that visualize the redistribution of sentiment scores throughout the proposed pipeline, highlighting transitions categorized as desirable, moderately undesirable, or undesirable. While the unrefined scores exhibit a clear neutrality bias, strongly polarized instances—typically driven by explicit linguistic cues—are generally accurate and should remain stable post-refinement. Consequently, transitions between sentiment extremes (e.g., from positive to negative) are considered undesirable. These are marked in red (severely undesirable), pink (undesirable), and yellow (moderately undesirable) within the matrices. The analysis confirms that the proposed VADER–Fuzzy–Transformer pipeline effectively avoids severe polarity reversals (no red zones), with pink and yellow transitions being rare and limited to 0.5–1%. This demonstrates that the refinement process introduces minimal classification noise while preserving valid semantic shifts. Importantly, all datasets exhibit a consistent trend of movement from neutral to polarized classes, validating the system’s ability to overcome VADER’s neutrality bias. Moreover, positive-to-positive transitions significantly outnumber negative-to-negative ones across all matrices, reaffirming that the base system is inherently more effective in detecting positive sentiment. Overall, the confusion matrix analysis demonstrates that the proposed framework enhances sentiment resolution while maintaining a low error profile across diverse domains.

3.5. Case-Based Qualitative Validation

Table 3 presents selected comments from multiple datasets alongside the sentiment scores generated by the proposed framework. Due to length constraints, TripAdvisor entries are excluded. The results demonstrate the pipeline’s effectiveness in producing refined sentiment scores that accurately reflect user sentiment and correspond with assigned ratings. For instance, rows 1, 2, 6, 8, 9, and 12 illustrate how the system corrects neutrality bias by assigning appropriately negative scores. In contrast, rows 7, 11, and 14 showcase cases where positive sentiments are further intensified, enhancing alignment with user feedback. Nonetheless, certain limitations persist. In row 13, the refined score contradicts both the textual content and the user rating, indicating a misadjustment. Similarly, row 15 exhibits a near-correct outcome—the system assigns a moderately low score (0.32) for a negative 1-star review—yet a slightly lower value would have better captured the intensity of dissatisfaction.

