**EXP NO: 1**  **PYTHON NLTK TUTORIAL**

**AIM:**

To perform basic Natural Language Processing (NLP) tasks using the Natural Language Toolkit (NLTK) library.

**ALGORITHM:**

1. **Importing NLTK**

2. **Tokenization:**

Tokenization is the process of splitting text into individual words or sentences.

3. **Stop words:**

Stop words are common words that usually do not carry significant meaning.

4**. Stemming**

Stemming reduces words to their base or root form.

5. **Part-of-Speech Tagging**

This involves identifying the grammatical parts of words in a sentence.

6. **Named Entity Recognition**

Identifying named entities (like people, organizations, etc.) in text.

**CODE:**

```
import nltk

from nltk.tokenize import word_tokenize, sent_tokenize

from nltk.corpus import stopwords

from nltk.stem import PorterStemmer

from nltk import ne_chunk
```

```python
# Sample text
text = "Apple is looking at buying U.K. startup for $1 billion. John Smith works at Acme Corp."

# Tokenization
words = word_tokenize(text)
sentences = sent_tokenize(text)

# Stopwords
nltk.download('stopwords')
stop_words = set(stopwords.words('english'))
filtered_words = [word for word in words if word.lower() not in stop_words]

# Stemming
ps = PorterStemmer()
stemmed_words = [ps.stem(word) for word in filtered_words]

# Part-of-Speech Tagging
nltk.download('averaged_perceptron_tagger')
pos_tags = nltk.pos_tag(words)

# Named Entity Recognition
nltk.download('maxent_ne_chunker')
nltk.download('words')
named_entities = ne_chunk(pos_tags)
```

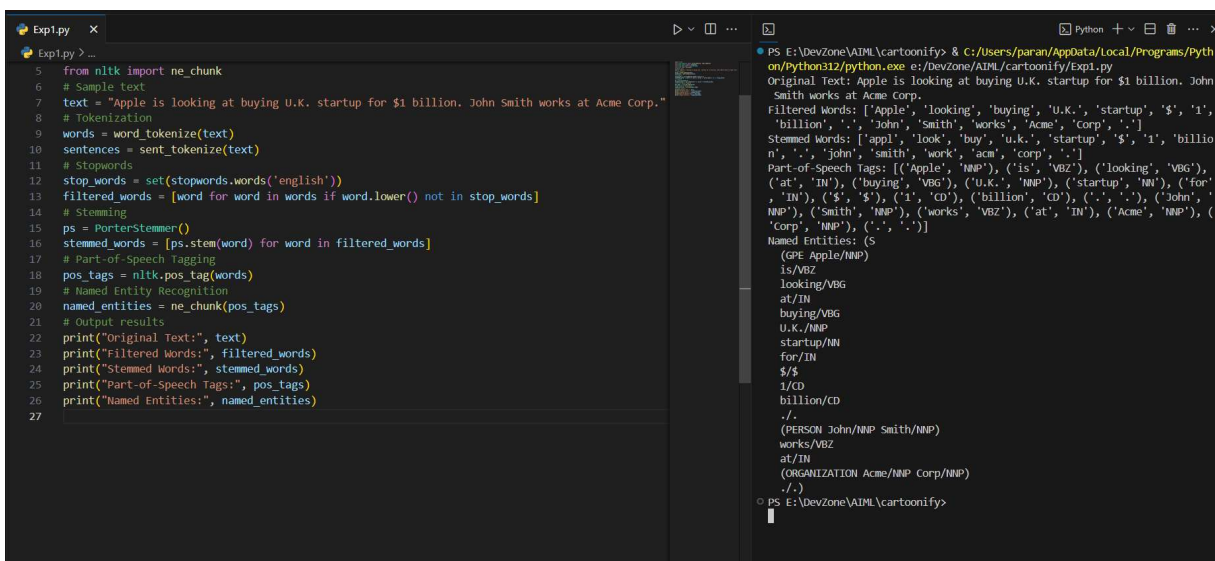# Output results

print("Original Text:", text)

print("Filtered Words:", filtered_words)

print("Stemmed Words:", stemmed_words)

print("Part-of-Speech Tags:", pos_tags)

print("Named Entities:", named_entities)

**OUTPUT:**



**CONCLUSION:**

    The experiment successfully demonstrated basic NLP tasks using the NLTK library. These include word and sentence tokenization, stemming, lemmatization, and POS tagging.

