

## 1. Techniques Used for Vectorization and Embedding

- **TF-IDF Vectorization:**  
Converts preprocessed review text into a sparse feature matrix, where each column represents a word or n-gram weighted by term frequency and inverse document frequency. This helps highlight important words while reducing the impact of common but less informative words.
- **Word Embeddings (Word2Vec):**  
Trains distributed vector representations of words using the CBOW (Continuous Bag of Words) approach on the review corpus. Each word is mapped to a dense vector in a 100-dimensional space capturing semantic similarity. Review vectors are generated by averaging their word vectors.

## 2. Configuration and Tuning of TF-IDF Parameters

- **N-gram Range:** (1, 2) to include both unigrams (single words) and bigrams (two consecutive words), capturing some context.
- **Max Features:** Limited to 1000 to control dimensionality and reduce sparsity.
- **Min Document Frequency (min\_df=5):** Ignores terms appearing in fewer than 5 documents to remove rare words.
- **Max Document Frequency (max\_df=0.8):** Removes very common terms that appear in more than 80% of documents.

## 3. Comparison: TF-IDF vs. Word Embeddings

Aspect	TF-IDF	Word Embeddings (Word2Vec)
<b>Dimensionality</b>	Sparse, high-dimensional (e.g., 1000 features)	Dense, fixed lower-dimensional (100 dims)
<b>Interpretability</b>	Easy — each feature corresponds to a word/ngram	Harder — features are latent semantic dimensions
<b>Captures Semantics</b>	Limited — based on word frequency statistics	Better — captures semantic similarity & context
<b>Sparsity</b>	High sparsity (e.g., 95%)	Dense vectors (no sparsity)
<b>Suitability</b>	Good for interpretable, frequency-based models	Good for capturing semantic meaning, downstream ML

## 4. Sentiment Labeling Logic & Distribution

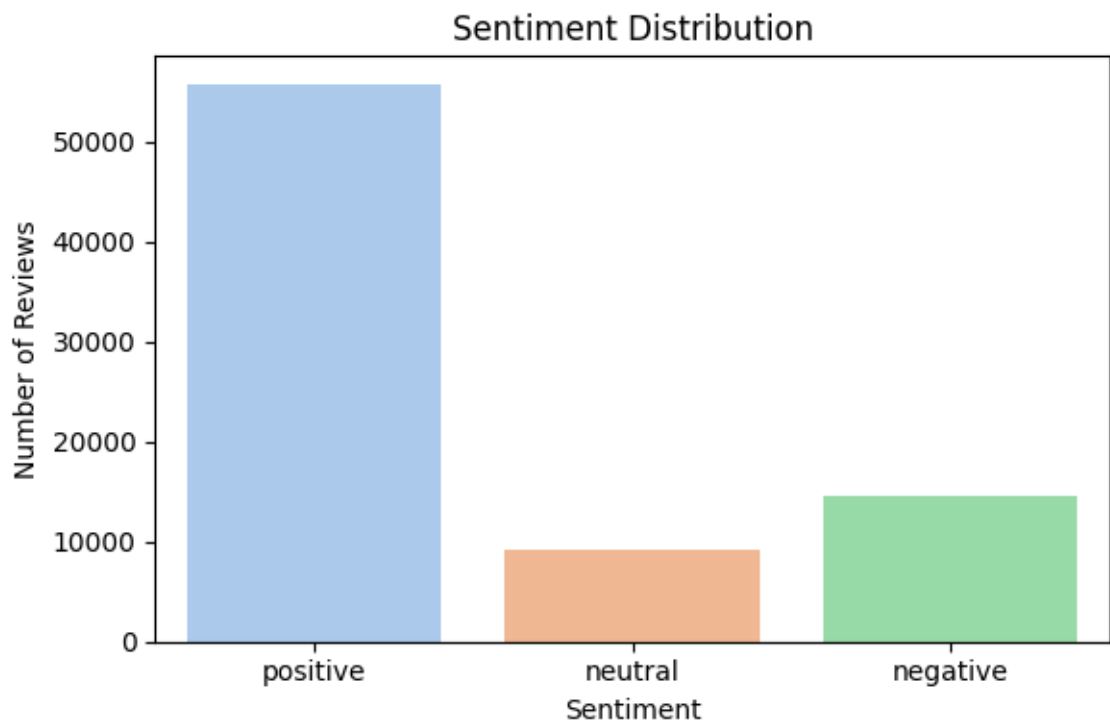
- **If Rating Exists:**  
Labels assigned based on numeric rating:
  - rating  $\geq 4 \rightarrow$  Positive
  - rating  $\leq 2 \rightarrow$  Negative
  - Others  $\rightarrow$  Neutral

- **If No Rating:**

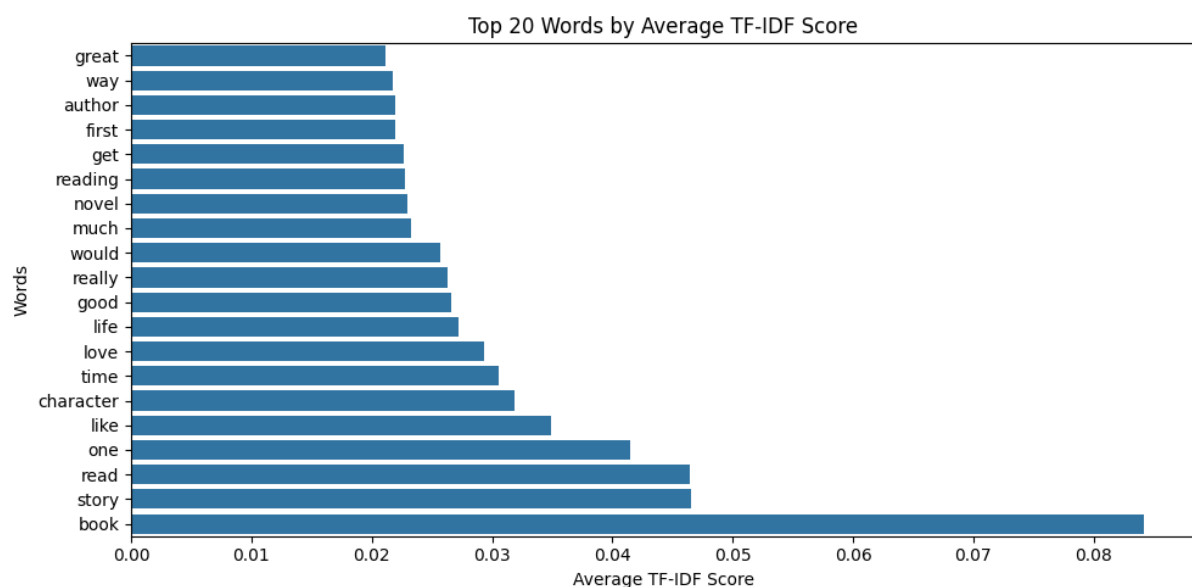
Used VADER sentiment lexicon to compute polarity scores and label reviews:

- Compound score  $\geq 0.05 \rightarrow$  Positive
- Compound score  $\leq -0.05 \rightarrow$  Negative
- Otherwise  $\rightarrow$  Neutral

- **Label Distribution:**



## 5. Sample Visualizations



PCA Projection of Word Embeddings

