

AM604PC- Natural Language Processing Laboratory Manual

[ANURAG ENGINEERING COLLEGE]

Contents

1	Tokenization	4
2	Stop Word Removal	7
3	Word Analysis	9
4	Word Generation	11
5	Word Sense Disambiguation	14
6	Install NLTK Toolkit and Perform Stemming	18
7	POS Tagging	21
8	Morphological Analysis and N-Grams	24
9	Implementation of Porter Stemmer Algorithm for Stemming	28
10	Conversion of Audio to Text and Text to Audio	30

NLTK: Natural Language Toolkit

The Natural Language Toolkit (NLTK) is a powerful Python library widely used in Natural Language Processing (NLP). It provides tools for text processing, corpus handling, linguistic analysis, and statistical modeling. NLTK is designed to be user-friendly, making it a popular choice for educational and research purposes in NLP.

Features of NLTK

Text Processing Tools

- **Tokenization:** NLTK can divide text into smaller units, such as words or sentences.
- **Stemming:** This feature reduces words to their root form, such as *running* to *run*.
- **Lemmatization:** Unlike stemming, lemmatization returns words to their dictionary form while considering their context.

Part-of-Speech (POS) Tagging

NLTK assigns grammatical categories to words in a sentence, such as nouns, verbs, and adjectives, aiding syntactic and semantic analysis.

Named Entity Recognition (NER)

This feature identifies entities such as persons, locations, and organizations within text. For example, in the sentence “*Steve Jobs founded Apple,*” Steve Jobs is tagged as a *Person*, and Apple as an *Organization*.

Syntactic Parsing

NLTK provides tools to analyze the grammatical structure of sentences using Context-Free Grammar (CFG) and Dependency Parsing.

Corpus Support

NLTK includes access to over 50 preloaded corpora such as the Brown Corpus and the Gutenberg Corpus, enabling easy experimentation with real-world text data. In *Natural Language Processing (NLP)*, *corpus support refers to the ability of tools, frameworks, or libraries to provide access to collections of text data (known as corpora) for analysis, research, and application development. A corpus is a structured set of text, often annotated with metadata, that serves as a foundational resource for NLP tasks.*

WordNet Interface

WordNet is a semantic lexicon that allows exploration of relationships between words, such as synonyms, antonyms, hypernyms, and hyponyms.

Text Classification

NLTK supports text classification with both supervised and unsupervised learning methods, including algorithms such as Naive Bayes.

Collocations and Bigrams

The library can detect commonly co-occurring words, such as *New York*, using collocation and bigram techniques.

Language Modeling

It provides functionality for building unigram and bigram models for language analysis and generation.

Visualization

Visualization tools in NLTK allow for displaying parse trees and analyzing word frequency distributions.

Applications of NLTK

- **Sentiment Analysis:** Analyze the polarity of reviews or feedback.
- **Document Classification:** Categorize text documents, such as spam detection.
- **Grammar and Syntax:** Explore and implement parsing techniques.
- **Chatbot Development:** Build chatbots with tokenization, lemmatization, and classification.

NLTK is an indispensable tool for Natural Language Processing, offering a comprehensive suite of text processing, linguistic analysis, and visualization capabilities. Its ease of use and versatility make it an excellent choice for both educational and research purposes.

Tokenization

Program-1

Write a Python Program to perform Tokenization tasks on text

Tokenization in NLP

Tokenization is the process of breaking down a text into smaller units called *tokens*. These tokens can be words, phrases, or characters, depending on the application. For example, the sentence *"Natural Language Processing is fun"* can be tokenized into individual words: [*Natural*, *Language*, *Processing*, *is*, *fun*].

Following are the uses of Tokenization:

- **Text Representation:** Tokenization is a fundamental step in preparing text data for computational processing. Most NLP algorithms and models require input in the form of tokens.
- **Feature Extraction:** Tokens serve as features for models, enabling the representation of text in numerical forms such as vectors for further processing.
- **Context Preservation:** Tokenization ensures that the structure and meaning of text are preserved to an extent while breaking it into manageable parts.
- **Language Understanding:** It enables tasks like syntactic parsing, sentiment analysis, machine translation, and more by providing a manageable representation of the text.
- **Efficiency:** Working with tokens rather than raw text makes computations faster and more efficient, as the raw text might contain noise, punctuation, and other elements that could hinder analysis.

Tokenization is particularly useful in tasks such as search engines, chatbots, language modeling, and more, as it allows the system to process and understand language in a structured way.

