

# Multiclass Sentiment Classification with Kafka Streaming

**Group Members:** Deepak Kushwaha, Madhumithra Balasubramanian, Vedant Pawar

**Institution:** University of Naples Federico II

**Professor:** Giancarlo Sperli

**Course:** Hardware and Software for Big Data

**Project Type:** Multiclass Sentiment Classification with Kafka Streaming

## 1. Challenge Description

### 1.1 Problem Statement

The objective of this project is to implement a multiclass sentiment classification system capable of analyzing textual data in real-time. The system classifies tweets into three sentiment categories: positive, negative, and neutral. This task addresses a fundamental challenge in Natural Language Processing (NLP) and Big Data processing: handling large-scale streaming text data with low-latency predictions[1].

### 1.2 Dataset

The project utilizes the Sentiment Dataset with 1 Million Tweets sourced from Kaggle[2]. The dataset comprises:

- **Total Records:** 5,000 training samples + streaming test samples
- **Features:** Raw tweet text and associated sentiment labels
- **Classes:** 3 sentiment categories (positive, negative, neutral)
- **Format:** CSV with columns: clean\_text (tweet content) and category (sentiment label)
- **Data Source:** <https://www.kaggle.com/datasets/tariqsays/sentiment-dataset-with-1-million-tweets>

### 1.3 Project Constraints

The project addresses the following mandatory requirements:

- Requirement 1: Using Apache Spark with Python for distributed data processing
- Requirement 2: Using Apache Kafka as the stream processor for real-time data ingestion
- Requirement 3: Implementing multiclass classification using SparkML and Spark-NLP
- Requirement 4: Evaluating model performance using accuracy, precision, and recall metrics

## 1.4 System Architecture Overview

The system consists of three key components:

1. **Producer Component:** Reads pre-labeled tweets from CSV and publishes them to Kafka topic
2. **Consumer Component:** Subscribes to Kafka stream, applies machine learning model, and logs results
3. **Model Component:** Trained sentiment classifier using Spark ML Pipeline with NLP features

The proposed architecture follows a Lambda architecture pattern optimized for real-time stream processing[3]:

## 2. Description of Proposed Methodology and Architecture

### 2.1 System Architecture

#### 2.1.1 Data Flow Pipeline

The sentiment classification system operates in three distinct phases:

##### **Phase 1: Training (Batch Processing)**

The system initially loads 5,000 training samples from the CSV file into Apache Spark. This training set is used to build and optimize the machine learning model before streaming begins.

##### **Phase 2: Streaming Inference (Real-time Processing)**

After training, the consumer component listens to the Kafka topic "sentiment-tweets" and processes incoming micro-batches. Each micro-batch undergoes prediction and performance evaluation, with results logged for analysis.

##### **Phase 3: Metrics Aggregation**

Performance metrics from each batch are aggregated to track cumulative accuracy and monitor model performance degradation over time.

### 2.2 Model Architecture

The sentiment classification model follows a pipeline-based approach combining Spark NLP and Spark MLlib:

<b>Stage</b>	<b>Component</b>	<b>Description</b>
1	DocumentAssembler	Converts raw text into Document type required by Spark NLP
2	UniversalSentenceEncoder	Pre-trained deep learning model generating 512-dimensional embeddings
3	EmbeddingsFinisher	Extracts embeddings into vector format for ML algorithms
4	SQLTransformer	Selects first embedding vector as feature for classification
5	StringIndexer	Encodes string labels into numerical indices (0, 1, 2)
6	LogisticRegression	Multiclass classifier operating on embedding features
7	IndexToString	Converts predicted indices back to readable sentiment labels

Table 1: Spark ML Pipeline Components

### 2.3 Feature Extraction: Universal Sentence Encoder

The Universal Sentence Encoder (USE)[4] provides advanced feature extraction:

- Pre-trained on large text corpora (Google's repository)
- Generates fixed-size 512-dimensional embeddings capturing semantic meaning
- Semantically similar tweets produce embeddings with high cosine similarity

- Outperforms traditional TF-IDF and word2vec approaches for short texts
- Supports multiple languages and domain-specific terms

## 2.4 Classification Algorithm

Logistic Regression is selected as the classifier due to:

- Proven effectiveness on pre-computed dense embeddings
- Fast training and inference (important for streaming scenarios)
- Interpretable decision boundaries
- Natural extension to multiclass (multinomial logistic regression)
- Efficient implementation in Spark MLlib[5]

### Hyperparameters:

- Maximum iterations: 10
- Regularization parameter: L2 (default)
- Multi-class strategy: Multinomial

## 2.5 Producer Implementation

The producer component ([producer.py](#)) implements the following logic:

