



## COMSATS University Islamabad, Lahore Campus

Course Title:	Introduction to Data Science			Course Code:	CSC461	Credit Hours:	3(3,0)
Resource Person:	Dr. Muhammad Sharjeel			Programme Name:	BSSE		
Semester:	5 <sup>th</sup>	Batch:	FA21	Section:	C	Max Marks:	10

**HAMNA ASHRAF**  
**SP20-BSE-047**

### **Assignment 4**

**Due Date: 11-12-2023**

**Submission: Upload the assignment solution (PDF file and Python code, preferably iPython notebook) to your GitHub account (private repository).**

Important instructions: Please write the following information at the start of your ipython file.

```
# Date  
# CSC461 – Assignment4 – NLP  
# Your Full Name  
# You Complete Registration Number #  
A brief description of the task
```

*Important Instruction:*

*Solve the following questions manually as well as implement the solution using Python. Submit both.*

Q1. Compute BoW, TF, IDF, and then TF.IDF values for each term in the following three 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 scientists perform data analysis”

Q2. Compute the similarity between S1, S2, and S3 using cosine, manhattan, and euclidean distances.

Hamma Ashraf  
SP20-BSE-047  
Section - A

Assignment : 4

Q: 1 Compute Bow, TF, IDF &  
TF-IDF.

Ans:

Vocabulary (Unique Terms) :

data, science, is, one, of, the, most,  
important, courses, in, computer, this, best,  
scientists, performing, analysis

Bag of words (Bow) :

Term	S1	S2	S3
data	1	1	2
science	2	1	0
is	1	1	0
one	1	1	0
of	1	1	0
the	1	1	1
most	1	0	0
important	1	0	0
courses	1	1	0
in	1	0	0
computer	1	0	0
this	0	1	0
best	0	1	0

## Bag of words.

Term	S1	S2	S3
Scientists	0	0	1
perform	0	0	1
Analysis	0	0	1

Vector S1 : [1.211111111000000]

Vector S2 : [1111110010011000]

Vector S3 : [20000100000000111]

TF (Term Frequency):

Term	S1	S2	S3
tf(data)	$\frac{1}{12}$	$\frac{1}{9}$	$\frac{2}{6}$
tf(science)	$\frac{2}{12}$	$\frac{1}{9}$	0
tf(is)	$\frac{1}{12}$	$\frac{1}{9}$	0
tf(one)	$\frac{1}{12}$	$\frac{1}{9}$	0
tf(of)	$\frac{1}{12}$	$\frac{1}{9}$	0
tf(the)	$\frac{1}{12}$	$\frac{1}{9}$	$\frac{1}{6}$
tf(most)	$\frac{1}{12}$	0	0
tf(important)	$\frac{1}{12}$	0	0
tf(courses)	$\frac{1}{12}$	$\frac{1}{9}$	0
tf(in)	$\frac{1}{12}$	0	0
tf(computer)	$\frac{1}{12}$	0	0
tf(this)	0	$\frac{1}{9}$	0
tf(best)	0	$\frac{1}{9}$	$\frac{1}{6}$
tf(scientists)	0	0	$\frac{1}{6}$
tf(perform)	0	0	$\frac{1}{6}$
tf(analysis)	0	0	$\frac{1}{6}$

## Inverse Document Frequency (IDF)

$$\text{idf}(\text{data}) = \log\left(\frac{3}{3}\right) = 0$$

$$\text{idf}(\text{science}) = \log\left(\frac{3}{2}\right) = 0.18$$

$$\text{idf}(\text{is}) = \log\left(\frac{3}{2}\right) = 0.18$$

$$\text{idf}(\text{one}) = \log\left(\frac{3}{2}\right) = 0.18$$

$$\text{idf}(\text{of}) = \log\left(\frac{3}{2}\right) = 0.18$$

$$\text{idf}(\text{the}) = \log\left(\frac{3}{3}\right) = 0$$

$$\text{idf}(\text{most}) = \log\left(\frac{3}{1}\right) = 0.48$$

$$\text{idf}(\text{important}) = \log\left(\frac{3}{1}\right) = 0.48$$

$$\text{idf}(\text{courses}) = \log\left(\frac{3}{2}\right) = 0.18$$

$$\text{idf}(\text{in}) = \log\left(\frac{3}{1}\right) = 0.48$$

$$\text{idf}(\text{computer}) = \log\left(\frac{3}{1}\right) = 0.48$$

$$\text{idf}(\text{this}) = \log\left(\frac{3}{1}\right) = 0.48$$

$$\text{idf}(\text{best}) = \log\left(\frac{3}{1}\right) = 0.48$$

$$\text{idf}(\text{scientists}) = \log\left(\frac{3}{1}\right) = 0.48$$

$$\text{idf}(\text{perform}) = \log\left(\frac{3}{1}\right) = 0.48$$

$$\text{idf}(\text{analysis}) = \log\left(\frac{3}{1}\right) = 0.48$$

### TF-IDF :-

Term	$+f \times \text{idf}(S1)$	$+f \times \text{idf}(S2)$	$+f \times \text{idf}(S3)$
data	$\frac{1}{12} \times 0 = 0$	$\frac{1}{9} \times 0 = 0$	$\frac{2}{6} \times 0 = 0$
science	$\frac{2}{12} \times 0.18 = 0.03$	$\frac{1}{9} \times 0.18 = 0.02$	0
is	$\frac{1}{12} \times 0.18 = 0.015$	$\frac{1}{9} \times 0.18 = 0.02$	0
one	$\frac{1}{12} \times 0.18 = 0.015$	$\frac{1}{9} \times 0.18 = 0.02$	0
of	$\frac{1}{12} \times 0.18 = 0.015$	$\frac{1}{9} \times 0.18 = 0.02$	0
the	$\frac{1}{12} \times 0 = 0$	$\frac{1}{9} \times 0 = 0$	$\frac{1}{6} \times 0 = 0$
most	$\frac{1}{12} \times 0.48 = 0.04$	0	0
important	$\frac{1}{12} \times 0.48 = 0.04$	0	0
courses	$\frac{1}{12} \times 0.18 = 0.015$	$\frac{1}{9} \times 0.18 = 0.02$	0
in	$\frac{1}{12} \times 0.48 = 0.04$	0	0
computer	$\frac{1}{12} \times 0.48 = 0.04$	0	0
this	0	$\frac{1}{9} \times 0.48 = 0.053$	0
best	0	$\frac{1}{9} \times 0.48 = 0.053$	0
scientists	0	0	$\frac{1}{6} \times 0.48 = 0.08$
perform	0	0	$\frac{1}{6} \times 0.48 = 0.08$
analysis	0	0	$\frac{1}{6} \times 0.48 = 0.08$