Python Code

```
1      import nltk
2      from nltk.tokenize import word_tokenize, sent_tokenize
3
4      # Download NLTK resources if not already downloaded
5      nltk.download('punkt')
6
7      def perform_tokenization(text):
8          """
9          Performs sentence and word tokenization on the input text.
10
11          Args:
12          text (str): The input text to be tokenized.
13
```

```

14     Returns:
15     tuple: A tuple containing a list of sentences and a list of words.
16     """
17     # Sentence tokenization
18     sentences = sent_tokenize(text)
19
20     # Word tokenization
21     words = word_tokenize(text)
22
23     return sentences, words
24
25     def display_tokenization_results(sentences, words):
26         """
27         Displays the results of sentence and word tokenization.
28
29         Args:
30         sentences (list): List of tokenized sentences.
31         words (list): List of tokenized words.
32         """
33         print("\nSentence Tokenization:")
34         for i, sentence in enumerate(sentences, start=1):
35             print(f"Sentence {i}: {sentence}")
36
37         print("\nWord Tokenization:")
38         print(words)
39
40     def main():
41         """
42         Main function to demonstrate tokenization.
43         """
44         # Sample text for demonstration
45         text = """Natural Language Processing is an amazing field of
46                 Artificial Intelligence.
47                 It enables machines to understand human language effectively and
48                 accurately."""
49
50         print("Original Text:")
51         print(text)
52
53         # Perform tokenization
54         sentences, words = perform_tokenization(text)
55
56         # Display tokenization results
57         display_tokenization_results(sentences, words)
58
59         if __name__ == "__main__":
60             main()

```

Listing 1.1: Word Tokenization using NLTK

Input:

Original Text:

Natural Language Processing is an amazing field of Artificial Intelligence.
It enables machines to understand human language effectively and accurately.

Output

Sentence Tokenization:

1: Natural Language Processing is an amazing field of Artificial Intelligence.
2: It enables machines to understand human language effectively and accurately.

Word Tokenization:

['Natural', 'Language', 'Processing', 'is', 'an', 'amazing', 'field', 'of',
'Artificial', 'Intelligence', '.', 'It', 'enables', 'machines', 'to', 'understand',
'human', 'language', 'effectively', 'and', 'accurately', '.']

Viva Questions

1. What is tokenization in NLP?
2. What are the different types of tokenization?
3. How does whitespace-based tokenization differ from rule-based tokenization?
4. Why is tokenization important for text preprocessing?
5. What are some challenges in tokenization?

Stop Word Removal

Program-2

Write a Python program to perform stop word Removal tasks on text.

Stopwords

Stop word removal is a fundamental preprocessing step in Natural Language Processing (NLP). Stop words are common words (e.g., "is", "and", "the", "a", "an") that frequently occur in text but typically do not contribute to its meaning. Removing these words helps reduce the dimensionality of text data, streamline analysis, and minimize noise.

By eliminating stop words, NLP models can focus on the most relevant features, improving their performance and reducing computational overhead. This step enhances the efficiency and effectiveness of various NLP tasks. The NLTK library provides a predefined list of stop words for multiple languages, making it a convenient tool for implementing stop word removal.

Python Code

```
1  import nltk
2  from nltk.tokenize import word_tokenize
3  from nltk.corpus import stopwords
4
5  # Download NLTK resources if not already downloaded
6  nltk.download('punkt')
7  nltk.download('stopwords')
8
9  def remove_stopwords(text):
10     """
11     Removes stop words from the tokenized text.
12
13     Args:
14     text (str): The input text to be processed.
15
16     Returns:
17     list: A list of words after removing stop words.
18     """
19     # Tokenize the text into words
20     tokens = word_tokenize(text)
21
22     # Get English stop words
23     stop_words = set(stopwords.words('english'))
24
25     # Remove stop words
```



```

26     filtered_tokens = [word for word in tokens if word.lower() not in
27                         stop_words]
28
29     return filtered_tokens
30
31     def main():
32         """
33         Main function to demonstrate stop word removal.
34         """
35         # Sample text for demonstration
36         text = """This is a simple example to demonstrate how stop word removal
37                 works in text processing.
38                 It removes commonly used words like 'is', 'a', 'to', and so on."""
39
40         print("Original Text:")
41         print(text)
42
43         # Perform stop word removal
44         filtered_words = remove_stopwords(text)
45
46         # Display results
47         print("\nText After Stop Word Removal:")
48         print(filtered_words)
49
50     if __name__ == "__main__":
51         main()

```

Listing 2.1: Stop Word Removal using NLTK

Input:

This is a simple example to demonstrate how stop word removal works in text processing. It removes commonly used words like 'is', 'a', 'to', and so on.

Output:

```
['simple', 'example', 'demonstrate', 'stop', 'word', 'removal',
'works', 'text', 'processing', '.', 'removes', 'commonly',
'used', 'words', 'like', ',', 'so', '.']
```

Viva Questions

1. What are stop words, and why are they removed in NLP tasks?
2. How does case sensitivity affect stop word removal?
3. What are the advantages and disadvantages of removing stop words in text processing?
4. Can you list some applications where stop word removal is particularly useful?

Word Analysis

Program-3

Write a Python Program for Word Analysis..

Word Analysis in NLP

Word analysis in **Natural Language Processing (NLP)** involves breaking down a given text into its fundamental components to extract meaningful insights. It helps in understanding word distribution, frequency, uniqueness, and statistical properties of words in a given text corpus. Word analysis forms the foundation for more advanced NLP techniques such as *Part-of-Speech (POS) tagging*, *Named Entity Recognition (NER)*, and *Topic Modeling*.

