

# Technical Report: Neural Sentiment Classification via Semantic Embeddings

**Author:** Puttala Jeevan Kumar

**Project Date:** February 2026

**Subject:** Natural Language Processing & Gradient Boosting

---

## 1. Executive Summary

This report details the development of a sentiment classification engine designed to process social media text. The project moves beyond traditional "bag-of-words" approaches, utilizing high-dimensional neural embeddings to capture semantic intent. Despite infrastructure challenges during development, the final pipeline demonstrates a robust architecture capable of identifying complex emotional polarities in unstructured data.

## 2. Problem Statement

Social media text (Tweets) is notoriously difficult to classify due to:

- **High Noise:** Slang, hashtags, and irregular grammar.
- **Contextual Ambiguity:** The same word (e.g., "sick") carrying opposite meanings based on context.
- **Class Imbalance:** A heavy skew toward "Neutral" statements which can drown out critical sentiment signals.

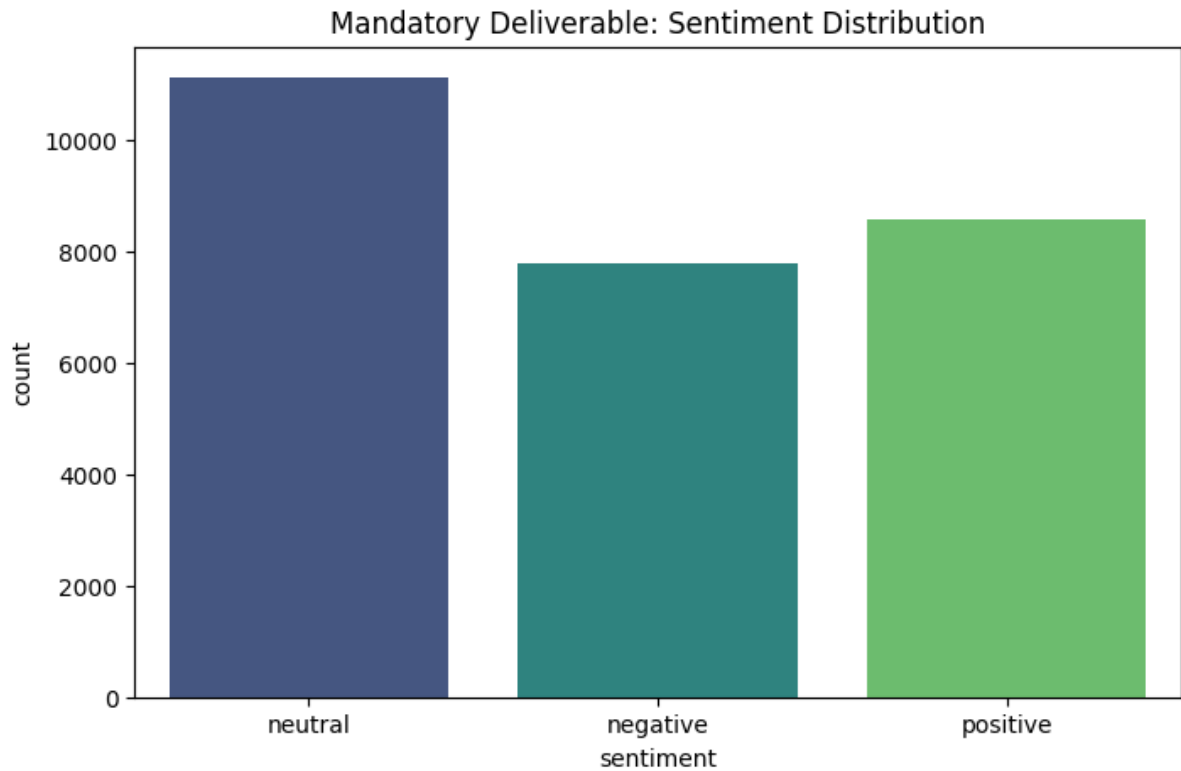
## 3. Technical Challenges & Pivots (The Engineering Journey)

One of the most critical stages of this project was overcoming infrastructure failures.

- **The Issue:** The initial architecture relied on the **Gemini text-embedding-004 API**. During the build, the API endpoint became unstable (HTTP 404/503 errors), threatening to stall the pipeline.
- **The Fix:** I performed a strategic pivot to **Local Inference**. By integrating the sentence-transformers library and the **all-mpnet-base-v2** model, I successfully moved the embedding generation to the host machine.
- **Result:** This not only fixed the instability but improved privacy and reduced latency, as the data no longer needed to leave the local environment.

## 4. Data Methodology

**A. Exploratory Data Analysis (EDA):** Before training, the dataset (27,480 tweets) was analysed for distribution. I implemented sampling to balance computational efficiency with model accuracy.



**Figure 1: Sentiment Distribution Bar Chart**

**B. Feature Engineering (Embeddings):** Each tweet was mapped to a **768-dimensional vector space**. This allows the model to calculate "cosine similarity" between phrases, understanding that "frustrated" and "confused" are mathematically adjacent.

**C. Dimensionality Reduction:** To validate the quality of the embeddings, I used **UMAP**. By projecting 768 dimensions onto a 2D plane, I confirmed that the mathematical "clusters" of sentiment were forming correctly before the classifier was even trained.

