

<b>Practical Assignment No. 5</b>	
Title:	<b>Text preprocessing</b>
Problem Statement:	Build sentiment analysis model using techniques like Bag-of Words, TF-IDF, or word embedding (Word2Vec, GloVe).
Objective:	To develop applications for text processing using natural language processing techniques.
Outcome:	CO606.2: Develop text processing application using natural language processing techniques.
Software or Hardware Requirements:	Anaconda/Java/GCC
Theory:	<p><b>Sentiment Analysis:</b></p> <p>Sentiment Analysis (also known as Opinion Mining) is a natural language processing (NLP) technique used to determine the emotional tone behind a body of text. It helps in identifying whether the expressed opinion in a text is positive, negative, or neutral.</p> <p><b>It is widely used in:</b></p> <ul style="list-style-type: none"> <li>• Social media monitoring</li> <li>• Product and movie reviews</li> <li>• Customer feedback systems</li> <li>• Market analysis and brand monitoring</li> </ul> <p><b>The process involves:</b></p> <ol style="list-style-type: none"> <li>1. Text preprocessing</li> <li>2. Feature extraction using Bag-of-Words, TF-IDF, or word embeddings</li> <li>3. Model training and classification using machine learning algorithms like Naïve Bayes, Logistic Regression, or SVM</li> </ol>

## **Steps in Sentiment Analysis**

### **1. Data Collection:**

Text data is collected from sources such as Twitter, movie reviews, or product feedback.

### **2. Text Preprocessing:**

Cleaning text by removing special characters, punctuation, stop-words, and performing tokenization, stemming, and lemmatization.

### **3. Feature Extraction:**

Converting text into numerical form using one of the following methods:

- Bag-of-Words
- TF-IDF
- Word Embedding (Word2Vec / GloVe)

### **4. Model Building:**

Training a machine learning classifier on the vectorized data.

### **5. Performance Evaluation:**

Evaluating the model using metrics such as accuracy, precision, recall, and F1-score.

## **Feature Extraction Techniques**

### **(a) Bag-of-Words (BoW) Model**

#### **Definition:**

Bag-of-Words represents a document as a collection (or "bag") of individual words, disregarding grammar and word order but keeping word frequency.

#### **Process:**

1. Tokenize the text into words
2. Build a vocabulary of all unique words
3. Represent each document as a vector showing the frequency of each word

**Example:**

Sentence	Bag-of-Words Representation
“I love Python”	[0, 1, 1, 0]
“Python is great”	[0, 0, 1, 1]

Vocabulary: [‘I’, ‘love’, ‘Python’, ‘great’]

**Advantages:**

- Simple and easy to implement
- Works well for small and medium-sized datasets

**Disadvantages:**

- Ignores context and word order
- Large feature space with sparse vectors

**Python Example:**

```
from sklearn.feature_extraction.text import CountVectorizer  
vectorizer = CountVectorizer()  
X = vectorizer.fit_transform(corpus)
```

**(b) TF-IDF (Term Frequency - Inverse Document Frequency)****Definition:**

TF-IDF improves upon BoW by reducing the importance of frequently occurring words (like “the”, “is”) and giving higher weight to rarer, more informative words.

**Formula:**

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \times \text{IDF}(t)$$

**Where,**

- $\text{TF}(t, d)$  = Number of times term  $t$  appears in document  $d$

- $IDF(t) = \log (\text{Total number of documents} / \text{Number of documents containing } t)$

#### **Advantages:**

- Captures importance of words
- Reduces impact of common terms

#### **Disadvantages:**

- Still ignores word order and context
- Produces sparse matrices for large datasets

#### **Python Example:**

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer()
X = vectorizer.fit_transform(corpus)
```

### **(c) Word Embeddings**

#### **Definition:**

Word Embedding is a dense vector representation of words that captures their semantic relationships and context.

Unlike BoW or TF-IDF, it understands that similar words have similar meanings.

#### **Common Embedding Techniques:**

##### **1. Word2Vec**

- Developed by Google.
- Based on Neural Networks using either CBOW (Continuous Bag of Words) or Skip-Gram model.
- Captures contextual relationships between words.

#### **Example:**

The model learns that “king” – “man” + “woman”  $\approx$  “queen”.

### **Python Example:**

```
from gensim.models import Word2Vec  
model = Word2Vec(sentences, vector_size=100, window=5, min_count=1,  
sg=0)
```

## **2. GloVe (Global Vectors for Word Representation)**

- Developed by Stanford NLP Group.
- Combines the benefits of global matrix factorization and local context window methods.
- Learns word meanings based on co-occurrence statistics.

### **Mathematical Idea:**

Words appearing together frequently have similar vector representations.

### **Example:**

```
from gensim.models import KeyedVectors  
model = KeyedVectors.load_word2vec_format('glove.6B.100d.txt',  
binary=False)
```

### **Advantages:**

- Captures semantic meaning and context
- Words with similar meanings have closer vector representations
- Suitable for deep learning models

### **Disadvantages:**

- Requires large data and training time
- More complex compared to BoW or TF-IDF

### **Model Building**

	<p>Once text is converted into numerical form using one of the above techniques, various machine learning algorithms can be applied:</p> <ul style="list-style-type: none"> <li>● Naïve Bayes Classifier</li> <li>● Logistic Regression</li> <li>● Support Vector Machine (SVM)</li> <li>● Random Forest</li> <li>● Deep Learning models (e.g., LSTM, CNN)</li> </ul> <p>The trained model predicts sentiment labels like Positive, Negative, or Neutral.</p> <h3>Evaluation Metrics</h3> <p>Performance of the sentiment analysis model is evaluated using:</p> <ul style="list-style-type: none"> <li>● Accuracy = <math>(TP + TN) / (TP + TN + FP + FN)</math></li> <li>● Precision = <math>TP / (TP + FP)</math></li> <li>● Recall = <math>TP / (TP + FN)</math></li> <li>● F1-Score = <math>2 \times (Precision \times Recall) / (Precision + Recall)</math></li> <li>● Confusion Matrix to visualize performance</li> </ul> <h3>Applications</h3> <ul style="list-style-type: none"> <li>● Product review and movie review analysis</li> <li>● Social media sentiment monitoring</li> <li>● Customer service and chatbot systems</li> <li>● Financial market prediction based on news sentiment</li> </ul>
Input/Datasets/Test Cases:	<p>Dataset- Name of the Dataset:          Description of the Dataset:          Dataset Characteristics:</p>

	<p>Subject Area:            Associated Tasks:            Feature Type:            # Instances:            # Features:</p>
Results:	Execute code to use any one standard dataset and any NLP technique for sentiment analysis. Take a print of this code with output for submission as part of results.
Analysis and conclusion:	Write your own analysis of output and conclusion( Minimum 1 statement of Analysis, Minimum 1 Statement Conclusion)
References:	Reference /Links(min Any 2, include dataset ref.). Write references in IEEE format.