Key Steps in Word Analysis

1. **Tokenization**: Splitting the text into individual words or tokens.
2. **Normalization**: Converting words to lowercase and removing punctuation or special characters.
3. **Frequency Analysis**: Counting the occurrences of each word to identify the most commonly used words.
4. **Unique Word Count**: Determining how many distinct words are present in the text.
5. **Statistical Analysis**: Computing measures like average word length and word density.

Program

```
1  import nltk
2  from nltk.tokenize import word_tokenize
3  from collections import Counter
4  import string
5
6  nltk.download('punkt')
7
8  def word_analysis(text):
9      words = word_tokenize(text.lower())
10     words = [word for word in words if word.isalpha()]
11
12     word_freq = Counter(words)
13
14     unique_words = len(set(words))
15     avg_word_length = sum(len(word) for word in words) / len(words) if words
16     else 0
17
18     print("\nWord Frequency Distribution:")
```

```

18     for word, freq in word_freq.most_common(10):
19         print(f"{word}: {freq}")
20
21     print("\nBasic Text Statistics:")
22     print(f"Total Words: {len(words)}")
23     print(f"Unique Words: {unique_words}")
24     print(f"Average Word Length: {avg_word_length:.2f}")
25
26     if __name__ == "__main__":
27         sample_text = "Natural Language Processing helps computers understand
28             and generate human language."
29         word_analysis(sample_text)

```

Output

Word Frequency Distribution:

```

language: 2
natural: 1
processing: 1
helps: 1
computers: 1
understand: 1
generate: 1
human: 1

```

Basic Text Statistics:

```

Total Words: 8
Unique Words: 8
Average Word Length: 6.12

```

Viva Questions

1. What is the significance of word frequency analysis in NLP?
2. How does tokenization work in NLTK?
3. Why do we remove punctuation and convert text to lowercase?
4. What is the importance of computing unique word count?
5. How can this technique be applied in real-world applications?

Word Generation

Program-4

Write a Python Program for Word Generation.
(using a bigram-based Markov model from a given text.)

Word Generation in NLP

Word generation is a fundamental concept in **Natural Language Processing (NLP)** that involves creating new text sequences by predicting the next word based on learned patterns from existing text. It is widely used in applications such as *chatbots, story generation, machine translation, and predictive text input*.

Bigram Markov Model

Here we perform word generation using **bigram-based Markov model**. A **bigram** is a sequence of two consecutive words, and a **bigram language model** calculates the probability of a word occurring given the previous word. The **Markov model**, a probabilistic model, assumes that the occurrence of a word depends only on the previous word, making it computationally efficient for text generation.

- **Training Phase:**
 - The model first processes a given text corpus and extracts bigrams (word pairs).
 - It builds a probability distribution where each word is mapped to the possible words that can follow it.
- **Generation Phase:**
 - The process starts with a given word (seed word).
 - The model probabilistically selects the next word based on the learned bigram probabilities.
 - This process continues until a desired number of words are generated or an end condition is met.

While bigram models provide a basic approach to text generation, more advanced techniques such as **trigram models, n-gram smoothing, and deep learning models** (e.g., *Recurrent Neural Networks (RNNs) and Transformers*) are used for generating more coherent and context-aware text.

Program

```

1 import random
2 import nltk
3 from nltk.tokenize import word_tokenize
4 from collections import defaultdict
5
6 # Download necessary NLTK resources
7 nltk.download('punkt')
8
9 def train_bigram_model(text):
10     """
11     Trains a simple bigram-based word generation model.
12
13     Args:
14     text (str): The input text.
15
16     Returns:
17     dict: A dictionary mapping words to their possible next words.
18     """
19     words = word_tokenize(text.lower())
20     bigram_model = defaultdict(list)
21
22     for i in range(len(words) - 1):
23         bigram_model[words[i]].append(words[i + 1])
24     print(bigram_model)
25     return bigram_model
26
27 def generate_text(bigram_model, start_word, num_words=10):
28     """
29     Generates a sequence of words based on a trained bigram model.
30
31     Args:
32     bigram_model (dict): Trained bigram word map.
33     start_word (str): The word to start generation.
34     num_words (int): Number of words to generate.
35
36     Returns:
37     str: Generated sentence.
38     """
39     current_word = start_word.lower()
40     generated_words = [current_word]
41
42     for _ in range(num_words - 1):
43         if current_word in bigram_model:
44             next_word = random.choice(bigram_model[current_word])
45             generated_words.append(next_word)
46             current_word = next_word
47         else:
48             break # Stop if no next word exists
49     return ' '.join(generated_words)
50
51 # Example Usage
52 if __name__ == "__main__":
53     sample_text = """John is a generous human being.
54     He loves to help others and always shares his knowledge.
55     Generosity and kindness are his true nature.
```

```
56 People respect him for his wisdom and humility."""
57
58 model = train_bigram_model(sample_text)
59 print("\nGenerated Text:")
60 print(generate_text(model, start_word="is", num_words=10))
```

Sample Output

```
start_word: is
Generated Text:
is a generous human being . he loves to help
```

```
start_word: help
Generated Text:
help others and humility . people respect him for his
```

Viva Questions

1. What is a Markov model, and how is it used in word generation?
2. What is the role of bigrams in text prediction?
3. How does randomness influence generated text?
4. What are the advanced models used for text generation?
5. How does a Markov model differ from neural network-based text generation?