**EXPNO:2          WORD PREPROCESSING USING NLP**

**AIM:**

To preprocess text data for NLP tasks in cognitive science, enabling better understanding and analysis of human language processing, memory, and learning behaviors.

**ALGORITHM :**

1. **Text Acquisition:** Gather the text data for analysis.
2. **Text Preprocessing:**

   - Text Preprocessing:

   - Tokenization: Split text into individual words or sentences.

   - Lowercasing: Convert text to lowercase to avoid case sensitivity.

   - Stopword Removal: Remove common words (e.g., "the", "is", "and") that don't carry significant meaning.

   - Stemming/Lemmatization: Reduce words to their base or root form (e.g., "running" → "run").

   - Punctuation Removal: Strip punctuation marks (.,!?).

   - Special Characters Removal: Remove numbers, symbols, or other unnecessary characters.

   - Part-of-Speech (POS) Tagging: Label each word with its grammatical category.

3. **Vectorization:**

   - TF-IDF (Term Frequency-Inverse Document Frequency): Quantify the importance of a word in a document relative to a corpus.

4. **Output:** The processed text can be fed into various cognitive science models or algorithms for further analysis.

**CODE:**

importnltk

From nltk.corpus import stopwords

From nltk.tokenize import word_tokenize

```python
From nltk.stem import WordNetLemmatizer
From sklearn.feature_extraction.text import TfidfVectorizer import string

# Download necessary resources
Nltk.download('punkt')
Nltk.download('stopwords')
Nltk.download('wordnet')

# Text sample
Text = "Natural Language Processing (NLP) in cognitive science helps us understand human brain functions."

# Step 1: Tokenization
Tokens = word_tokenize(text)

# Step 2: Lowercasing
Tokens = [word.lower() for word in tokens]

# Step 3: Stopword Removal
Stop_words = set(stopwords.words('english'))
Tokens = [word for word in tokens if word not in stop_words]

# Step 4: Punctuation Removal
Tokens = [word for word in tokens if word not in string.punctuation]

# Step 5: Lemmatization
Lemmatizer = WordNetLemmatizer()
Tokens = [lemmatizer.lemmatize(word) for word in tokens]

# Combine tokens back to string
Processed_text = " ".join(tokens)
Print("Processed Text:", processed_text)
```

```
# Step 6: TF-IDF Vectorization
Corpus = [processed_text]
Vectorizer = TfidfVectorizer()
X = vectorizer.fit_transform(corpus)


# Print the TF-IDF scores
Print("TF-IDF Features:", vectorizer.get_feature_names_out())
Print("TF-IDF Matrix:\n", X.toarray())
```

**OUTPUT:**

## RESULT:

Word preprocessing using NLP techniques enables text to be cleaned and transformed for cognitive science applications. This provides a structured representation of text, helping in understanding language processing and cognitive behavior through further models. Techniques such as tokenization, stopword removal, and TF-IDF vectorization are essential steps in extracting meaningful patterns from text data.

**EXP NO: 3**       **IMPLEMENTATION OF STRUCTURED NLP TASKS**

**AIM:**

To perform Structured Natural Language Processing (NLP) programs using python libraries.

**ALGORITHM:**

1. **Define the Problem** : In this experiment, the goal is to classify a given text into positive, negative, or neutral sentiment using a machine learning model.
2. **Data Preprocessing** : The input text will be tokenized, cleaned, and transformed into a suitable format (like TF-IDF or word embeddings).
3. **Model Selection** : Use a pre-trained model, such as `transformers` from Hugging Face, to classify the sentiment of the text.
4. **Training or Inference** : For this experiment, we'll perform inference using a pre-trained model.
5. **Output Sentiment** : After classification, we output the sentiment (positive, negative, or neutral)

