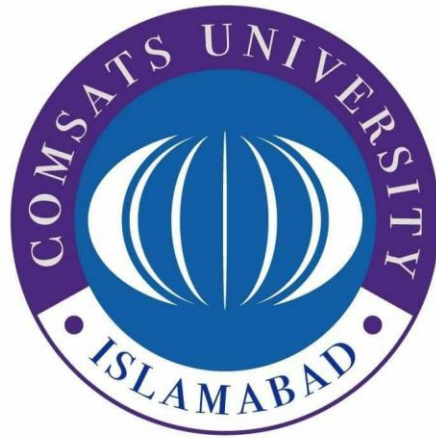


Assignment No. 4



Name: M.Sarmad Aslam

Roll No: FA21-BSE-093 - C

Subject: Data Science

COMSATS University Islamabad, Lahore Campus

11-December-2023

11-12-23

CSC461 – Assignment 4 – NLP

Muhammad Sarmad Aslam

Fa21-BSE-093

Question No. 1

i) Bag of Words (BoW):

Vocabulary: {data, science, is, one, of, the, most, important, courses, in, computer, this, best, scientists, perform, analysis}

S1: [1, 2, 1, 1, 2, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0]

S2: [1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 0, 0]

S3: [2, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1]

ii) Term Frequency (TF):

$TF = (\text{Number of times a word appears in a document}) / (\text{Total number of words in the document})$

TF(S1): [1/16, 2/16, 1/16, 1/16, 2/16, 1/16, 1/16, 1/16, 1/16, 1/16, 1/16, 0, 0, 0, 0, 0]

TF(S2): [1/16, 1/16, 1/16, 1/16, 1/16, 0, 1/16, 0, 1/16, 1/16, 0, 1/16, 1/16, 0, 0, 0]

TF(S3): [2/11, 1/11, 0, 0, 1/11, 1/11, 0, 0, 0, 0, 0, 0, 0, 0, 1/11, 1/11]

iii) Inverse Document Frequency (IDF):

$IDF = \log(\text{total number of documents} / \text{number of documents with word (term)})$

$IDF(\text{science}) = \log(3 / 3) = 0$

$IDF(\text{is}) = \log(3 / 3) = 0$

$IDF(\text{one}) = \log(3 / 2) = 0.1761$

$IDF(\text{of}) = \log(3 / 3) = 0$

$$\text{IDF}(\text{the}) = \log(3 / 3) = 0$$

$$\text{IDF}(\text{most}) = \log(3 / 1) = 1.0986$$

$$\text{IDF}(\text{important}) = \log(3 / 1) = 1.0986$$

$$\text{IDF}(\text{courses}) = \log(3 / 2) = 0.1761$$

$$\text{IDF}(\text{in}) = \log(3 / 2) = 0.1761$$

$$\text{IDF}(\text{computer}) = \log(3 / 1) = 1.0986$$

$$\text{IDF}(\text{this}) = \log(3 / 1) = 1.0986$$

$$\text{IDF}(\text{best}) = \log(3 / 1) = 1.0986$$

$$\text{IDF}(\text{scientists}) = \log(3 / 1) = 1.0986$$

$$\text{IDF}(\text{perform}) = \log(3 / 1) = 1.0986$$

$$\text{IDF}(\text{analysis}) = \log(3 / 1) = 1.0986$$

TF.IDF:

$$\text{TF.IDF} = \text{TF} * \text{IDF}$$

$$\text{TF.IDF}(S1): [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]$$

$$\text{TF.IDF}(S2): [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]$$

$$\text{TF.IDF}(S3): [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1.0986/11, 1.0986/11, 1.0986/11]$$

Question No. 2

$$\text{Cosine}(S1, S2) = \frac{\text{dot product}(\text{TF-IDF}(S1), \text{TF-IDF}(S2))}{(\text{magnitude}(\text{TF-IDF}(S1)) * \text{magnitude}(\text{TF-IDF}(S2)))}$$

$$\text{Manhattan}(S1, S2) = \text{sum}(\text{abs}(\text{TF-IDF}(S1) - \text{TF-IDF}(S2)))$$

$$\text{Euclidean}(S1, S2) = \text{sqrt}(\text{sum}((\text{TF-IDF}(S1) - \text{TF-IDF}(S2))^2))$$