Word Sense Disambiguation

Program-5

create a sample list of atleast five words with ambiguous sense and write a python program to implement Word Sense Disambiguation

Word Sense Disambiguation

Word Sense Disambiguation (WSD) is a fundamental problem in **Natural Language Processing (NLP)** that involves determining the correct meaning of a word based on its surrounding context. Many words in English have multiple meanings (senses), and resolving their ambiguity is crucial for applications such as:

- **Machine Translation** (e.g., Google Translate needs to correctly translate polysemous words).
- **Speech Recognition** (e.g., “light” as brightness vs. “light” as not heavy).
- **Chatbots and Virtual Assistants** (e.g., distinguishing meanings of “bank” in finance vs. river-bank).
- **Information Retrieval and Search Engines** (e.g., distinguishing “rock” as a music genre vs. a physical object).

Lesk Algorithm

In this program, we implement WSD using the **Lesk algorithm**, a dictionary-based method that selects the meaning of a word by maximizing the overlap between its dictionary definition and the surrounding context in a sentence.

The **Lesk Algorithm** is a classic approach to WSD and works as follows:

1. For a given ambiguous word in a sentence, retrieve all its possible meanings (senses) from WordNet, which is a large lexical database of English that serves as a key resource in Natural Language Processing (NLP).
2. Compute the overlap between each sense’s definition and the words in the given sentence.
3. Select the sense with the highest overlap as the correct meaning.

Though simple, the Lesk algorithm forms the basis for more advanced techniques like *supervised learning-based WSD*, *deep learning models (e.g., BERT)*, and *contextual embeddings*.

Sample List of Ambiguous Words

The following words with multiple meanings have been considered for WSD:

- **Bank:** (Financial institution / Riverbank)
- **Bat:** (Flying mammal / Equipment for playing sports)

- **Crane**: (A type of bird / A machine used for lifting heavy objects)
- **Date**: (A calendar day / A romantic meeting / A type of fruit)
- **Light**: (Not heavy / Illumination)
- **Light**: (Electromagnetic radiation / Not serious)
- **Spring**: (A season / A coiled object / A natural water source)
- **Bass**: (A type of fish / A low musical tone)
- **Rock**: (A stone / A genre of music)
- **Watch**: (A timepiece / The act of looking at something attentively)

Program

```

1  import nltk
2  from nltk.wsd import lesk
3  from nltk.tokenize import word_tokenize
4
5  # Download necessary resources
6  nltk.download('wordnet')
7  nltk.download('omw-1.4')
8  nltk.download('punkt')
9
10 # Sample list of ambiguous words and sentences
11 sample_sentences = {
12     "bank": "I deposited money in the bank.",
13     "bat": "A bat flew across the night sky.",
14     "crane": "The crane lifted heavy steel beams at the construction
15             site.",
16     "date": "Our first date was at a coffee shop.",
17     "light1": "This bag is very light to carry.", # Light (Not heavy)
18     "light2": "Please turn on the light in the room.", # Light
19             (Illumination)
20     "spring": "The flowers bloom in spring.",
21     "bass": "He caught a huge bass while fishing.",
22     "rock": "She enjoys listening to rock music.",
23     "watch": "I wear a watch on my left wrist."
24 }
25
26 def perform_wsd(word, sentence):
27     """
28     Implements Word Sense Disambiguation (WSD) using the Lesk algorithm.
29
30     Args:
31     word (str): The ambiguous word for which WSD is performed.
32     sentence (str): The sentence containing the ambiguous word.
33
34     Returns:
35     str: The best-matching sense of the word based on context.
36     """
37     tokens = word_tokenize(sentence) # Tokenize the sentence
38     actual_word = word.replace("1", "").replace("2", "") # Remove numbering
39                     from "light1", "light2"

```



```

37     sense = lesk(tokens, actual_word)  # Apply the Lesk algorithm
38
39     return sense.definition() if sense else "No suitable sense found"
40
41     # Run WSD on all sample words
42     for word, sentence in sample_sentences.items():
43         sense = perform_wsd(word, sentence)
44         print(f"Word: {word.replace('1', '').replace('2', '')}") # Remove
            numbering from print output
45         print(f"Sentence: {sentence}")
46         print(f"Predicted Sense: {sense}\n")

```

Listing 5.1: Implementation of WSD using the Lesk Algorithm

Sample Output

Word: bank
Sentence: I deposited money in the bank.
Predicted Sense: a financial institution that accepts deposits and channels the money into

Word: bat
Sentence: A bat flew across the night sky.
Predicted Sense: nocturnal mouselike mammal with forelimbs modified to form membranous wings

Word: crane
Sentence: The crane lifted heavy steel beams at the construction site.
Predicted Sense: lifts and moves heavy objects; lifting tackle is suspended from a pivot

Word: date
Sentence: Our first date was at a coffee shop.
Predicted Sense: an appointment to meet at a specified time and place

Word: light
Sentence: This bag is very light to carry.
Predicted Sense: of comparatively little physical weight or density