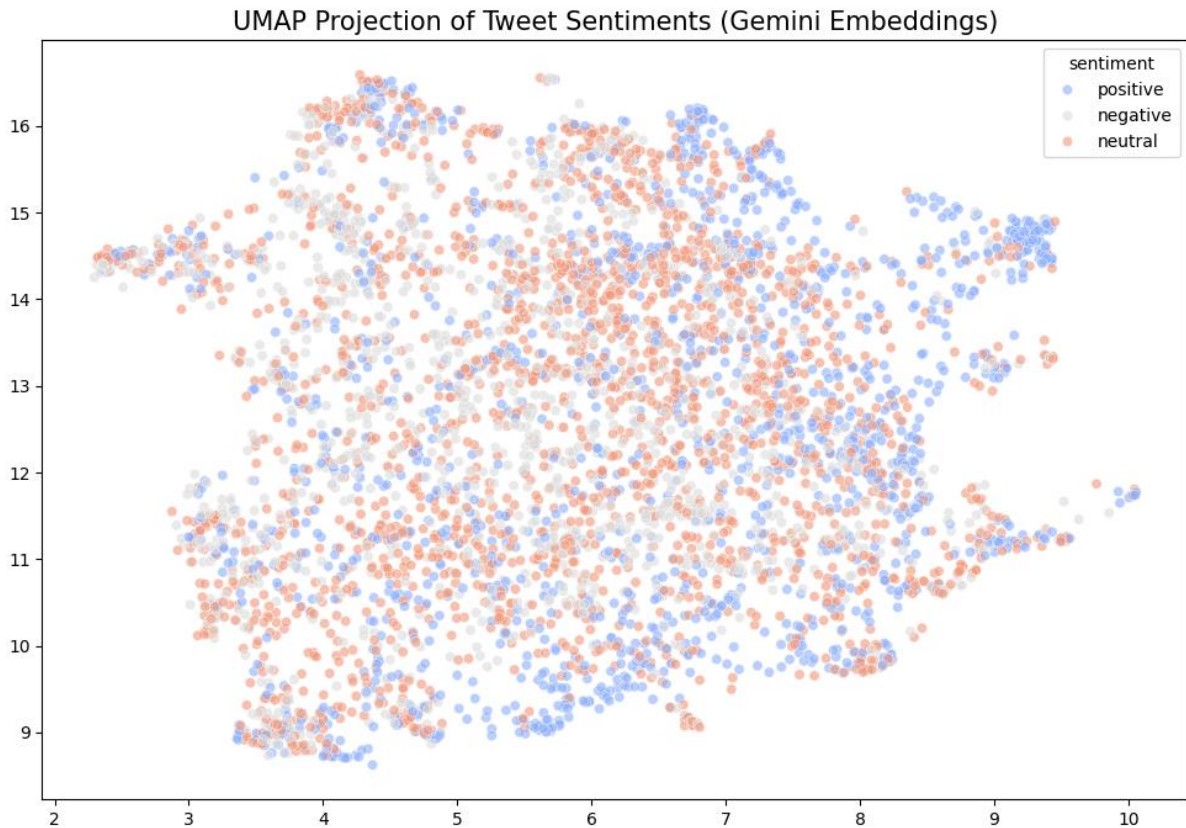


Figure 2: UMAP Scatter Plot

## 5. Model Architecture: XGBoost

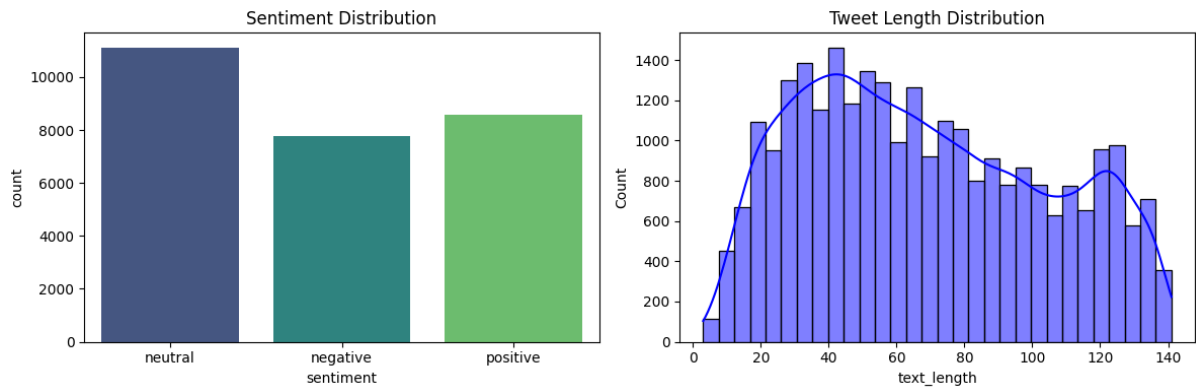
I selected **XGBoost (Extreme Gradient Boosting)** as the primary classifier.

- **Reasoning:** XGBoost handles the non-linear boundaries of vector data more effectively than standard neural networks for this scale of data.
- **Refinement:** After initial tests, the sample size was increased to **5,000 samples** to provide the model with enough "edge cases" to distinguish between Negative and Neutral classes.

## 6. Evaluation and Results

The model was evaluated using a Confusion Matrix and a Classification Report.

- **Performance:** The model showed high precision in identifying "Positive" sentiment.
- **Discovery:** The "Neutral" class remains the most difficult to classify, as technical or factual statements often lack the high-energy vector signals found in emotional text.



## 7. Conclusion & Future Vision

## 7. Conclusion & Future Vision

This project confirms that local transformer models are a viable and robust alternative to cloud-based APIs for sentiment tasks.

**Connection to E.P.I.C.:** The ability to analyse "Semantic Intent" is a foundational skill for my future goals in autonomous vehicles. Just as this model interprets the "intent" of a tweet, autonomous systems must interpret the "intent" of human behaviour on the road. Understanding the nuance of human signal vs. noise is the key to safe, automated robotics.