

* Assignment :- 4 *

Page No.:

Date:

youva

Title:- Consider suitable text dataset, Remove stop words, apply stemming and feature selection techniques to represent document as vectors, classify documents & evaluate precision, recall.

* Objective:-

- * Implementation of the problem statement using python.
- * Remove stop words, apply stemming and feature selection.

* Theory:-

In Computing stop words are the words which are filtered out before and after processing of natural language data. Though stop words refer to most common words in a language, there is no single universal list of stop words used by all natural language processing tools and indeed not all tools even use such a list. Any group of word can be chosen as the stop words for a given purpose.

Code to remove stop words,

```
from nltk.tokenize import sent_tokenize, word_tokenize.
```

```
from nltk.corpus import stopwords.
```

```
data = "data mining is a subject."
```



```

stopwords = set(stopwords.words('english'))
words = word_tokenize(data)
wordsFiltered = []
for w in words:
    if w not in stopwords:
        wordsFiltered.append(w)
print(wordsFiltered)

```

* stemming:-

stemming is the process of reducing inflected (or sometimes derived) words to their word stem, base or root form - generally a written word form. The stem need not be identical to the morphological root of the word. It is usually sufficient that related words map to same stem, even if this stem is not in itself a valid root.

Code for stemming:-

```

from nltk.stem import PorterStemmer
from nltk.tokenize import sent_tokenize, word_tokenize

```

```

ps = PorterStemmer()
ex- words = ["python", "data", "mining"]

```

```

for w in ex-words:
    print(ps.stem(w))

```


* Feature Extraction:

In machine learning and statistics feature selection, also known as variable selection, attribute selection or variable subset selection is the process of selecting a subset of relevant feature for use in model construction.

- ① Feature selection techniques are used for four reasons.
- ② Simplification of models to take them easier to interpret by researchers/users.
- ③ Shorter training times.
- ④ To avoid the Curse of dimensionality.
- ⑤ Enhanced generalization by reducing over fitting.

Conclusion:

We successfully removed stop words, apply stemming and feature extraction techniques to represent documents as vectors.