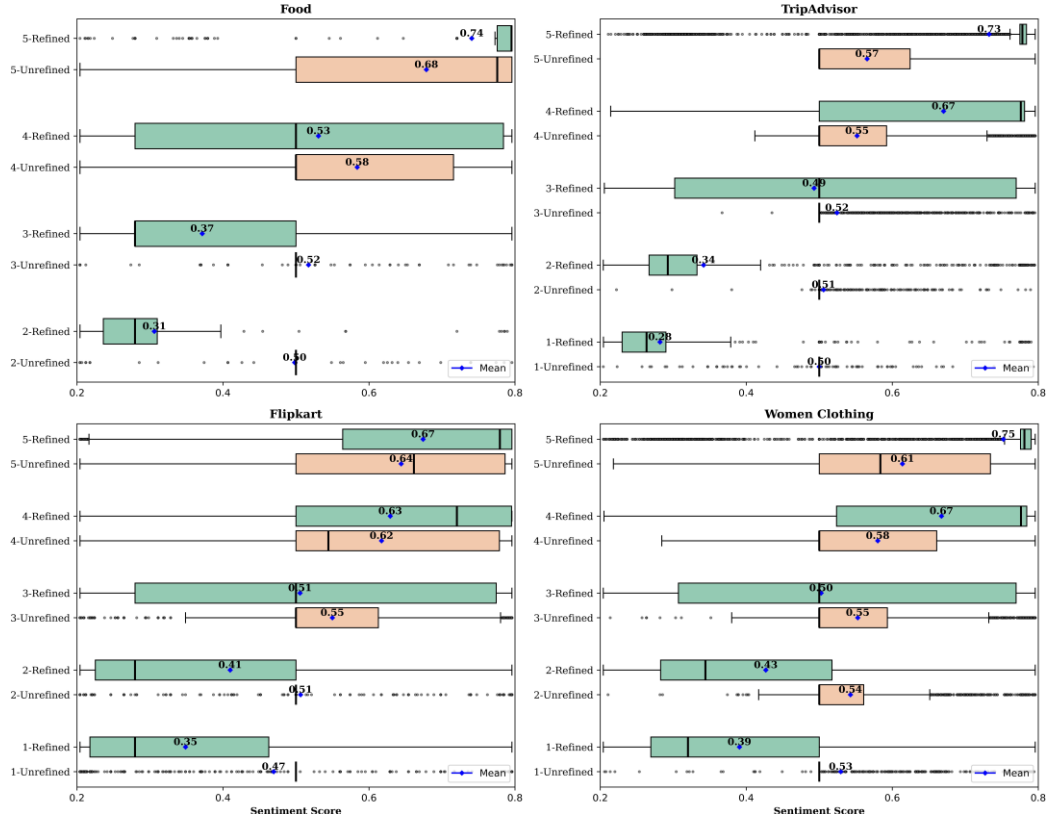


Figure 4: Sentiment Scores Distributions Stratified by User-assigned Rating

Table 3: Generated Scores & Categories for Selected Sample Reviews

Row	Comment	Dataset	Final SA (Unrefined)	Final SA (Refined)	Unrefined Category	Refined Category	Score
1	The salads sent were all bitter and clearly not fresh. The number of fillets was also fewer than usual.	Food	0.50	0.22	Neutrality	Negativity	3
2	The food was excessively greasy.	Food	0.50	0.28	Neutrality	Negativity	2
3	The fillet pieces were small as they had been halved, and they were overly breaded. Overall, it wasn't good and was as usual.	Food	0.64	0.40	Slight Positiv-ity	Neutrality	3
4	The quality was excellent, and the free salad they included gave a very good feeling. But the portion seemed decreased.	Food	0.77	0.61	Positivity	Slight Positiv-ity	4
5	Hello, I highly recommend it. The burger was incredibly delicious and very fresh. It's been a long time since I had such a meal at a restaurant.	Food	0.59	0.78	Neutrality	Positivity	5
6	The fillets were very thin and small, mostly breading. Also, if exactly 12 pieces are to be served, they should be proper portions.	Food	0.50	0.20	Neutrality	Negativity	3
7	Awesome product, super quality, clear high volume sound. Cushion quality is good.	Flipkart	0.63	0.79	Slight Positiv-ity	Positivity	5
8	In Bluetooth connection sound, bass is all good, but using aux cable, bass & sound is low. Volume button not working, and ears hurt.	Flipkart	0.50	0.30	Neutrality	Negativity	2
9	Product not working properly since purchase. Poor build, average sound quality. Will not recommend.	Flipkart	0.50	0.22	Neutrality	Negativity	2
10	Wow, superb product with good effect and sound clarity. Works well with mo- bile and system.	Flipkart	0.65	0.78	Slight Positiv-ity	Positivity	5
11	Writing this review after 9 days of usage. Sound quality is rich and impressive. Noise cancellation is good.	Flipkart	0.66	0.77	Slight Positiv-ity	Positivity	5
12	Beautiful fabric, but poor fit and proportions. Returned it.	Women's Cloth- ing	0.50	0.30	Neutrality	Negativity	2
13	Really cute piece, but it's huge. Not flattering. Returned.	Women's Cloth- ing	0.71	0.71	Positivity	Positivity	2
14	Absolutely wonderful – silky and sexy and comfortable.	Women's Cloth- ing	0.65	0.78	Slight Positiv-ity	Positivity	5
15	I ordered this 3 months ago, finally came off back order. Quality disappointing.	Women's Cloth-	0.50	0.32	Neutrality	Slight Nega-	1

3.6. Execution Time

All experiments were conducted in the Google Colab environment, utilizing both the default virtual machine and an NVIDIA T4 GPU with 16 GB RAM. Table 4 reports the total runtime of the complete pipeline across all datasets. Runtime was influenced not only by dataset size but also by average review length, which affected the Transformer component most significantly, followed by VADER. For example, although the TripAdvisor and Women's Clothing datasets have record counts similar to the Food dataset, they exhibited higher runtimes due to

longer text inputs. The reported times are also subject to variability caused by factors such as internet speed, memory allocation, and other system-level conditions. Consequently, runtime values may differ under alternative computational environments.