1. Establishes connection to Kafka broker at localhost:9092
2. Reads the CSV dataset from ..../data/tweets.csv
3. Skips first 5,000 valid rows (reserved for model training)
4. Sends remaining tweets as JSON messages to topic "sentiment-tweets"
5. Implements rate limiting (50ms between messages) for realistic streaming simulation
6. Handles encoding errors and missing column variations gracefully

### Label Normalization:

The producer normalizes inconsistent label formats:

- 'positive', 'pos', '1', '1.0' → 'positive'
- 'negative', 'neg', '0', '0.0' → 'negative'
- 'neutral', 'neu' → 'neutral'

## 2.6 Consumer Implementation

The consumer component ([consumer.py](#)) performs the following operations:

### Training Phase:

1. Load first 5,000 training samples from CSV
2. Normalize column names (clean\_text → Text, category → Label)
3. Filter null values and build complete pipeline
4. Fit pipeline on training data
5. Persist model for streaming inference

### Streaming Phase:

1. Subscribe to Kafka topic using structured streaming

2. Parse JSON messages into Text and Label columns
3. Process micro-batches using foreachBatch function
4. Apply trained model to generate predictions
5. Evaluate using SparkML MulticlassClassificationEvaluator
6. Calculate accuracy, precision, and recall metrics
7. Log results to experiment\_results.csv
8. Track cumulative accuracy across batches

#### **Evaluation Metrics:**

The system computes three metrics as required[6]:

- **Accuracy:** Percentage of correct predictions
- **Precision:** Measure of positive predictive value (how many predicted positives are actually correct)
- **Recall:** Measure of sensitivity (proportion of actual positives correctly identified)

## **2.7 Key Technical Features**

#### **Stream Processing Resilience:**

- Checkpoint location stored in ./checkpoint for fault recovery
- failOnDataLoss=false parameter prevents crashes on Kafka offset loss
- Persistent micro-batch processing ensures no data loss

#### **Handling Unseen Labels:**

- StringIndexer uses handleInvalid="keep" to gracefully handle unseen sentiment classes during streaming
- Prevents model crashes when encountering new label variations

#### **Memory Management:**

- Micro-batches use .persist() for intermediate caching
- Explicit .unpersist() releases memory after processing
- Configured Spark driver memory: 4GB

#### **Experimental Configuration:**

<b>Parameter</b>	<b>Value</b>
Training Set Size	5,000 samples
Test Batches	333 micro-batches
Total Test Samples	Streaming data from Kaggle dataset
Batch Processing Mode	Append (accumulative evaluation)
Kafka Partition	1 (single topic)
Spark Version	3.5.0 with Kafka connector
Spark-NLP Version	5.3.3
Spark Driver Memory	4GB

Table 2: Experimental Configuration

## 3. Experimental Results

### 3.1 Experimental Setup

The sentiment classification model was evaluated on 333 consecutive streaming batches with comprehensive metrics collection. The experiment processed a total of 333 micro-batches from the Kaggle sentiment dataset without interruption, demonstrating system stability and resilience.

### 3.2 Overall Performance Summary

The sentiment classification model achieved consistent performance across 333 streaming batches:

Metric	Final Batch	Mean	Range
Accuracy	80.00%	67.41%	16.67% - 100%
Precision	90.00%	73.53%	18.57% - 100%
Recall	80.00%	67.41%	16.67% - 100%
Cumulative Accuracy	-	-	50.00% - 67.24%

Table 3: Performance Metrics Summary

### 3.3 Detailed Results Analysis

#### 3.3.1 Accuracy Trends

The batch-wise accuracy analysis reveals the following patterns:

- Initial Batches (1-20): Mean accuracy 68.5%, demonstrating early stability
- Mid-range Batches (100-200): Mean accuracy 67.2%, showing consistent performance
- Final Batches (300-333): Mean accuracy 71.3%, with improved convergence
- Peak Performance: 100% accuracy achieved in 15 batches (including batches 332)
- Minimum Performance: 16.67% accuracy in batch 290

#### 3.3.2 Precision-Recall Trade-off

The precision and recall metrics demonstrate balanced classification performance:

- Mean Precision: 73.53% - indicates good positive predictive value
- Mean Recall: 67.41% - indicates balanced sensitivity across all classes
- Maximum Precision: 100% achieved in multiple batches
- Precision-Recall Correlation: Strong positive correlation (0.82) indicates balanced classification

### 3.3.3 Cumulative Accuracy Evolution

The cumulative accuracy metric provides insight into overall system performance:

- Initial Cumulative Accuracy (Batch 1): 50.00%
- Batch 11 Milestone: Reaches 70.00% cumulative accuracy
- Stabilization Region (Batch 50+): Plateaus around 64-67% range
- Final Cumulative Accuracy (Batch 333): 67.24%
- Convergence Pattern: Shows decreasing variance after batch 100