Word: light
Sentence: Please turn on the light in the room.
Predicted Sense: (physics) electromagnetic radiation that can produce a visual sensation

Word: spring
Sentence: The flowers bloom in spring.
Predicted Sense: the season of growth

Word: bass
Sentence: He caught a huge bass while fishing.
Predicted Sense: the lean flesh of a saltwater fish

Word: rock
Sentence: She enjoys listening to rock music.
Predicted Sense: a genre of popular music originating in the 1950s

Word: watch

Sentence: I wear a watch on my left wrist.
Predicted Sense: a small portable timepiece

Viva Questions

1. What is Word Sense Disambiguation (WSD) and why is it important in NLP?
2. How does the Lesk algorithm determine the correct sense of a word?
3. What are some real-world applications of WSD?
4. What are the limitations of the Lesk algorithm?
5. How can supervised machine learning improve WSD performance?

Install NLTK Toolkit and Perform Stemming

Program-6

Install the NLTK toolkit and perform stemming on sample text.

NLTK (Natural Language Toolkit)

The **Natural Language Toolkit (NLTK)** is a widely used Python library for Natural Language Processing (NLP). It provides tools for text processing, tokenization, stemming, lemmatization, part-of-speech tagging, named entity recognition, and other linguistic operations. NLTK includes various corpora, lexical resources, and models that help perform complex NLP tasks.

Installing NLTK

To install NLTK, the following command is executed in the terminal or command prompt:

```
pip install nltk
```

Once installed, additional datasets and resources can be downloaded using:

```
import nltk
nltk.download('all')
```

This ensures that all necessary corpora, models, and tools required for NLP tasks are available.

Key Features of NLTK

- **Tokenization**: Splitting text into words or sentences.
- **Stemming**: Reducing words to their root form.
- **Lemmatization**: Converting words to their dictionary base form.
- **POS Tagging**: Assigning grammatical tags to words.
- **Parsing & Syntax Analysis**: Understanding sentence structure.
- **Text Classification**: Classifying text documents using machine learning techniques.
- **Named Entity Recognition (NER)**: Identifying entities like names, locations, and organizations.

Stemming in NLP

Stemming is the process of reducing a word to its root or base form by removing suffixes. Unlike **lemmatization**, stemming does not necessarily produce a valid word but rather a root form that can be used for text analysis.

Various Stemming Algorithms

- **Porter Stemmer**: One of the most commonly used stemming algorithms, it applies a series of rules to iteratively remove common suffixes from English words.
- **Lancaster Stemmer** : A more aggressive stemming algorithm compared to Porter, often producing shorter root words.
- **Snowball Stemmer**: An improvement over the Porter stemmer that supports multiple languages and applies a more refined set of rules.
- **Regexp Stemmer**: Uses regular expressions to define stemming rules, allowing for customized word reductions based on specific patterns

Example of Stemming

Original Word	Stemmed Form
Running	Run
Studying	Study
Happily	Happi
Flying	Fly

Table 6.1: Examples of Stemming

Stemming is particularly useful in:

- Search engines for query optimization.
- Text normalization in NLP tasks.
- Reducing vocabulary size for machine learning models.

Python Implementation

Following Python program **install NLTK and performs stemming** using the **Lancaster Stemmer** .

```
1  # Import necessary libraries
2  import nltk
3  from nltk.stem import LancasterStemmer
4  from nltk.tokenize import word_tokenize
5
6  # Download required NLTK resources
7  nltk.download('punkt')
8
9  def perform_stemming(text):
10     """
11     Performs stemming on the input text using the Lancaster Stemmer.
12
13     Args:
14     text (str): The input text to be stemmed.
15
16     Returns:
17     list: A list of stemmed words.
18     """
19     # Initialize Lancaster Stemmer
20     stemmer = LancasterStemmer()
```

```

21
22     # Tokenize the text into words
23     words = word_tokenize(text)
24
25     # Apply stemming to each word
26     stemmed_words = [stemmer.stem(word) for word in words]
27
28     return stemmed_words
29
30 def main():
31     """
32     Main function to demonstrate stemming.
33     """
34     # Sample text for demonstration
35     text = 'Focusing is most important for reading writing
36            understanding and for any skill development'
37
38     print("Original Text:")
39     print(text)
40
41     # Perform stemming
42     stemmed_text = perform_stemming(text)
43
44     # Display results
45     print("\nStemmed Text:")
46     print(stemmed_text)
47
48     if __name__ == "__main__":
49         main()

```

Listing 6.1: Installing NLTK and perform stemming

Output

Original Text:

Focusing is most important for reading writing
understanding and for any skill development

Stemmed Text:

['focus', 'is', 'most', 'import', 'for', 'read', 'writ', 'understand',
'and', 'for', 'any', 'skil', 'develop']

Viva Questions

1. What is stemming in NLP?
2. How does the Lancaster Stemmer algorithm work?
3. What is the difference between stemming and lemmatization?
4. Why do some words not retain their original form after stemming?
5. Name four stemming algorithms.