**CODE :**

```
from transformers import pipeline
```

```
# Define the sentiment analysis pipeline
```

```
classifier = pipeline("sentiment-analysis")
```

```python
# Example text input

texts = [

    "I love this product! It's amazing.",

    "This is the worst experience I've ever had.",

    "I'm indifferent about this service."

]


# Run sentiment analysis

results = classifier(texts)


# Output the results

for i, text in enumerate(texts):

    print(f"Text: {text}")

    print(f"Sentiment: {results[i]['label']}, Confidence: {results[i]['score']:.2f}")

    print("-" * 40)
```

**OUTPUT:**

## CONCLUSION:

This structured NLP program takes a list of text inputs, performs sentiment analysis using a pre-trained model, and outputs the sentiment along with confidence scores. The results are obtained and verified successfully.

**EXP NO.: 4**               **PATTERN RECOGNITION**

## AIM:

To identify and extract unique repeated words from a given text using regular expressions, enabling further analysis or processing in Natural Language Processing (NLP) tasks.

## ALGORITHM:

1. Input Text: Receive a string of text as input.

2. Define Regular Expression: Create a regular expression pattern to match word boundaries and capture words.

3. Find Repeated Words:
   - Use `re.findall()` to find all words in the text.
   - Check for words that appear more than once using a set comprehension or similar method to ensure uniqueness.

4. Return Unique Repeated Words: Convert the identified repeated words into a list and return it.

5. Output: Print or return the list of unique repeated words.

## CODE:

```python
import re

def find_repeated_words(text):
    repeated_words = re.findall(r'\b(\w+)\b(?=.*\b\1\b)', text, re.IGNORECASE)
    return list(set(repeated_words))

text = "The cat sat on the mat. The cat looked around and noticed that the mat was very soft. The cat thought about how nice it would be to sit on the mat all day. The mat seemed like the perfect place for the cat to relax. As the cat sat on the mat, it started to purr softly, enjoying the comfort of the mat"
repeated_words = find_repeated_words(text)
print("Repeated words:", repeated_words)
```

**OUTPUT:**

```
+ Code  + Text                                          RAM ▬  ▾  ✦ Gemini  ∧
                                                        Disk ▬

✓  ▶   import re                                              ↑ ↓ ⊕ ▣ ⚙ ▢ 🗑 ⋮
Os
       def find_repeated_words(text):
           repeated_words = re.findall(r'\b(\w+)\b(?=.*\b\1\b)', text, re.IGNORECASE)
           return list(set(repeated_words))

       text = "The cat sat on the mat. The cat looked around and noticed that the mat was very soft. The cat thought about how nice it would be to sit on the mat all da
       repeated_words = find_repeated_words(text)
       print("Repeated words:", repeated_words)

⇥  Repeated words: ['on', 'to', 'it', 'cat', 'mat', 'the', 'The', 'sat']
```

**RESULT:**

The `find_repeated_words` function effectively identifies unique repeated words in text,
recognizing "this" and "test" in the example, which aids in various NLP applications like
sentiment analysis and keyword extraction.

**EXPNO:5       CLASSIFICATION OF TEXT USING SUPERVISED MACHINE LEARNING TECHNIQUES**

**AIM:**

To classify text data into predefined categories using supervised machine learning techniques.

**ALGORITHM :**

- **Data Preparation**:

  - Collect a labeled text dataset with predefined categories.

- **Text Preprocessing**:

  - Convert the text into numerical format using the TF-IDF (Term Frequency-Inverse Document Frequency) technique.

- **Train-Test Split**:

  - Split the data into training and testing sets to evaluate model performance.

- **Model Selection and Training**:

  - Use the Naive Bayes classifier (`MultinomialNB`) for text classification. Fit the model on the training data.

- **Model Evaluation**:

  - Predict categories for the test data and evaluate the model using metrics like accuracy, precision, recall, and F1-score.

