# NLP-Driven Fake Review Detection System

## 1. Project Objective

This project aimed to construct an **NLP-driven classification framework** to identify fake reviews by categorizing them as either:

- **Computer Generated (CG)**

- **Original (OR)**

The system integrates advanced natural language processing (NLP) techniques with machine learning to preprocess textual data, derive discriminative features, and develop an accurate predictive model.

## 2. Dataset Overview

### Dataset Description

- **Dataset File:** fake reviews dataset.csv

- **Columns:**

  - text_: Contains the review text subjected to classification.

  - label: Binary target variable with classes CG (Computer Generated) and OR (Original).

  - category: Dropped during preprocessing due to irrelevance to the analytical objective.

### Data Preprocessing

- The category column was removed to streamline the dataset.

- Distribution of values in the label column was analyzed for class balance.

## 3. Text Preprocessing

A comprehensive preprocessing pipeline was implemented to refine the raw text data for vectorization.

### Methodology

1. **Case Normalization:** Text was converted to lowercase for uniformity.

2. **Punctuation Elimination:** Punctuation marks were stripped using Python's string.punctuation module.

3. **Tokenization:** Reviews were segmented into individual tokens using the word_tokenize method from NLTK.

4. **Stopword Removal:** Non-informative words (e.g., "and," "the," "is") were excluded using NLTK's predefined English stopword list.

5. **Lemmatization:** Tokens were reduced to their canonical forms using the WordNetLemmatizer.

**Automation**

A custom function preprocess encapsulated these steps to ensure efficient and consistent text processing. This function outputs clean, tokenized, and lemmatized text.

---

## 4. Feature Extraction

The textual data was transformed into a numerical representation using the **TF-IDF Vectorizer**.

**TF-IDF Methodology**

- **TF-IDF (Term Frequency-Inverse Document Frequency):** Quantifies the importance of terms in individual documents relative to the entire corpus.

- Implemented using Scikit-learn's TfidfVectorizer.

- **Output:** A sparse matrix representation with rows corresponding to documents and columns representing unique terms.

**Dimensional Characteristics**

- The resultant matrix had a shape of (number_of_documents, number_of_unique_terms).

---

## 5. Data Partitioning

The dataset was split into training and testing subsets for model training and evaluation.

**Specifications**

- **Split Ratio:** 80% training and 20% testing.

- **Random State:** Fixed at 42 for reproducibility.

- **Library:** Scikit-learn's train_test_split.

**Label Encoding**

The label column was binarized as follows:

- CG → 0 (Computer Generated)

- OR → 1 (Original)

---

## 6. Model Development and Training

**Chosen Model: Logistic Regression**

- **Rationale:** Logistic Regression is well-suited for binary classification due to its simplicity, interpretability, and computational efficiency.

- **Implementation:**
    - Employed Scikit-learn's LogisticRegression class.
    - Trained using the TF-IDF-transformed x_train and corresponding y_train labels.

---

## 7. Model Evaluation

**Performance Metrics**

Evaluation was conducted on the held-out x_test set using standard classification metrics:

- **Precision:** Accuracy of positive class predictions.
- **Recall:** Sensitivity in detecting relevant instances.
- **F1-Score:** Harmonic mean of precision and recall.
- **Accuracy:** Overall percentage of correct predictions.

**Results**

The classification report detailed robust performance across both classes (0 and 1), demonstrating the model's efficacy in differentiating computer-generated and authentic reviews.

---

## 8. Fake Review Prediction Functionality

A utility function fake_pred was designed for real-time review classification. The workflow includes:

1. Preprocessing the input text using the preprocess function.
2. Vectorizing the cleaned text via the TF-IDF model.
3. Predicting the class using the trained logistic regression model.

**Example Predictions**

1. **Input:**

"The wireless Bluetooth headphones offer superior sound quality and a seamless connection."
**Output:** Computer Generated Review

2. **Input:**

"I recently purchased the XYZ Mobile and it has exceeded my expectations in every way." **Output:** Original Review

3. **Input:**

"The iPhone 14 is a top-tier smartphone that combines sleek design with powerful performance."
**Output:** Computer Generated Review

---

## 9. Observations and Insights

**Strengths**

1. **High Predictive Accuracy:** The model demonstrated commendable precision and recall on the test dataset.

2. **Effective Preprocessing Pipeline:** The cleaning and tokenization workflow ensured the text data was primed for feature extraction.

3. **User-Friendly Application:** The fake_pred function simplifies interaction and enables practical deployment.

**Limitations**

1. **Dataset Dependency:** Model performance is inherently tied to the dataset's quality and diversity.

2. **Model Simplicity:** Logistic Regression may fail to capture nuanced relationships in highly complex datasets.

---

## 10. Future Enhancements

1. **Model Diversification:**

   o Experiment with advanced algorithms like Random Forest, Gradient Boosting, or neural networks.

   o Integrate pretrained NLP models (e.g., BERT, GPT) for superior contextual understanding.

2. **Feature Enrichment:**

   o Explore n-grams and embeddings (e.g., Word2Vec, FastText) to enhance feature representation.

3. **Data Augmentation:**

   o Expand the dataset with additional samples and leverage synthetic techniques to balance class distributions.

4. **Model Explainability:**

   o Use frameworks like SHAP or LIME to interpret and explain model predictions effectively.

5. **Real-Time Deployment:**

   o Package the model into an API or integrate it into a web application for real-world usability.

---

## 11. Conclusion

This project successfully implemented a full-stack NLP pipeline to address the challenge of fake review detection. The model's ability to differentiate between computer-generated and authentic reviews provides a valuable resource for maintaining trust on online platforms. With robust preprocessing and a straightforward logistic regression model, this study establishes a strong foundation for future advancements in automated review analysis.