POS Tagging

Program-7:

Create Sample list of at least 10 words POS tagging and find the POS for any given word

Introduction to POS Tagging

Part-of-Speech (POS) tagging is a fundamental task in Natural Language Processing (NLP) that involves assigning grammatical categories (such as noun, verb, adjective, etc.) to words in a sentence. POS tagging helps in understanding sentence structure and meaning, making it essential for applications such as text processing, information retrieval, and machine translation.

Each word in a sentence belongs to a specific grammatical category based on its function. The commonly used POS tags in the Penn Treebank Tagset include:

- **NN (Noun, singular)** – e.g., *dog, car, tree*
- **NNS (Noun, plural)** – e.g., *dogs, cars, trees*
- **VB (Verb, base form)** – e.g., *run, eat, sleep*
- **VBD (Verb, past tense)** – e.g., *ran, ate, slept*
- **VBG (Verb, gerund/present participle)** – e.g., *running, eating, sleeping*
- **JJ (Adjective)** – e.g., *happy, blue, large*
- **RB (Adverb)** – e.g., *quickly, silently, well*
- **IN (Preposition)** – e.g., *in, on, at, by*
- **DT (Determiner)** – e.g., *the, a, an*
- **PRP (Pronoun)** – e.g., *he, she, they*

POS tagging is typically performed using rule-based, statistical, or machine learning-based approaches. The **Natural Language Toolkit (NLTK)** provides a built-in function for POS tagging using the `pos_tag()` function, which employs a pre-trained model.

Applications of POS Tagging

- **Information Retrieval:** Helps improve search engine results by understanding query intent.
- **Machine Translation:** Assists in translating words accurately by considering their POS.
- **Text-to-Speech Systems:** Improves pronunciation by distinguishing between homonyms (e.g., "lead" as a noun vs. verb).

- **Grammar Checking:** Identifies grammatical errors in writing tools.

POS tagging is an essential task in NLP that aids in text analysis by providing grammatical information. The use of NLTK's `pos_tag()` function simplifies the process, making it easier to analyze and understand sentence structures.

Implementation in Python

The following Python program performs POS tagging on a sample list of words and allows users to find the POS of any given word.

```
1  import nltk
2  from nltk import pos_tag
3  from nltk.tokenize import word_tokenize
4
5  # Download required resources
6  nltk.download('averaged_perceptron_tagger')
7  nltk.download('punkt')
8
9  def pos_tagging(words_list):
10     """
11     Performs POS tagging on a list of words.
12     """
13     return pos_tag(words_list)
14
15  def find_pos(word, tagged_list):
16     """
17     Finds the POS tag for a given word.
18     """
19     for w, pos in tagged_list:
20         if w.lower() == word.lower():
21             return pos
22     return "POS not found"
23
24  def main():
25     """
26     Main function to demonstrate POS tagging and finding POS for a
27     given word.
28     """
29     # Sample list of words
30     words_list = ["The", "quick", "brown", "fox", "jumps", "over",
31                  "the", "lazy", "dog"]
32
33     # Perform POS tagging
34     tagged_words = pos_tagging(words_list)
35
36     # Display POS tagging results
37     print("POS tagging results:")
38     for word, pos in tagged_words:
39         print(f"{word}: {pos}")
40
41     # Take input from user
```

```

40     search_word = input("\nEnter a word to find its POS: ")
41     pos_result = find_pos(search_word, tagged_words)
42     print(f"POS for '{search_word}': {pos_result}")
43
44     if __name__ == "__main__":
45         main()

```

Listing 7.1: POS Tagging using NLTK

Output

POS tagging results:

The: DT

quick: JJ

brown: JJ

fox: NN

jumps: VBZ

over: IN

the: DT

lazy: JJ

dog: NN

Enter a word to find its POS: fox

POS for 'fox': NN

Viva Questions

1. What is POS tagging, and why is it important in NLP?
2. Which Python library is commonly used for POS tagging?
3. What does the POS tag 'NN' represent? give an example?
4. How does POS tagging help in real-world applications like search engines or chatbots?
5. What does the POS tag 'VBZ' represent? give an example?

Morphological Analysis and N-Grams

Program-8:

Write a Python program to

- a) Perform Morphological Analysis using NLTK library
- b) Generate n-grams using NLTK N-Grams library
- c) Implement N-Grams Smoothing

Morphological Analysis

Morphological analysis involves breaking down words into their base forms, typically using:

- **Stemming**: Reducing words to their root form by removing affixes (e.g., *running* → *run*).
- **Lemmatization**: Mapping words to their dictionary base forms while considering context (e.g., *better* → *good*).

In this program, we use **NLTK's WordNet Lemmatizer** to perform lemmatization and analyze a given sentence.

N-Grams and Laplace Smoothing

An *n-gram* is a contiguous sequence of *n* items (words, characters, etc.) extracted from a text. Common examples include:

- **Unigram (n=1)**: ["The", "quick", "brown", "fox"]
- **Bigram (n=2)**: [("The", "quick"), ("quick", "brown")]
- **Trigram (n=3)**: [("The", "quick", "brown"), ("quick", "brown", "fox")]

However, traditional n-gram models suffer from data sparsity—many word combinations may not appear in the training data. To handle this, **Laplace (add-one) smoothing** is applied, which assigns a small probability to unseen n-grams, preventing zero probabilities.

