# Report: Tweet Sentiment Analysis

## 1. Problem Statement:

The goal of this project is to build a sentiment analysis model that can classify tweets into four categories:

- **Positive**
- **Negative**
- **Neutral**
- **Irrelevant**

Given the vast amount of opinions and feedback shared on social media, analyzing sentiment in tweets can help businesses and organizations gauge public opinion, customer satisfaction, and sentiment trends.

## 2. Data Collection:

The dataset used for this project consists of tweets that are categorized into the mentioned sentiment labels. Each tweet is represented as a text document that requires preprocessing and feature extraction for analysis.

Since this project focuses on sentiment analysis, we assume the dataset contains tweets labeled with their respective sentiment (positive, negative, neutral, irrelevant).

## 3. Approach:

### 3.1 Preprocessing:

Text data often requires significant preprocessing before it can be used for machine learning models. The following steps were taken:

1. **Lowercasing**: All text was converted to lowercase to ensure uniformity.
2. **Tokenization**: The text was split into individual words (tokens).
3. **Removing Stop Words**: Common words (e.g., "the", "is", "in") that do not contribute to sentiment analysis were removed.
4. **Stemming/Lemmatization**: Words were reduced to their base form (e.g., "running" -> "run").
5. **Removing Special Characters**: Punctuation, special characters, and URLs were removed.

### 3.2 Feature Extraction:

To convert the textual data into a format that can be fed into machine learning algorithms, **TF-IDF (Term Frequency-Inverse Document Frequency)** was used as a feature extraction method.

- **TF-IDF** helps capture the importance of each word in the context of the entire corpus, giving higher weights to words that are frequent in a specific document but less frequent across the entire dataset.
- The vectorizer was trained using the tweet data to create numerical representations of the tweets.

### 3.3 Model Selection:

**Logistic Regression** was chosen as the model for classification due to its simplicity, interpretability, and good performance in binary and multi-class classification tasks.

- Logistic regression is a linear model that predicts probabilities for each class. In the case of multi-class classification, the model uses a one-vs-rest strategy to classify tweets into one of the four sentiment categories.

### 3.4 Evaluation Metric:

The model's performance was evaluated using **accuracy**, which measures the percentage of correctly classified tweets in the dataset. Accuracy is chosen here as the primary metric due to its simplicity and the balanced nature of the dataset.

## 4. Model Performance:

- **Accuracy**: The model achieved an accuracy of **77%** on the test dataset. This means that the model correctly classified 77% of the tweets in the dataset.

The accuracy metric indicates that while the model performs reasonably well, there is still room for improvement, especially in handling tweets that may fall under the "Neutral" or "Irrelevant" categories, which can be more challenging to classify.

## 5. Challenges:

Several challenges were encountered during the development of the model:

1. **Ambiguity in Sentiment**: Some tweets are inherently ambiguous, such as sarcasm, which poses a challenge for traditional sentiment analysis models.
2. **Irrelevant Tweets**: Tweets that don't convey clear sentiment (e.g., asking a question or posting an image) were harder for the model to classify, affecting accuracy.
3. **Preprocessing Variations**: Variations in spelling, abbreviations, and slang in tweets required robust preprocessing steps to improve the model's effectiveness.

## 6. Results and Discussion:

The model was able to classify tweets with decent accuracy, but some areas for improvement include:

- **Handling sarcasm and humor**: These forms of language are difficult for any machine learning model to understand.
- **Improving handling of neutral tweets**: Neutral tweets, which contain neither positive nor negative sentiment, might require a separate class or additional features to be more accurately classified.

## 7. Future Work:

To improve the performance of this model, several strategies could be explored:

1. **Deep Learning Models**: Transitioning to more advanced models like **LSTM (Long Short-Term Memory)** or **BERT** (Bidirectional Encoder Representations from Transformers) could improve sentiment understanding by capturing contextual relationships.
2. **Handling Sarcasm**: Incorporating more sophisticated text representations or training on datasets that include sarcasm could improve the model's handling of ambiguous sentiment.
3. **Hyperparameter Tuning**: Further tuning of the logistic regression model's hyperparameters, such as regularization strength and solver options, could improve its generalization ability.

## 8. Conclusion:

This project demonstrates how machine learning can be applied to classify tweets into sentiment categories. By using **TF-IDF** for feature extraction and **Logistic Regression** for classification, the model achieved an accuracy of 77%. While the model works well for clear sentiment tweets, it could benefit from further refinement, especially for ambiguous or neutral tweets.