Table 4: Execution Time

Dataset	Number of Records	Time (Colab CPU)	Time (Colab T4 GPU)
Food	1266	3.3 min	1.4 min
Flipkart	9976	14.2 min	3.4 min
Women Clothing	23486	75.9 min	12.7 min
TripAdvisor	20491	116.2 min	19.8 min

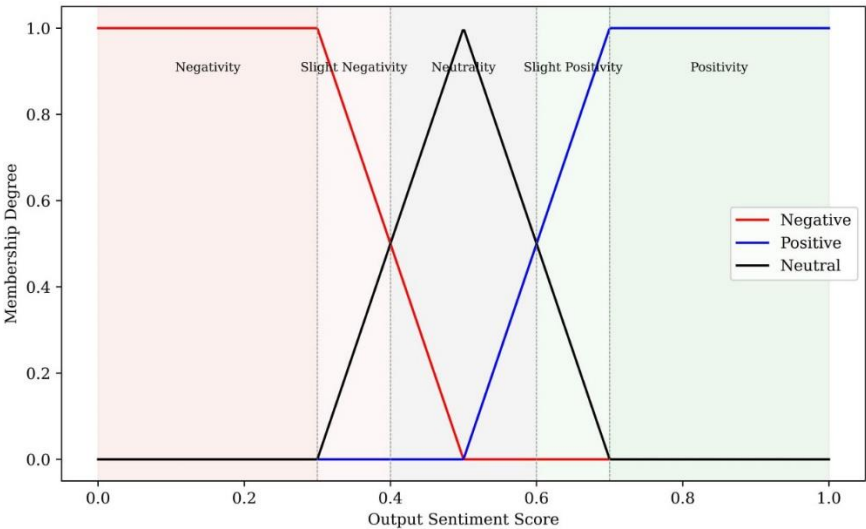


