

LP 2.

Assignment 4

Date of Completion:-
24.9.2020

Date of Submission:-
29.9.2020

Title:- Text analysis.

Problem Statement:- Consider a suitable text dataset. Remove stop words, apply stemming and feature selection techniques to represent documents as vectors. Classify documents and evaluate precision recall.

Learning Objectives:- To understand the process of stemming, calculating precision and recall.

Learning Outcomes:- Students will be able to understand stop words, stemming, feature selection techniques and calculate precision and recall.

Software / Hardware requirements:- ~~NLP~~ ~~NLP~~ NLP data / packages, python, Anaconda IDE.

Theory:-

Stop Words

- 1) The most commonly used words in a language.
- 2) These words are filtered out before or after the

- natural language data (txt) is processed.
- 3) They do not add much meaning to a sentence.
eg. as, the, at, a etc.

Stemming:-

- 1) A process of producing morphological variant of the root/base word.
- 2) A word is reduced to its root form.
eg loved \rightarrow love
played \rightarrow play
eating \rightarrow eat.

Precision and Recall:-

- 1) precision = $\frac{\text{Number of correct triples}}{\text{Number of triples retrieved}} = \frac{T_p}{T_p + F_p}$
 $T_p \Rightarrow$ true positive
 $F_p \Rightarrow$ false positive
- 2) recall = $\frac{\text{Number of correct triples}}{\text{Number of triples in gold set}} = \frac{T_p}{T_p + F_n}$
 $F_n \Rightarrow$ false negative
- 3) A measure of success of prediction when classes are very imbalanced.
- 4) Precision is measure of result relevancy
- 5) Recall is measure of how many truly relevant results are returned.

Algorithm:-

1. Import python different packages numpy, pandas, nltk (natural language toolkit), re (regEx)
2. read the dataset. and perform stemming on the words
3. remove the stopwords from the dataset and form a

- corpus (dataset with no stop word and stemmed words)
4. Prepare or vectorize the words. dataset. (feature extraction)
 5. Split the data in training and test set.
 6. Fit classifier model on the dataset (Naive Bayes)
 7. Predict value on the test data.
 8. Build a confusion matrix
 9. Calculate precision and recall.

Conclusion:-

Dataset Used:- Restaurant Reviews.

Conclusion:- Thus I have completed this assignment and understood how to calculate precision recall and analyze text.