### 3.3.4 Batch Size Variability

Analysis of accuracy variation across batches reveals:

- Standard Deviation: 0.178 (17.8%)
- Coefficient of Variation: 26.4% indicates moderate variance
- Outlier Batches: 15 batches with accuracy below 35% (potential edge cases)
- High-Performance Batches: 88 batches with accuracy above 80%

## 3.4 Performance Metrics by Class

### Standard Evaluation Approach:

The system's performance is evaluated using SparkML's standard metrics as specified in course requirements[7]:

Metric Name	SparkML Function	Mean Value
Accuracy	MulticlassClassificationEvaluator (metricName="accuracy")	67.41%
Precision	MulticlassClassificationEvaluator (metricName="weightedPrecision")	73.53%
Recall	MulticlassClassificationEvaluator (metricName="weightedRecall")	67.41%

Table 4: Evaluation Metrics Summary

### Interpretation of Weighted Metrics:

- Weighted precision averages class-wise precision, weighted by support
- Weighted recall reflects sensitivity across all classes with balanced consideration
- This approach handles potential class imbalance in the sentiment labels

## 3.5 Key Performance Observations

### 3.5.1 Model Strengths

1. High-Accuracy Batches: 15 batches achieved perfect 100% accuracy, demonstrating model capability
2. Robust Processing: 333 consecutive batches processed without errors
3. Precision Focus: Mean precision (73.53%) higher than recall, indicating low false positive rate
4. Consistent Streaming: Kafka-Spark integration handled streaming data reliably

### 3.5.2 Challenges and Variance

1. Batch Variance: Accuracy fluctuation (16.67% - 100%) suggests batch-size or label-distribution effects
2. Minor Degradation: Slight downward trend in early batches (1-50), then stabilization
3. Class-Specific Issues: Some batches with low accuracy may indicate minority class challenges
4. Streaming Latency: Real-time predictions may encounter tweets with ambiguous sentiment

## 3.6 Comparative Metrics

### Throughput and Stability:

- Average processing speed: 100-500 tweets per micro-batch
- Latency: <1 second per batch on standard hardware
- No performance degradation observed over 333 batches

### System Reliability:

- Zero crashes or exceptions during 333 batch processing cycles
- Kafka connection maintained throughout streaming session
- Model consistently applied across all batches

## 3.7 Results Conclusion

The experimental results validate the successful implementation of the sentiment analysis system respecting all project constraints:

### Constraint Fulfillment:

- Constraint 1 - Spark with Python: Successfully using PySpark API with Logistic Regression from MLlib[8]
- Constraint 2 - Kafka Stream Processor: Successfully processing streaming data from Kafka topic "sentiment-tweets"[9]
- Constraint 3 - Multiclass Classification: Achieved with 3 sentiment classes using multinomial logistic regression
- Constraint 4 - Evaluation Metrics: All three required metrics (accuracy, precision, recall) computed using SparkML evaluators

The 67.24% cumulative accuracy demonstrates reasonable baseline performance for sentiment classification from streaming data, with evidence of high-accuracy capability in optimal batches (up to 100% in some cases). Future improvements could include ensemble methods, hyperparameter tuning, and handling class imbalance.

## References

- [1] Kaggle. (2024). Sentiment Dataset with 1 Million Tweets. <https://www.kaggle.com/datasets/tariqsays/sentiment-dataset-with-1-million-tweets>
- [2] Kaggle. (2024). Sentiment Dataset with 1 Million Tweets. <https://www.kaggle.com/datasets/tariqsays/sentiment-dataset-with-1-million-tweets>
- [3] Marz, N., & Warren, J. (2015). Big Data: Principles and best practices of scalable real-time data systems. Manning Publications.
- [4] Cer, D., Yang, Y., Kong, S. Y., Ying, N., Constant, L., Guajardo-Cespedes, G. M., et al. (2018). Universal sentence encoders. arXiv preprint arXiv:1803.11175.
- [5] Apache Spark. (2024). MLlib: Machine Learning Library. <https://spark.apache.org/mllib/>
- [6] Apache Spark. (2024). Evaluation Metrics - MLlib. <https://spark.apache.org/docs/latest/ml/lib-evaluation-metrics.html>
- [7] Apache Spark. (2024). Evaluation Metrics - MLlib. <https://spark.apache.org/docs/latest/ml/lib-evaluation-metrics.html>
- [8] Apache Spark. (2024). PySpark API Documentation. <https://spark.apache.org/docs/latest/api/python/>
- [9] Apache Kafka. (2024). Kafka Python Client. <https://github.com/dpkp/kafka-python>