Figure 5: Conventional Categorization of Sentiment Scores

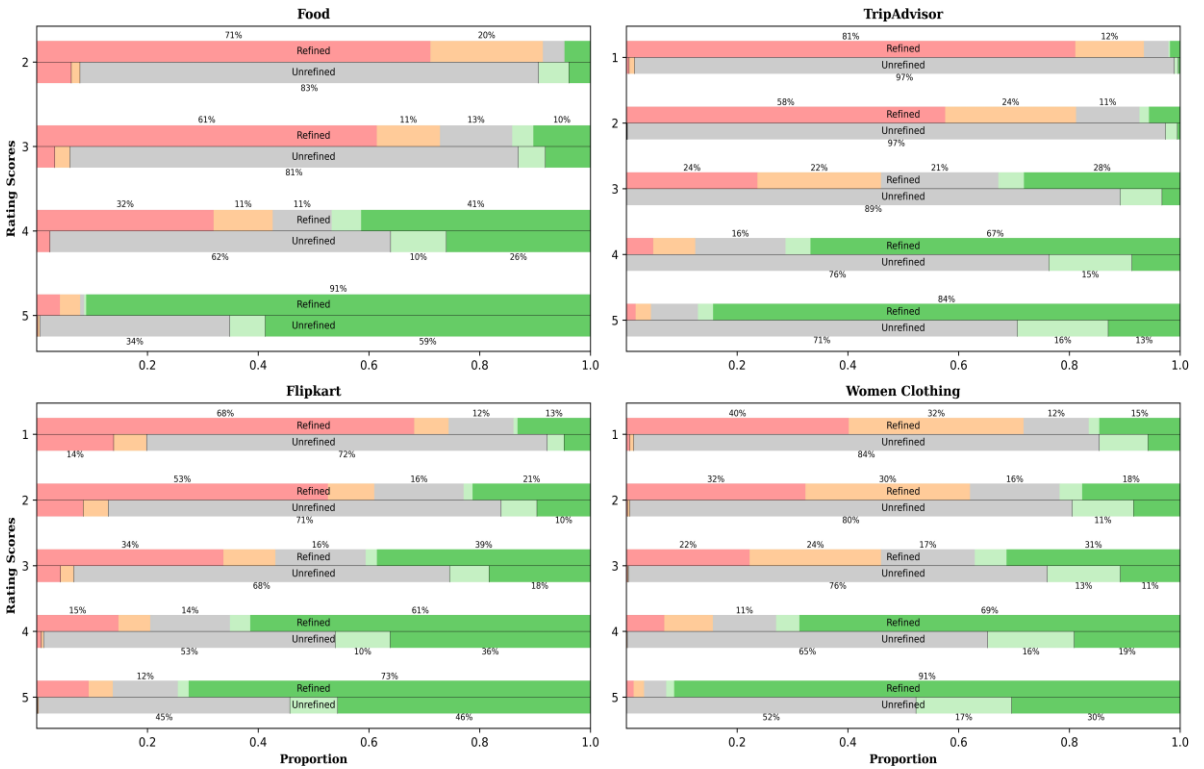


Figure 6: Distribution of Conventional Sentiment Categories across User Rating Levels

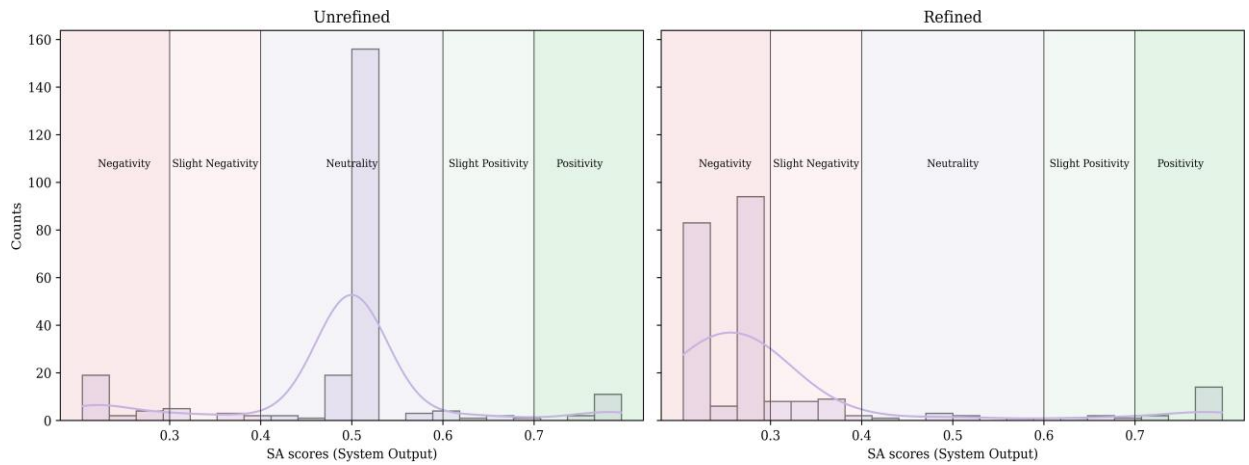


Figure 7: Redistribution and Evaluation of Missing-rated Data in the Food dataset

