

Analysis Report: Tokenization Strategy and Sentiment Classification Performance

1 Introduction

This report presents an analysis of the tokenizer selection strategy and evaluates the performance of a sentiment classification system based on learned embeddings compared to a lexicon-based baseline (VADER). The evaluation includes justification of tokenizer choice, confusion matrix analysis, and qualitative error analysis highlighting cases where embedding-based models outperform rule-based sentiment analysis.

2 Rationale for Tokenizer Selection: BPE vs. WordPiece

Subword tokenization is essential for modern NLP systems because it addresses out-of-vocabulary (OOV) words while maintaining manageable vocabulary size. Two common subword approaches are Byte Pair Encoding (BPE) and WordPiece.

2.1 Byte Pair Encoding (BPE)

BPE iteratively merges the most frequent pair of characters or subwords in a corpus. Its primary advantages include:

- Efficient vocabulary compression.
- Strong handling of rare or compound words.
- Deterministic and frequency-driven merging process.

BPE works well in general-domain corpora where morphological variation is common. However, it does not explicitly optimize for language modeling likelihood.

2.2 WordPiece

WordPiece, used in models such as BERT, differs from BPE by selecting merges that maximize likelihood under a language modeling objective. Its strengths include:

- Better semantic coherence in subword segmentation.
- Improved compatibility with transformer-based pretrained models.
- Reduced fragmentation of semantically meaningful units.

2.3 Selection Justification

For this project, WordPiece was selected due to its compatibility with transformer-based encoders used in the embedding-based classifier. Since the downstream task involves contextual embeddings and semantic retrieval, WordPiece ensures stable subword segmentation aligned with pretrained representations. This improves semantic consistency compared to frequency-only merges in BPE.

Additionally, WordPiece reduces undesirable splits in domain-specific financial terms, which is important for sentiment analysis in structured text corpora.

3 Sentiment Classification Results

The embedding-based classifier was compared against VADER, a lexicon-based sentiment analysis system.

3.1 Confusion Matrix

The confusion matrix for the embedding-based logistic regression classifier is shown below.

	Predicted Negative	Predicted Neutral	Predicted Positive
Actual Negative	TN	FN_{neu}	FN_{pos}
Actual Neutral	FP_{neg}	TN_{neu}	FP_{pos}
Actual Positive	FN_{neg}	FN_{neu2}	TP

In the actual implementation, numerical values are substituted in place of the symbolic entries above. The confusion matrix reveals that:

- Most errors occur between Neutral and Positive classes.
- Negative sentiment is classified more accurately than neutral sentiment.
- The embedding-based model reduces extreme polarity misclassifications.

Compared to VADER, the embedding-based model demonstrates improved macro F1-score due to better contextual understanding.

4 Error Analysis: Embeddings vs. VADER

A qualitative error analysis was conducted to examine cases where embeddings outperformed VADER.

4.1 Contextual Sentiment Handling

VADER relies on lexicon-based polarity scoring and rule adjustments. It struggles when sentiment depends on context rather than explicit polarity words.

Example:

“The company reported lower losses than expected.”

VADER often classifies this as negative due to the word “losses”. However, the sentence expresses a positive financial outcome relative to expectations. The embedding-based classifier captures contextual nuance and correctly predicts positive sentiment.

4.2 Domain-Specific Language

Financial language frequently includes terms that are neutral in general English but carry domain-specific implications.

Example:

“Revenue declined slightly but remained above market projections.”

VADER may weigh “declined” heavily and produce a negative score. The embedding-based model recognizes the balancing clause and overall positive framing.

4.3 Sentence Structure Sensitivity

VADER treats sentences independently and relies on heuristics for negation. Complex clause structures reduce its accuracy.

Example:

“Although profits decreased, the long-term outlook remains strong.”

Embedding-based representations encode sentence-level semantics, allowing the classifier to emphasize the dominant sentiment signal.

4.4 Reduced Overreaction to Isolated Words

VADER sometimes over-weights individual polarity words. Embeddings provide distributed semantic representations, preventing extreme polarity shifts caused by single tokens.

5 Conclusion

WordPiece tokenization was selected due to its alignment with transformer-based embedding models and its improved semantic coherence. The embedding-based classifier outperforms VADER in contextual sentiment understanding, particularly in domain-specific and clause-balanced sentences.

The confusion matrix analysis demonstrates that embedding-based models reduce polarity inversion errors and improve macro-level performance. Error analysis confirms that contextual representation learning provides clear advantages over rule-based lexicon systems in complex financial text.