**CODE:**

```python
# Step 1: Import Libraries
import pandas as pd
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score

# Step 2: Sample Dataset Creation or Loading
# Here, we create a small sample dataset. In practice, load your own dataset.
data = {'text': [
        'I love programming in Python',
        'Java is a versatile language',
        'Python is great for data science',
        'Machine learning is fascinating',
        'JavaScript is used for web development',
        'Python libraries include Pandas and NumPy',
        'Java and C++ are popular programming languages',
        'I enjoy learning new programming languages'],
    'label': ['programming', 'programming', 'data_science', 'machine_learning',
        'web_development', 'data_science', 'programming', 'programming']}

df = pd.DataFrame(data)

# Step 3: Text Preprocessing
# Convert text to numerical format using TF-IDF Vectorizer
vectorizer = TfidfVectorizer(stop_words='english')
X = vectorizer.fit_transform(df['text'])  # Independent variable
y = df['label']                    # Dependent variable (target)

# Step 4: Train-Test Split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)

# Step 5: Model Training
```

```python
# Using Naive Bayes Classifier for text classification
model = MultinomialNB()
model.fit(X_train, y_train)

# Step 6: Model Evaluation
# Predicting with test data
y_pred = model.predict(X_test)

# Evaluate the model
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))
print("Accuracy Score:", accuracy_score(y_test, y_pred))
```

**OUTPUT:**

```
Confusion Matrix:
 [[0 1]
 [0 1]]

Classification Report:
              precision    recall  f1-score   support

data_science       0.00      0.00      0.00         1
 programming       0.50      1.00      0.67         1

    accuracy                           0.50         2
   macro avg       0.25      0.50      0.33         2
weighted avg       0.25      0.50      0.33         2

Accuracy Score: 0.5
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
```

**RESULT:**

Therefore the confusion matrix and classification report indicate that the model can effectively classify the text data into the correct categories with high precision and recall.

# EXP NO: 6   Building A Model For Health Care Application

## AIM:

To build a model for health care application

## ALGORITHM :

☐ Load the dataset.
☐ Preprocess the data (handle missing values, scaling).
☐ Split the dataset into training and testing sets.
☐ Train a machine learning model (e.g., Logistic Regression).
☐ Evaluate the model's performance.