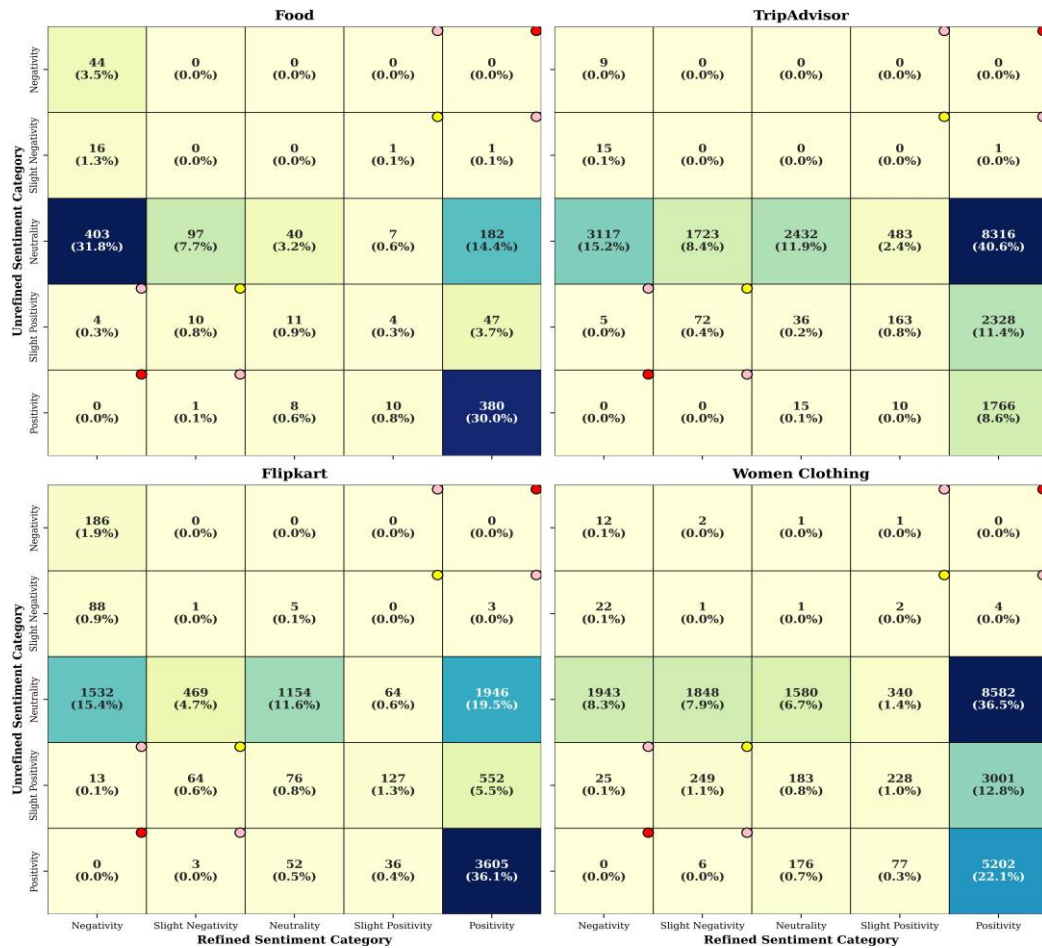


Figure 8: Confusion Matrices Showing Data Redistribution within the Proposed Pipeline

3.7. Strengths, Limitations, and Future Work Directions

3.7.1. Strengths

One of the principal strengths of this study lies in the design of an end-to-end sentiment analysis (SA) pipeline that unifies four complementary knowledge paradigms: rule-based sentiment detection (VADER), deep learning-based contextual understanding (DistilBERT), fuzzy logic for uncertainty management, and a fuzzy inference system (FIS) to simulate expert reasoning. Unlike most prior studies (Table 1) that treat sentiment intensity as discrete categories, the proposed pipeline generates continuous sentiment scores within the $[0, 1]$ range, effectively capturing both polarity and nuanced gradations, as illustrated in Sections 3.1–3.5. Another notable strength is computational efficiency. By incorporating DistilBERT—a compressed Transformer model that preserves approximately 97% of BERT’s accuracy while being 40% smaller—the framework maintains strong contextual performance with reduced memory and runtime requirements. This makes it feasible for medium-scale deployments using T4 GPUs, while larger systems equipped with A100 or H100 GPUs can achieve further acceleration. Finally, the framework maintains theoretical clarity by avoiding the conflation of probabilistic and possibilistic reasoning. DistilBERT’s confidence estimates are employed solely to recalibrate VADER’s outputs, ensuring consistent information flow across probabilistic and fuzzy components. Although probability and possibility are not formally equivalent, their alignment through sentiment-bearing word distributions provides a practically sound foundation for such cross-domain integration in NLP.

