



Bansilal Ramnath Agarwal Charitable Trust's
Vishwakarma Institute of Information Technology

**Department of
Artificial Intelligence and Data Science**

Name: Siddhesh Dilip Khairnar

Class: TY

Division: B

Roll No: 372028

Semester: 6th

Academic Year: 2023-24

Subject Name & Code: Natural Language and Processing & ADUA32203

Title of Assignment: Perform PoS tagging using regular expressions and inbuilt PoS taggers

Date of Performance: 09-02-2024

Date of Submission: 16-03-2024

ASSIGNMENT NO: - 4

Aim: Perform PoS tagging using regular expressions and inbuilt PoS taggers

❖ **THEORY:**

1) Brief Discussion on PoS Tag Sets for Different Languages:

- Part-of-speech (PoS) tagging is a fundamental task in natural language processing (NLP) that involves labelling each word in a sentence with its corresponding part of speech, such as noun, verb, adjective, etc.
- PoS tag sets vary across different languages due to linguistic differences, morphological complexities, and syntactic structures.
- Each language presents its own set of challenges and nuances for PoS tagging, and researchers continue to develop and refine PoS tag sets and tagging algorithms to improve accuracy and efficiency across different languages.

2) Brief Discussion on Various Approaches for PoS Tagging:

- Part-of-speech (PoS) tagging is a crucial task in natural language processing (NLP) that involves assigning grammatical categories (such as noun, verb, adjective, etc.) to each word in a text.
- Various approaches have been developed over the years to tackle PoS tagging, each with its own advantages and limitations.
- Here's a brief discussion on some of these approaches:

1) **Rule-Based Approach:**

- Rule-based PoS tagging relies on hand-crafted linguistic rules to assign tags to words based on patterns in their context.
- These rules may consider features such as word suffixes, prefixes, and neighbouring words.
- While rule-based approaches can be transparent and interpretable, They often require extensive linguistic expertise and may not generalize well to diverse text genres or languages.

2) Probabilistic Approach:

- Probabilistic methods, such as Hidden Markov Models (HMMs) and Conditional Random Fields (CRFs), model the likelihood of observing a sequence of tags given the input words.
- These models use training data to learn the probabilities of transitions between tags and the emission probabilities of words given tags.
- Probabilistic approaches are widely used for PoS tagging due to their flexibility, scalability, and ability to capture complex dependencies in language.

3) Hybrid Approaches:

- Hybrid approaches combine the strengths of rule-based, probabilistic, and deep learning methods to improve PoS tagging accuracy.
- For example, a hybrid system may use rule-based pre-processing to handle specific linguistic phenomena, followed by a probabilistic or deep learning model for tagging.
- These hybrid systems aim to leverage the advantages of different approaches while mitigating their individual weaknesses.

Each approach has its own set of trade-offs in terms of accuracy, efficiency, interpretability, and resource requirements. Researchers continue to explore novel techniques and combinations of methods to further advance PoS tagging performance across different languages and domains.

A) POS tagging using regular expression (English):

```
1  import nltk
2  from nltk.tokenize import word_tokenize
3
4  # Define a dictionary to map POS tags to their descriptions
5  pos_tag_descriptions = {
6      'CC': 'coordinating conjunction',
7      'CD': 'cardinal digit',
8      'DT': 'determiner',
9      'EX': 'existential there',
10     'FW': 'foreign word',
11     'IN': 'preposition/subordinating conjunction',
12     'JJ': 'adjective',
13     'JJR': 'adjective, comparative',
14     'JJS': 'adjective, superlative',
15     'LS': 'list marker',
16     'MD': 'modal',
17     'NN': 'noun, singular or mass',
18     'NNS': 'noun plural',
19     'NNP': 'proper noun, singular',
20     'NNPS': 'proper noun, plural',
21     'PDT': 'predeterminer',
22     'POS': 'possessive ending',
23     'PRP': 'personal pronoun',
24     'PRP$': 'possessive pronoun',
25     'RB': 'adverb',
26     'RBR': 'adverb, comparative',
27     'RBS': 'adverb, superlative',
28     'RP': 'particle',
29     'SYM': 'symbol',
30     'TO': 'to',
```

```

31     'UH': 'interjection',
32     'VB': 'verb, base form',
33     'VBD': 'verb, past tense',
34     'VBG': 'verb, gerund/present participle',
35     'VBN': 'verb, past participle',
36     'VBP': 'verb, sing. present, non-3d',
37     'VBZ': 'verb, 3rd person sing. present',
38     'WDT': 'wh-determiner',
39     'WP': 'wh-pronoun',
40     'WP$': 'possessive wh-pronoun',
41     'WRB': 'wh-adverb'
42 }

44 def pos_tagging(sentence):
45     # Tokenize the sentence into words
46     words = word_tokenize(sentence)
47
48     # Perform POS tagging
49     pos_tags = nltk.pos_tag(words)
50
51     # Map POS tags to descriptions
52     tagged_words = [(word, pos, pos_tag_descriptions.get(pos,
53         'Unknown')) for word, pos in pos_tags]
54
55     return tagged_words
56
57 # Take input from the user
58 sentence = input("Enter a sentence: ")
59
60 # Perform POS tagging on the user input
61 tags = pos_tagging(sentence)
62
63 # Print tagged words with descriptions
64 print("Tagged Words with Descriptions:")
65 for word, pos, description in tags:
66     print(f"Word: {word}, POS Tag: {pos}, Description:
67         {description}")