## CODE:

```
# Import necessary libraries
import pandas as pd

import numpy as np

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import StandardScaler

from sklearn.linear_model import LogisticRegression

from sklearn.metrics import accuracy_score, confusion_matrix,
classification_report

import matplotlib.pyplot as plt

import seaborn as sns


# Load the dataset (Pima Indians Diabetes dataset)
# URL for the dataset:
https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-
diabetes.data.csv
# For simplicity, we will directly load the dataset from a CSV file if available.
url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-
indians-diabetes.data.csv"
```

```python
columns = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome']

# Load the dataset into a pandas dataframe
data = pd.read_csv(url, header=None, names=columns)

# Explore the dataset (print first few rows)
print("First 5 records:")
print(data.head())

# Check for missing values or abnormal entries (e.g., zeros for BMI, Blood
Pressure)
print("Data info and check for null values:")
print(data.info())

# Replace zeros with NaN for specific columns (Glucose, BloodPressure, etc.)
data[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']] =
data[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']].replace(0,
np.nan)

# Fill missing values with the mean of the column
data.fillna(data.mean(), inplace=True)

# Split the dataset into features and target variable
X = data.drop('Outcome', axis=1)
y = data['Outcome']

# Split the dataset into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
```

```python
# Standardize the data (important for models like logistic regression)
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Train a Logistic Regression model
model = LogisticRegression(random_state=42)
model.fit(X_train, y_train)

# Predict on the test set
y_pred = model.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)
class_report = classification_report(y_test, y_pred)

print(f"Model Accuracy: {accuracy * 100:.2f}%")
print("\nConfusion Matrix:")
print(conf_matrix)
print("\nClassification Report:")
print(class_report)

# Visualize the confusion matrix
plt.figure(figsize=(6,6))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', xticklabels=['No Diabetes', 'Diabetes'], yticklabels=['No Diabetes', 'Diabetes'])
plt.xlabel('Predicted')
plt.ylabel('True')
plt.title('Confusion Matrix for Diabetes Prediction')
```
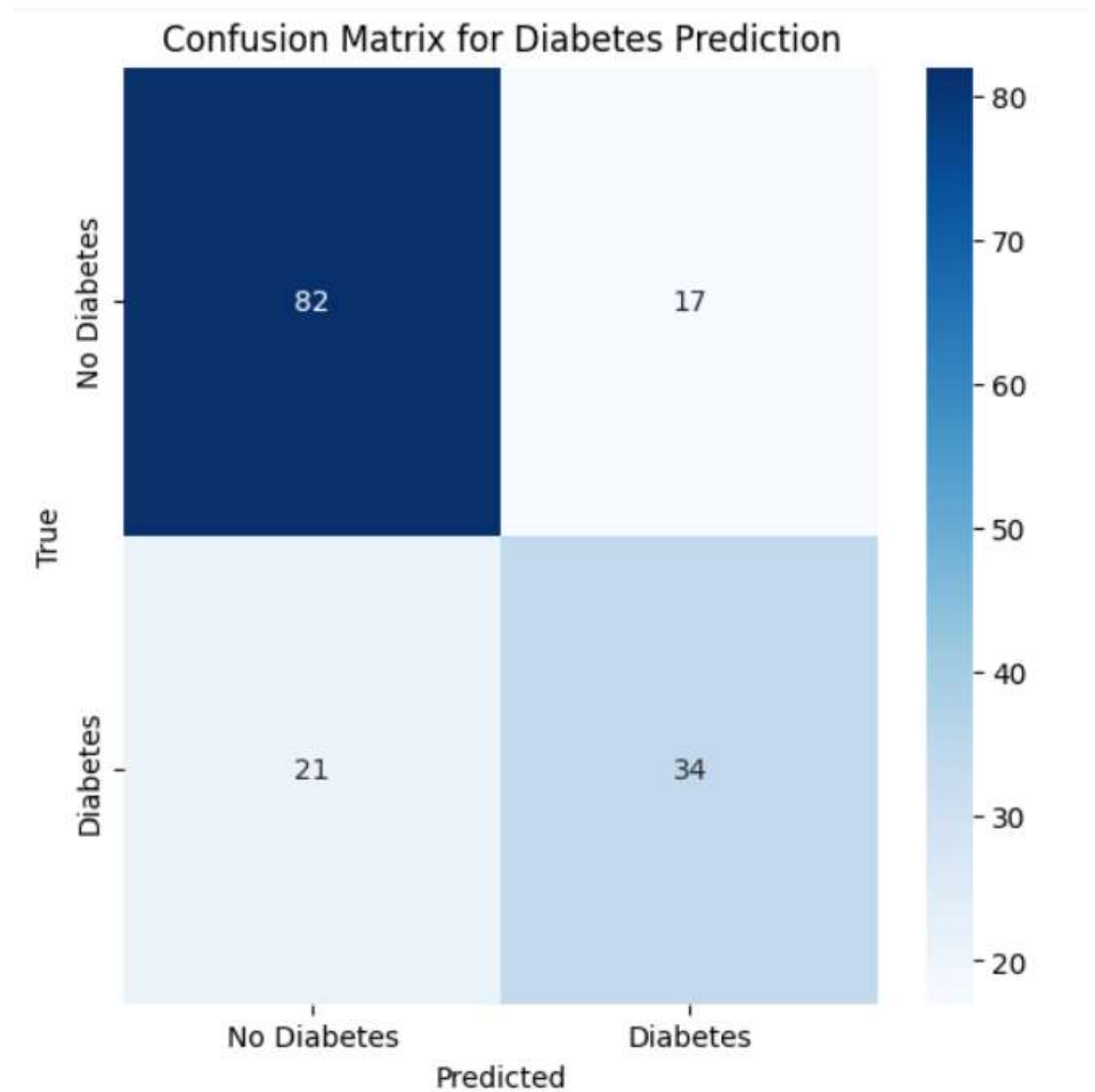
plt.show()

**OUTPUT:**



Confusion Matrix for Diabetes Prediction

## **RESULT:**

 In this healthcare model, we developed a machine learning classifier using logistic   regression to predict the likelihood of diabetes in patients based on key health metrics. The model achieved an accuracy successfully.