Implementation in Python

```
1 import nltk
2 from nltk.tokenize import word_tokenize
3 from nltk.stem import WordNetLemmatizer
4 from nltk.util import ngrams
5 from nltk.corpus import wordnet
6 from nltk import pos_tag
```

```

8      from collections import Counter
9
10     # Download required resources
11     nltk.download('punkt')
12     nltk.download('wordnet')
13     nltk.download('averaged_perceptron_tagger')
14
15     def get_wordnet_pos(word):
16         """
17         Converts POS tag to format compatible with WordNet Lemmatizer.
18         """
19         tag = pos_tag([word])[0][1][0].upper()
20         tag_dict = {"J": wordnet.ADJ, "N": wordnet.NOUN, "V":
21                     wordnet.VERB, "R": wordnet.ADV}
22         return tag_dict.get(tag, wordnet.NOUN) # Default to NOUN if
23         not found
24
25     def morphological_analysis(sentence):
26         """
27         Performs morphological analysis on a sentence using
28         lemmatization.
29         """
30         lemmatizer = WordNetLemmatizer()
31         words = word_tokenize(sentence)
32         lemmatized_words = [lemmatizer.lemmatize(word,
33             get_wordnet_pos(word)) for word in words]
34         return " ".join(lemmatized_words)
35
36     def generate_ngrams(text, n):
37         """
38         Generates n-grams from a given text.
39         """
40         tokens = word_tokenize(text)
41         return list(ngrams(tokens, n))
42
43     def laplace_smoothing(ngrams_list):
44         """
45         Implements Laplace (add-one) smoothing for n-grams.
46         """
47         ngram_counts = Counter(ngrams_list)
48         total_ngrams = sum(ngram_counts.values())
49         vocab_size = len(set(ngrams_list))
50         smoothed_probs = {ngram: (count + 1) / (total_ngrams +
51             vocab_size) for ngram, count in ngram_counts.items()}
52         return smoothed_probs, ngram_counts
53
54     def main():
55         """
56         Main function to demonstrate morphological analysis, n-gram
57         generation, and n-gram smoothing.
58         """
59         # Morphological Analysis

```

```

54     sentence = "The running dogs are barking loudly."
55     lemmatized_sentence = morphological_analysis(sentence)
56     print(f"Original sentence: {sentence}")
57     print(f"Morphological analysis (lemmatized):
58         {lemmatized_sentence}")
59
60     # N-gram Generation
61     text = "The quick brown fox jumps over the lazy dog"
62     n = 3 # Trigrams
63     trigrams_list = generate_ngrams(text, n)
64     print(f"\nGenerated {n}-grams:")
65     print(trigrams_list)
66
67     # N-gram Smoothing
68     smoothed_probs, _ = laplace_smoothing(trigrams_list)
69     print("\nN-gram probabilities after Laplace smoothing:")
70     for ngram, prob in smoothed_probs.items():
71         print(f"{ngram}: {prob:.4f}")
72
73     if __name__ == "__main__":
74         main()

```

Listing 8.1: Morphological Analysis and N-Grams

Output

Original sentence: The running dogs are barking loudly.

Morphological analysis (lemmatized): The run dog be barking loudly .

Generated 3-grams:

```

[('The', 'quick', 'brown'),
 ('quick', 'brown', 'fox'),
 ('brown', 'fox', 'jumps'),
 ('fox', 'jumps', 'over'),
 ('jumps', 'over', 'the'),
 ('over', 'the', 'lazy'),
 ('the', 'lazy', 'dog')]

```

N-gram probabilities after Laplace smoothing:

```

('The', 'quick', 'brown'): 0.1429
('quick', 'brown', 'fox'): 0.1429
('brown', 'fox', 'jumps'): 0.1429
('fox', 'jumps', 'over'): 0.1429
('jumps', 'over', 'the'): 0.1429
('over', 'the', 'lazy'): 0.1429
('the', 'lazy', 'dog'): 0.1429

```

Viva Questions

1. What is morphological analysis?
2. What is the difference between stemming and lemmatization?
3. What are n-grams, and how are they used in language modeling?
4. What is the purpose of Laplace smoothing in n-gram models?
5. How does tokenization impact n-gram generation?

Implementation of Porter Stemmer Algorithm for Stemming

Program-9

Write a Python program to implement Porter stemmer algorithm for stemming

Stemming

Stemming is a crucial preprocessing step in **Natural Language Processing (NLP)** that reduces words to their root or base form. It helps in **text normalization**, which improves the efficiency of text-based applications such as **search engines, text classification, and sentiment analysis**.

The **Porter Stemming Algorithm**, developed by **Martin Porter in 1980**, is one of the most widely used stemming techniques. It works through a series of heuristic rules to iteratively strip **suffixes** from words while maintaining the word's base meaning.

For example:

- running → run
- flies → fli
- better → better (unchanged, since it's already in its root form)

The algorithm applies the following **five steps**:

1. Removes **plural** and **-ed, -ing** suffixes.
2. Converts **y** to **i** if necessary.
3. Removes **double suffixes** (e.g., “-ization” → “-ize”).
4. Removes **derivational suffixes** (e.g., “-full” → “-ful”).
5. Handles **exceptions and irregular cases**.