3.7.2. Limitations and Future Work Directions

Despite its contributions, the study presents several areas for further enhancement. The current refinement stage employs empirically determined threshold triplets (30%, 20%, 10%) corresponding to confidence levels (90%, 80%, 70%). These values, identified through trial and error, yielded the most stable cross-dataset performance in terms of correlation with user ratings and interpretability. Future research could replace this heuristic calibration with a learning-based optimization strategy—such as multi-layer perceptrons (MLPs)—to adapt thresholds dynamically in a data-driven or domain-independent fashion. Moreover, while the Transformer-based recalibration module substantially improved sentiment prediction accuracy, occasional misclassifications (see Table 3, Row 13) highlight the potential of ensemble strategies. Employing multiple lightweight Transformer variants combined via hard or soft voting could improve reliability and resilience against outlier expressions, though this may entail increased computational complexity. Further investigations may focus on refining the fuzzy logic rule base, adopting adaptive rule-learning techniques, experimenting with alternative inference systems such as type-2 fuzzy or neuro-fuzzy hybrids. Additionally, expanding or domain-customizing the underlying lexical resources—for example, by combining VADER with SentiWordNet or constructing context-specific lexicons—could strengthen adaptability and interpretive accuracy across diverse textual domains, thereby extending the framework’s applicability to broader NLP and sentiment-driven analytics tasks.

4. Conclusion

This study proposed a hybrid SA framework combining rule-based scoring, deep contextual modeling, and fuzzy logic to generate continuous sentiment scores capturing polarity and intensity. Unlike traditional lexicon-based approaches, the pipeline mitigates neutrality bias through a two-stage refinement process guided by Transformer-derived confidence scores. Comprehensive evaluations across four domains demonstrated the model’s effectiveness in aligning sentiment scores with user ratings, enhancing interpretability, and accurately reflecting sentiment extremes and subtle nuances. The FIS further generated expert-like, interpretable outputs continuously. While threshold parameters were selected heuristically, the framework proved robust and adaptable. Future work may explore learning-based threshold optimization, ensemble Transformer models, and domain-specific lexicon integration to enhance scalability and generalization. By bridging symbolic reasoning with contextual learning, this approach advances interpretable and adaptive SA suited for nuanced, real-time applications.

