| Name: | Pranay Singhvi |
|---|---|
| UID: | 2021300126 |

# Experiment 7

| Aim: | Text Analysis: Exploring Chunking and Named Entity Recognition for Information Extraction from Text Corpora |
|---|---|
| Theory: | <ul><li>**Chunking**<ul><li>**Purpose:** Groups words into meaningful phrases (chunks). It's a step above part-of-speech (POS) tagging but less complex than full parsing.</li><li>**Types:**<ul><li>**Noun Phrase (NP):** A phrase built around a noun (e.g., "the big dog")</li><li>**Verb Phrase (VP):** A phrase built around a verb (e.g., "is running quickly")</li><li>**Adjective Phrase (ADJP):** A phrase built around an adjective (e.g., "very happy")</li><li>**Adverb Phrase (ADVP):** A phrase built around an adverb (e.g., "quite slowly")</li><li>**Prepositional Phrase (PP):** A phrase starting with a preposition and including its object (e.g., "on the table")</li></ul></li></ul></li><li>**Named Entity Recognition (NER)**<ul><li>**Purpose:** Locates and classifies named entities in text into predefined categories.</li><li>**Common NER Types:**<ul><li>**Persons (PER):** Names of people</li><li>**Organizations (ORG):** Companies, government bodies, etc.</li><li>**Locations (LOC):** Countries, cities, mountains, etc.</li><li>**Dates and Times:** Specific dates, time expressions</li><li>**Quantities:** Numerical values, measurements</li><li>**Miscellaneous (MISC):** Other entities (products, events)</li></ul></li></ul></li></ul>**Methods**<br><br>**1. Chunking**<br><br>**a) Regular Expressions**<ul><li>**Pros:** Flexible if you're familiar with regex.</li><li>**Cons:** Regexes can get complex; they might not be as accurate as library-based methods.</li></ul>**b) Libraries (like NLTK)**<ul><li>**NLTK (Natural Language Toolkit):** - Train a chunk parser for basic chunking. - Use rule-based grammars for more customized patterns.<ul><li>**Pros:** Generally more accurate and robust than simple regex.</li><li>**Cons:** Can require some setup and understanding of NLP libraries.</li></ul></li></ul>**2. Named Entity Recognition**<br><br>**a) Rule-Based with Regex**<ul><li>**Pros:** Simple if you need to target very specific entities.</li><li>**Cons:** Limited; less adaptable to variations in entity forms.</li></ul>**b) Libraries (like NLTK or spaCy)**<ul><li>**NLTK:** Offers a pre-trained NER model (`nltk.ne_chunk()`).</li><li>**spaCy:** Popular library with efficient and highly customizable NER models.</li><li>**Pros:** Powerful models, pre-trained on large datasets for better accuracy.</li><li>**Cons:** Might require more setup and resources.</li></ul> |

## 1. Installation of NLTK and downloading the required corpus

```python
import re
import nltk
import spacy
from prettytable import PrettyTable
```

## 2. Declare the sample text corpus

```python
# Sample text corpus
text_corpus = """
President Donald Trump visited France and met with French President Emmanuel Macron.
The movie was directed by Steven Spielberg and starred Tom Hanks.
The cat sat lazily on the mat.
"""
```

## 3. Regular expression patterns for chunking

```
grammar = r"""
    NP: {<DT|JJ|NN.*>+}          # Chunk sequences of DT, JJ, NN
    VP: {<VB.*><NP|PP|CLAUSE>+$}  # Chunk verbs and their arguments
    ADJP: {<JJ.*><PP>?}          # Chunk adjectives with optional prepositional phrase
    ADVP: {<RB.*>}                # Chunk adverbs
    PP: {<IN><NP>}                # Chunk prepositions and their objects
    """
```

## 4. Chunking using regular expressions

```python
def chunk_with_regex(text):
    sentences = nltk.sent_tokenize(text)
    for sentence in sentences:
        words = nltk.word_tokenize(sentence)
        tagged_words = nltk.pos_tag(words)
        chunk_parser = nltk.RegexpParser(grammar)
        chunks = chunk_parser.parse(tagged_words)
        print(chunks)
```

## 5. Chunking using NLTK library

```python
def chunk_with_nltk(text):
    sentences = nltk.sent_tokenize(text)
    for sentence in sentences:
        words = nltk.word_tokenize(sentence)
        tagged_words = nltk.pos_tag(words)
        chunk_parser = nltk.RegexpParser(grammar)
        chunks = chunk_parser.parse(tagged_words)
        print(chunks)
```

## 6. Named Entity Recognition using spaCy library

```python
def ner_with_spacy_pretty_table(text):
    nlp = spacy.load("en_core_web_sm")
    doc = nlp(text)
    table = PrettyTable(["Entity", "Label"])
    for ent in doc.ents:
        table.add_row([ent.text, ent.label_])
    print(table)
```

## 7. Display the results

```python
# Chunking with regular expressions
print("Chunking with regular expressions:")
chunk_with_regex(text_corpus)
```

```
Chunking with regular expressions:
(S
  (NP President/NNP Donald/NNP Trump/NNP)
  visited/VBD
  (NP France/NNP)
  and/CC
  met/VBN
  (PP with/IN (NP French/NNP President/NNP Emmanuel/NNP Macron/NNP))
  ./.)
(S
  (NP The/DT movie/NN)
  was/VBD
  directed/VBN
  (PP by/IN (NP Steven/NNP Spielberg/NNP))
  and/CC
  starred/VBD
  (NP Tom/NNP Hanks/NNP)
  ./.)
(S
  (NP The/DT cat/NN)
  sat/VBD
  (ADVP lazily/RB)
  (PP on/IN (NP the/DT mat/NN))
  ./.)
```

```python
# Chunking with NLTK library
print("\nChunking with NLTK library:")
chunk_with_nltk(text_corpus)
```

```
Chunking with NLTK library:
(S
  (NP President/NNP Donald/NNP Trump/NNP)
  visited/VBD
  (NP France/NNP)
  and/CC
  met/VBN
  (PP with/IN (NP French/NNP President/NNP Emmanuel/NNP Macron/NNP))
  ./.)
(S
  (NP The/DT movie/NN)
  was/VBD
  directed/VBN
  (PP by/IN (NP Steven/NNP Spielberg/NNP))
  and/CC
  starred/VBD
  (NP Tom/NNP Hanks/NNP)
  ./.)
(S
  (NP The/DT cat/NN)
  sat/VBD
  (ADVP lazily/RB)
  (PP on/IN (NP the/DT mat/NN))
  ./.)
```

In [ ]:
```python
# Named Entity Recognition with spaCy library
print("\nNamed Entity Recognition with spaCy:")
ner_with_spacy_pretty_table(text_corpus)
```

```
Named Entity Recognition with spaCy:
+------------------+--------+
|      Entity      | Label  |
+------------------+--------+
|   Donald Trump   | PERSON |
|      France      |  GPE   |
|      French      |  NORP  |
| Emmanuel Macron  | PERSON |
| Steven Spielberg | PERSON |
|    Tom Hanks     | PERSON |
+------------------+--------+
```

# Conclusion

Through this experiment, I discovered that combining regular expressions and libraries provides the most robust solution for extracting meaningful chunks and named entities from text. While regular expressions offer flexibility, libraries like NLTK and spaCy provide superior accuracy and efficiency for real-world natural language processing tasks.