

## Experiment No.4

**Aim:** To study POS tagging.

### Theory:

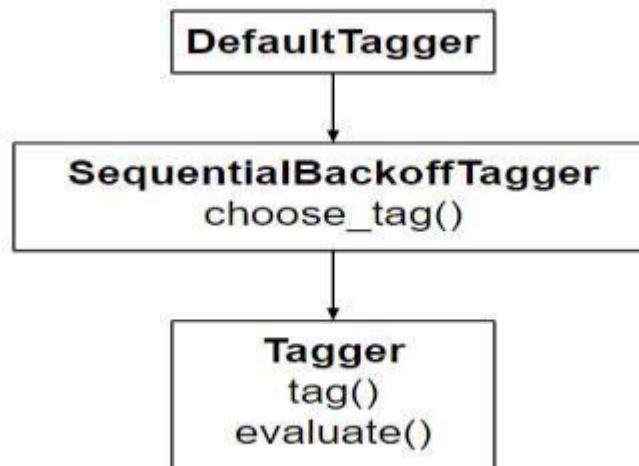
It is a process of converting a sentence to forms – list of words, list of tuples (where each tuple is having a form (word, tag)). The tag in case of is a part-of-speech tag, and signifies whether the word is a noun, adjective, verb, and so on.

Default tagging is a basic step for the part-of-speech tagging. It is performed using the DefaultTagger class. The DefaultTagger class takes 'tag' as a single argument. NN is the tag for a singular noun. DefaultTagger is most useful when it gets to work with the most common part-of-speech tag. That's why a noun tag is recommended.

Tagging is a kind of classification that may be defined as the automatic assignment of description to the tokens. Here the descriptor is called tag, which may represent one part-of-speech, semantic information and so on.

Now, if we talk about Part-of-Speech (PoS) tagging, then it may be defined as the process of assigning one of the parts of speech to the given word. It is generally called POS tagging. In simple words, we can say that POS tagging is a task of labeling each word in a sentence with its appropriate part of speech. We already know that parts of speech include nouns, verbs, adverbs, adjectives, pronouns, conjunction and their sub-categories.

Most of the POS tagging falls under Rule Base POS tagging, Stochastic POS tagging and



Transformation based tagging.

## Rule-based POS Tagging

One of the oldest techniques of tagging is rule-based POS tagging. Rule-based taggers use a dictionary or lexicon for getting possible tags for tagging each word. If the word has more than one possible tag, then rule-based taggers use hand-written rules to identify the correct tag. Disambiguation can also be performed in rule-based tagging by analyzing the linguistic features of a word along with its preceding as well as following words. For example, suppose if the preceding word of a word is article, then the word must be a noun.

## Stochastic POS Tagging

Another technique of tagging is Stochastic POS Tagging. Now, the question that arises here is which model can be stochastic. The model that includes frequency or probability (statistics) can be called stochastic. Any number of different approaches to the problem of part-of-speech tagging can be referred to as stochastic tagger.

The simplest stochastic tagger applies the following approaches for POS tagging –Word Frequency Approach

In this approach, the stochastic taggers disambiguate the words based on the probability that a word occurs with a particular tag. We can also say that the tag encountered most frequently with the word in the training set is the one assigned to an ambiguous instance of that word. The main issue with this approach is that it may yield an inadmissible sequence of tags.

## Tag Sequence Probabilities

It is another approach of stochastic tagging, where the tagger calculates the probability of a given sequence of tags occurring. It is also called the n-gram approach. It is called so because the best tag for a given word is determined by the probability at which it occurs with the n previous tags

## Transformation-based Tagging

Transformation based tagging is also called Brill tagging. It is an instance of the transformation based learning (TBL), which is a rule-based algorithm for automatic tagging of POS to the given text. TBL, allows us to have linguistic knowledge in a readable form, transforms one state to another state by using transformation rules. It draws inspiration from both the previous explained taggers – rule-based and stochastic. If we see similarity between rule-based and transformation tagger, then like rule-based, it is also based on the rules that specify what tags need to be assigned to what words. On the other hand, if we see similarity between stochastic and transformation tagger then like stochastic, it is a machine learning technique in which rules are automatically induced from data.

### HMM for POS Tagging

The POS tagging process is the process of finding the sequence of tags which is most likely to have generated a given word sequence. We can model this POS process by using a Hidden Markov Model (HMM), where tags are the hidden states that produced the observable output, i.e., the words.

#### Code:

```
import nltk  
  
nltk.download('averaged_perceptron_tagger')  
  
nltk.download('punkt')  
  
text = nltk.word_tokenize("And now for Everything completely Same")  
  
nltk.pos_tag(text)
```

#### Output:

```
[('And', 'CC'),  
( 'now', 'RB'),  
( 'for', 'IN'),  
( 'Everything', 'VBG'),  
( 'completely', 'RB'),  
( 'Same', 'JJ')]
```

#### Conclusion:

Thus, we have studied POS Tagging in the above experiment also learned regarding different types of POS Tagging and tried to implement the code for POS Tagging and successfully executed it