

# NLP Homework-4

## Part-2

**Question 2.1 :** When we lemmatize a word, we transform a word into its base form also known as lemma. For example, all the forms of the word like run, runs and running are transformed to the lemma run. Thus, when we lemmatize a sentence, we lose the form of the word. The form of the word plays an important role in determining the part of speech of the word. Therefore, if we use the lemmatized sentence instead of the actual sentence, the part of speech tagging will be incorrect as all the words are reduced to their base forms and we will be tagging the base forms of the word instead of the actual words.

## Part-3

**Question 3.1 :** Some other patterns that capture the hyponym-hypernym relation are :-

1.  $NP_H$ , like  $NP_1, NP_2, \dots$

Example Sentence : I am good at certain subjects like math and science.

In the above sentence 'subjects' is the hypernym and 'math' and 'science' are hyponyms.

2.  $NP_H$ , for example  $NP_1, NP_2, \dots$

Example Sentence : I can play quite a few musical instruments, for example, the flute, the guitar, and the piano.

In the above sentence, 'musical instruments' is the hypernym and 'flute', 'guitar' and 'piano' are hyponyms.

## Part-5

**Question 5.1 :** The values are as follows :

1. Precision : 1.000000
2. Recall : 0.148148
3. F-measure : 0.258065

**Question 5.2:** The test data only contains hypernym-hyponym pairs where the hypernym and the hyponym is mostly made up of a single word. However, the hypernyms and hyponyms that the system extracts are phrases consisting of multiple words that do not exist in the test data and are being skipped when we measure the performance of the system. This non-existence of multi-word hypernym and hyponym pairs in the test data is hurting the performance of the system. Further, when we calculate the false negatives we are assuming that all the hypernym-hyponym pairs in the gold-true set ( derived from the test data ) are present in the Wikipedia sentences dataset which in turn is effecting the recall and f-measure values. It is possible that

there are certain-hypernym-hyponym pairs that are not present in the dataset but are present in the test data. It would not be correct to count these pairs as false negatives. In order to fix this we need to use test data that contains hypernym-hyponym pairs that are similar to the ones being extracted by the system and are also present in the given dataset. In other words, the test data must contain multi-word or phrasal hypernym-hyponym pairs and all the hypernym-hyponym pairs in the gold-true set must be present in the given dataset in order to improve the performance of our system.

## Part-6

**Question 6.1 :** This assignment took me 15 hours to complete.

**Question 6.2 :** I did not discuss this homework with anyone.

## Part-7

**Question 7.1 :** The patterns that I came up with are :

1.  $NP_H$ , like  $NP_1, NP_2, \dots$

Example Sentence : I am good at certain subjects like math and science.

In the above sentence 'subjects' is the hypernym and 'math' and 'science' are hyponyms.

2.  $NP_H$ , for example  $NP_1, NP_2, \dots$

Example Sentence : I can play quite a few musical instruments, for example, the flute, the guitar, and the piano.

In the above sentence, 'musical instruments' is the hypernym and 'flute', 'guitar' and 'piano' are hyponyms.

3.  $NP_H$ , for instance  $NP_1, NP_2, \dots$

Example Sentence : There have been many leaders in history who have tried to rule the entire world, for instance, Julius Caesar and Alexander the Great.

In the above sentence 'leaders in history who have tried to rule the entire world' is the hypernym and 'Julius Caesar' and 'Alexander the Great' are hyponyms.

4.  $NP_H$ , notably  $NP_1, NP_2, \dots$

Example Sentence : He exercised a considerable influence over certain of its leaders, notably Mirabeau and Sieyes.

In the above sentence 'leaders' is the hypernym and 'Mirabeau' and 'Sieyes' are hyponyms.

**Question 7.2 :** The new measures are :

1. Precision : 0.983871
2. Recall : 0.161376
3. F-measure : 0.277273