References

- [1] H. N. Do, H. T. Phan, N. T. Nguyen, Multi-modal sentiment analysis using deep learning and fuzzy logic: A comprehensive survey, *Applied Soft Computing* 167 (2024) 112279. doi:10.1016/j.asoc.2024.112279.
- [2] J. Qin, M. Zeng, X. Wei, W. Pedrycz, Ranking products through online reviews: A novel data-driven method based on interval type-2 fuzzy sets and sentiment analysis, *Journal of the Operational Research Society* 75 (5) (2024) 860–873. doi:10.1080/01605682.2023.2215823.
- [3] E. B. Ramezani, Sentiment analysis applications using deep learning advancements in social networks: A systematic review, *Neurocomputing* (2025) 129862.
- [4] P. Bedi, P. Khurana, Sentiment analysis using fuzzy-deep learning, in: *Proceedings of ICETIT 2019*, Springer, 2020, pp. 246–257. doi:10.1007/978-3-030-30577-2_21.
- [5] D. Ji, W. Meng, H. Wang, Emotional analysis of tourism reviews based on long short-term memory and fuzzy control algorithm, *International Journal of Fuzzy Systems* (2025). doi:10.1007/s40815-024-01890-1.
- [6] X. Wang, J. Lyu, B.-G. Kim, B. D. Parameshachari, K. Li, Q. Li, Exploring multimodal multiscale features for sentiment analysis using fuzzy-deep neural network learning, *IEEE Transactions on Fuzzy Systems* 33 (1) (2025) 28–42. doi:10.1109/TFUZZ.2024.3419140.
- [7] H. T. Phan, N. T. Nguyen, A fuzzy graph convolutional network model for sentence-level sentiment analysis, *IEEE Transactions on Fuzzy Systems* 32 (5) (2024) 2953–2965. doi:10.1109/TFUZZ.2024.3364694.
- [8] S. Zhou, Q. Chen, X. Wang, Fuzzy deep belief networks for semi-supervised sentiment classification, *Neurocomputing* 131 (2014) 312–322.
- [9] S. Gurumoorthy, B. N. K. Rao, X.-Z. Gao, B. Vamshi Krishna, A. K. Pandey, A. Siva Kumar, Feature based opinion mining and sentiment analysis using fuzzy logic, *Cognitive Science and Artificial Intelligence: Advances and Applications* (2018) 79–89.
- [10] M. Alzaid, F. Fkih, Sentiment analysis of students' feedback on e-learning using a hybrid fuzzy model, *Applied Sciences* 13 (23) (2023) 12956. URL <https://www.mdpi.com/2076-3417/13/23/12956>
- [11] M. Shamami, H. Farahzadi, L. Amini, M. Ilani, Y. Banad, Advanced classification of drug-drug interactions for assessing adverse effect risks of fluvoxamine and curcumin using deep learning in covid-19, *Journal of Infrastructure, Policy and Development* 8 (15) (2024) 9734.
- [12] S. He, Y. Wang, Evaluating new energy vehicles by picture fuzzy sets based on sentiment analysis from online reviews, *Artificial Intelligence Review* 56 (3) (2023) 2171–2192. doi:10.1007/s10462-022-10217-1.
- [13] S. Abbasi, S. Rokhva, K. Farahmand, P. Ghasemi, E. Shahab, Designing sustainable closed-loop supply chain network based on a circular economy approach: under uncertainty during the post-covid era, *Circular Economy and Sustainability* 5 (3) (2025) 2229–2271.
- [14] F. Ali, D. Kwak, P. Khan, S. M. R. Islam, K. H. Kim, K. S. Kwak, Fuzzy ontology-based sentiment analysis of transportation and city feature reviews for safe traveling, *Transportation Research Part C: Emerging Technologies* 77 (2017) 33–48. doi:10.1016/j.trc.2017.01.014.
- [15] S. Vashishtha, S. Susan, Fuzzy rule based unsupervised sentiment analysis from social media posts, *Expert Systems with Applications* 138 (2019) 112834. doi:10.1016/j.eswa.2019.112834.
- [16] Y. A. Reddy, S. Agarwal, V. Parashar, A. Arora, Real-time sentiment insights from x using vader, distilbert, and web-scraped data, *arXiv preprint arXiv:2504.15448* (2025). doi:10.48550/arXiv.2504.15448.
- [17] M. Ágüa, N. Antonio, M. P. Carrasco, C. Rassal, Large language models powered aspect-based sentiment analysis for enhanced customer insights, *Tourism & Management Studies* 21 (1) (2025) 1. doi:10.18089/tms.20250101.
- [18] M. Isnan, G. N. Elwirehardja, B. Pardamean, Sentiment analysis for tiktok review using vader sentiment and svm model, *Procedia Computer Science* 227 (2023) 168–175. doi:10.1016/j.procs.2023.10.514.
- [19] R. R. Mukkamala, A. Hussain, R. Vatrappu, Fuzzy-set based sentiment analysis of big social data, in: *2014 IEEE 18th International Enterprise Distributed Object Computing Conference*, 2014, pp. 71–80. doi:10.1109/EDOC.2014.19.
- [20] A. Mahmoudi, D. Jemielniak, L. Ciechanowski, Assessing accuracy: A study of lexicon and rule-based packages in r and python for sentiment analysis, *IEEE Access* 12 (2024) 20169–20180. doi:10.1109/ACCESS.2024.3353692.
- [21] A. Sarabadani, M. A. Shamami, H. Sadeghsalehi, B. Asadi, S. Hesarakhi, Dkg-llm: A framework for medical diagnosis and personalized treatment recommendations via dynamic knowledge graph and large language model integration, *arXiv preprint arXiv:2508.06186* (2025).
- [22] D. Bao, W. Su, Optimizing deep learning-based natural language processing for sentiment analysis, *International Journal of High Speed Electronics and Systems* (2025) 2540304. doi:10.1142/S0129156425403043.

