

11-12-23

Assignment #4

Ayesha Zahid

FA21-BSE-003

## Question 1

Compute BOW, TF, IDF and TF.IDF values for each term in the given sentences.

S1: "data science is one of the most important courses in computer science."

S2: "this is one of the best data science courses."

S3: "the data scientist perform data analysis."

### 1- 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, 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, 1, 1, 1 ]

### 2- Term Frequency (TF):

TF = (No. of times a term appears in a sentence) / (Total no. of terms in the sentence)

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

TF(S2): [  $1/16$ ,  $1/16$ ,  $1/16$ ,  $1/16$ ,  $1/16$ , 0,  $1/16$ ,  $1/16$ , 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,  $1/11$ ,  $1/11$ ,  $1/11$  ]

### 3- Inverse Document Frequency (IDF) :

IDF =  $\log(N/n)$ , where  $N$  is the total no. of documents and  $n$  is the number of documents containing the term.

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

### 4- TF . IDF :

$$\text{TF} \cdot \text{IDF} = \text{TF} \times \text{IDF}$$

$$\text{TF} \cdot \text{IDF}(S1) : \{0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0\}$$

$$\text{TF} \cdot \text{IDF}(S2) : \{0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0\}$$

$$\text{TF} \cdot \text{IDF}(S2) : \{0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1.0986, 1.0986, 1.0986, 1.0986\}$$

## Question 2

Compute the similarities between  $S_1$ ,  $S_2$  and  $S_3$  using cosine, manhattan and euclidean distances.

Answer

$$\text{Cosine}(S_1, S_2) = \frac{\text{dot product}(\text{TF-IDF}(S_1), \text{TF-IDF}(S_2))}{(\text{magnitude}(\text{TF-IDF}(S_1)) * \text{magnitude}(\text{TF-IDF}(S_2)))}$$

$$\text{Manhattan}(S_1, S_2) = \text{sum}(\text{abs}(\text{TF-IDF}(S_1) - \text{TF-IDF}(S_2)))$$

$$\text{Euclidean}(S_1, S_2) = \sqrt{\text{sum}(\text{TF-IDF}(S_1) - \text{TF-IDF}(S_2))^2)}$$