```

Output:

```
Enter a sentence: The quick brown fox jumps over the lazy dog
Tagged Words with Descriptions:
Word: The, POS Tag: DT, Description: determiner
Word: quick, POS Tag: JJ, Description: adjective
Word: brown, POS Tag: NN, Description: noun, singular or mass
Word: fox, POS Tag: NN, Description: noun, singular or mass
Word: jumps, POS Tag: VBZ, Description: verb, 3rd person sing. present
Word: over, POS Tag: IN, Description: preposition/subordinating conjunction
Word: the, POS Tag: DT, Description: determiner
Word: lazy, POS Tag: JJ, Description: adjective
Word: dog, POS Tag: NN, Description: noun, singular or mass
```

B) POS tagging using inbuilt functions (English):

```
1  import spacy
2
3  def pos_tagging_spacy(sentence):
4      # Load the English language model in spaCy
5      nlp = spacy.load("en_core_web_sm")
6
7      # Process the sentence with the model
8      doc = nlp(sentence)
9
10     # Get the POS tags and their descriptions
11     tagged_words = [(token.text, token.pos_, spacy.explain(token.
12                        pos_)) for token in doc]
13
14     return tagged_words
15
16 # Take input from the user
17 sentence = input("Enter a sentence: ")
18
19 # Perform POS tagging on the user input
20 tags_spacy = pos_tagging_spacy(sentence)
21
22 # Print tagged words with descriptions
23 for word, pos, description in tags_spacy:
24     print(f"Word: {word}, POS Tag: {pos}, Description:
25           {description}")
```

Output:

```
Enter a sentence: The quick brown fox jumps over the lazy dog
Word: The, POS Tag: DET, Description: determiner
Word: quick, POS Tag: ADJ, Description: adjective
Word: brown, POS Tag: ADJ, Description: adjective
Word: fox, POS Tag: NOUN, Description: noun
Word: jumps, POS Tag: VERB, Description: verb
Word: over, POS Tag: ADP, Description: adposition
Word: the, POS Tag: DET, Description: determiner
Word: lazy, POS Tag: ADJ, Description: adjective
Word: dog, POS Tag: NOUN, Description: noun
```

C) POS tagging using regular expression (Marathi):

```
1  import nltk
2  from nltk.tokenize import word_tokenize
3  from nltk.tag import tnt
4
5  # Marathi POS Tagset
6  marathi_pos_tagset = {
7      'PRP': 'pronoun',
8      'NN': 'noun',
9      'QF': 'quantifier',
10     'VRB': 'verb',
11     'SYM': 'symbol',
12     'CC': 'conjunction',
13     'DET': 'determiner',
14     'INTF': 'interjection',
15     'NUM': 'numeral',
16     'PSP': 'postposition',
17     'RP': 'particle',
18     'WQ': 'wh-word',
19     'JJ': 'adjective',
20     'FW': 'foreign word'
21 }
22
```



```

23 # Train POS Tagger for Marathi
24 def train_marathi_tagger():
25     marathi_data = [[('मला', 'PRP'), ('बाल', 'NN'), ('पशू', 'NN'),
26                     ('आवडतात', 'VRB')],
27                     [('तुम्हाला', 'PRP'), ('कस', 'QF'), ('वाटत', 'VRB'),
28                     ('?', 'SYM')]]
29     tagger = tnt.TnT()
30     tagger.train(marathi_data)
31     return tagger
32
33 # Tag Marathi Sentences
34 def pos_tagging_marathi(sentence, tagger):
35     words = word_tokenize(sentence)
36     tags = tagger.tag(words)
37     return tags
38
39 # Sample Marathi Sentences
40 marathi_sentences = ["मला बाल पशूआवडतात", "तुमच्या कामाला जबाबदारी द्यावी
मी नको", "तुम्हाला कसंवाटत ?"]
41
42 # Train Marathi POS Tagger
43 marathi_tagger = train_marathi_tagger()
44
45 # Tag Marathi Sentences
46 for sent in marathi_sentences:
47     print(pos_tagging_marathi(sent, marathi_tagger))

```

Output:

```

[('मला', 'PRP'), ('बाल', 'NN'), ('पशूआवडतात', 'Unk')]
[('तुमच्या', 'Unk'), ('कसाला', 'Unk'), ('जबाबदारी', 'Unk'), ('दयाली', 'Unk'), ('मी', 'Unk'), ('नको', 'Unk')]
[('तुम्हाला', 'PRP'), ('कसंवाटत', 'Unk'), ('?', 'SYM')]

```

Conclusion: Thus, we have successfully performed PoS tagging using inbuilt PoS taggers(for English and Marathi), regular expression, dictionary based PoS tagging and N-gram model based PoS tagging.