Python Implementation

```
1 import nltk
2 from nltk.stem import PorterStemmer
3 from nltk.tokenize import word_tokenize
4
5 # Download required resources if not available
6 nltk.download('punkt')
7
8 def porter_stemming(text):
9     """
10    Applies Porter Stemming Algorithm to each word in the input text.
11
12    Args:
```

```

13 text (str): Input sentence or paragraph.
14
15 Returns:
16 list: List of words after applying stemming.
17 """
18 # Initialize Porter Stemmer
19 stemmer = PorterStemmer()
20
21 # Tokenize the input text
22 words = word_tokenize(text)
23
24 # Apply stemming to each token
25 stemmed_words = [stemmer.stem(word) for word in words]
26
27 return stemmed_words
28
29 def main():
30     # Example text
31     text = "Running jogging skipping   excercising workouts and dieting are good
32           for health.."
33
34     # Apply stemming
35     stemmed_output = porter_stemming(text)
36
37     # Display original and stemmed text
38     print("Original Text:")
39     print(text)
40     print("\nStemmed Output:")
41     print(" ".join(stemmed_output))
42
43     if __name__ == "__main__":
44         main()

```

Listing 9.1: Porter Stemming Algorithm in Python

Output

Original Text:

Running jogging skipping excercising workouts and dieting are good for health.

Stemmed Output:

run jog skip excercis workout and diet are good for health .

Viva Questions

1. What is stemming in NLP?
2. How does the Porter Stemming Algorithm work?
3. What is the difference between **stemming** and **lemmatization**?
4. Why is stemming useful in text processing?
5. What are the limitations of Porter Stemmer?

Conversion of Audio to Text and Text to Audio

Program-10:

Write a python program to

- Convert an audio file into text.
- Convert a text file into an audio file.
- (Using NLTK package)

Discription

Speech processing is a crucial component of NLP, allowing machines to understand spoken language and generate speech. There are two primary processes in this experiment:

Speech-to-Text (STT)

Speech recognition involves converting spoken language into text. This is achieved using libraries like:

- **SpeechRecognition**: A Python library that provides an interface to different speech recognition engines (Google Web Speech API, CMU Sphinx, etc.).
- **NLTK (Natural Language Toolkit)**: While NLTK itself does not provide direct speech recognition, it can process the converted text for further linguistic analysis.

Text-to-Speech (TTS)

Text-to-Speech (TTS) synthesis converts written text into spoken language. This is used in applications like:

- Voice assistants (e.g., Alexa, Google Assistant)
- Audiobooks and accessibility services

For this experiment, we use:

- **Google Text-to-Speech (gTTS)**: A Python library that allows converting text into speech using Google's TTS API.

Implementation

Prerequisites

Install the following libraries

- speechrecognition
- gtts
- pydub
- ffmpeg with chocolatey

```
pip install speechrecognition gtts nltk pydub
```

Python code

```

1
2 import speech_recognition as sr
3 from gtts import gTTS
4 import os
5
6 # Function to convert audio file to text
7 def audio_to_text(audio_file):
8     recognizer = sr.Recognizer()
9     with sr.AudioFile(audio_file) as source:
10         audio_data = recognizer.record(source)
11     try:
12         text = recognizer.recognize_google(audio_data)
13         print("Recognized Text:", text)
14         print("Speech-to-Text conversion successful")
15         return text
16     except sr.UnknownValueError:
17         print("Speech Recognition could not understand audio")
18         return None
19     except sr.RequestError:
20         print("Could not request results from Google Speech Recognition service")
21         return None
22
23 # Function to convert text to speech
24 def text_to_audio(text, output_file):
25     tts = gTTS(text=text, lang='en')
26     tts.save(output_file)
27     print("Text-to-Speech conversion successful. Saved as", output_file)
28
29 # Main Function
30 def main():
31     # Convert audio to text
32     audio_file = r"C:\Users\Woopa\Downloads\NLP LAB\abc.wav" # Provide the path
33     # to your audio file
34     recognized_text = audio_to_text(audio_file)
35
36     if recognized_text:
37         # Convert text to audio
38         output_audio_file = "output_audio.mp3"
39         text_to_audio(recognized_text, output_audio_file)
40
41 if __name__ == "__main__":
42     main()

```

Listing 10.1: Audio to Text and Text to Audio Conversion

-
1. The program will process the given audio file (`sample_audio.wav`) and convert it into text.
 2. The recognized text will be printed on the screen.
 3. The text will then be converted back into an audio file (`output_audio.mp3`), which can be played to hear the generated speech.

Output

```
Recognized Text: this is Natural Language Processing lab at Anurag Engineering College
Speech-to-Text conversion successful
Text-to-Speech conversion successful.
Saved as output_audio.mp3
```

Viva Questions

1. What is Speech Recognition, and how does it work?
2. What are the libraries used in this conversion?
3. How does Google Text-to-Speech (gTTS) work?
4. What are some real-world applications of text-to-speech and speech-to-text?