- [23] R. Gore, M. M. Safaei, C. J. Lynch, C. P. Ames, A spine-specific lexicon for the sentiment analysis of interviews with adult spinal deformity patients correlates with sf-36, srs-22, and odi scores: A pilot study of 25 patients, *Information* 16 (2) (2025) 2. doi:10.3390/info16020090.
- [24] M. S. Djalolovna, The evolution of informal speech: How language changes in everyday conversations, *Web of Teachers: Inderscience Research* 3 (1) (2025) 109–113.
- [25] K. Alahmadi, S. Alharbi, J. Chen, X. Wang, Generalizing sentiment analysis: a review of progress, challenges, and emerging directions, *Social Network Analysis and Mining* 15 (1) (2025) 1–28.
- [26] L. Mei, S. Liu, Y. Wang, B. Bi, X. Cheng, Slang: New concept comprehension of large language models, *arXiv preprint arXiv:2401.12585* (2024).
- [27] M. S. Islam, et al., Challenges and future in deep learning for sentiment analysis: a comprehensive review and a proposed novel hybrid approach, *Artificial Intelligence Review* 57 (3) (2024) 62. doi:10.1007/s10462-023-10651-9.
- [28] M. Dragoni, A. G. B. Tettamanzi, C. da Costa Pereira, A fuzzy system for concept-level sentiment analysis, in: *Semantic Web Evaluation Challenge*, Springer, 2014, pp. 21–27. doi:10.1007/978-3-319-12024-9_2.
- [29] Z. Raeisi, S. Rokhva, A. Roshanzamir, R. Ahmadi Lashaki, An accurate attention-based method for multi-tasking x-ray classification, *Multimedia Tools and Applications* (2025) 1–25.
- [30] M. Dragoni, G. Petrucci, A fuzzy-based strategy for multi-domain sentiment analysis, *International Journal of Approximate Reasoning* 93 (2018) 59–73. doi:10.1016/j.ijar.2017.10.021.
- [31] H. A. Hasan, Z. Alassedi, M. Adile, Z. Alsalami, A. Hussian, Movie review analysis using fuzzy logic-based natural language processing, in: *2024 International Conference on Smart Systems for Electrical, Electronics, Communication and Computer Engineering (ICSSECC)*, 2024, pp. 751–756. doi:10.1109/ICSSECC61126.2024.10649403.
- [32] A. Saualih, et al., Exploring the tourist experience of the majorelle garden using vader-based sentiment analysis and the latent dirichlet allocation algorithm: The case of tripadvisor reviews, *Sustainability* 16 (15) (2024) 15. doi:10.3390/su16156378.
- [33] M. Z. Mekonen, et al., An opinionated sentiment analysis using a rule-based method, *Bulletin of Electrical Engineering and Informatics* 14 (1) (2025) 1. doi:10.11591/eei.v14i1.8568.
- [34] I. Karabila, N. Darraz, A. El-Ansari, N. Alami, M. El Mallahi, A hybrid approach combining sentiment analysis and deep learning to mitigate data sparsity in recommender systems, *Neurocomputing* 636 (2025) 129886.
- [35] Z. Raeisi, S. Rokhva, F. Rahmani, A. Goodarzi, H. Najafzadeh, Multi-label diagnosis of dental conditions from panoramic x-rays using attention-enhanced deep learning, *Oral and Maxillofacial Surgery* 29 (2025) 166. doi: <https://link.springer.com/article/10.1007/s10006-025-01463-y>
- [36] E. Zangeneh, S. Rokhva, Application of machine learning in predictive maintenance scheduling: An industrial case study, *International Journal of Industrial Engineering and Operational Research* 7 (3) (2025) 19–27. doi:10.22034/ijieor.v7i3.169.