



12/30/2022

# INTRODUCTION TO DATA SCIENCE

ASSIGNMENT # 05

**NAME:**

MUHAMMAD SUBHAN  
KHAN

**REGISTRATION NUMBER:**

SP20-BCS-098

**SECTION & GROUP:**

(B)-4

**SEMESTER:**

6<sup>th</sup>

**ASSIGNMENT:**

Introduction to Data Science

SP20-BCS-098  
CUI LAHORE

## QUESTION

Compute the Bow, TF, and IDF model for each of the terms in the following sentences then calculate the  $tf \cdot idf$  Value

S1

Sunshine state enjoy Sunshine

S2

brown fox jump high, brown fox run

S3

Sunshine state fox run fast

## Bag of Words (Bow) :-

Documents	Sunshine	state	enjoy	brown	fox	jump	high	run	fast
S1	2	1	1	0	0	0	0	0	0
S2	0	0	0	2	2	1	1	1	0
S3	1	1	0	0	1	0	0	1	1

**Term Frequency (tf) :-** 
$$= \frac{\text{No. of times word appear in document}}{\text{Total Words}}$$

**S1 :-**

$$\text{Sunshine} = \frac{2}{4} \Rightarrow \frac{1}{2}$$

$$\text{State} = \frac{1}{4}$$

$$\text{enjoy} = \frac{1}{4}$$

**S2 :-**

$$\text{brown} = \frac{2}{7}$$

$$\text{fox} = \frac{2}{7}$$

$$\text{jump} = \frac{1}{7}$$

$$\text{high} = \frac{1}{7}$$

$$\text{Run} = \frac{1}{7}$$



**S<sub>3</sub> :-**

$$\begin{aligned}\text{Sunshine} &= \frac{1}{5} \\ \text{state} &= \frac{1}{5} \\ \text{fox} &= \frac{1}{5} \\ \text{run} &= \frac{1}{5} \\ \text{fast} &= \frac{1}{5}\end{aligned}$$

**Invers Document Frequency :-**  $\log_{10} \left( \frac{\text{Total Documents}}{\text{Word appear in Documents}} \right)$

$$\begin{aligned}\text{Sunshine} &= \log_{10} \left( \frac{3}{2} \right) = 0.1762 \\ \text{state} &= \log_{10} \left( \frac{3}{2} \right) = 0.1762 \\ \text{enjoy} &= \log_{10} \left( \frac{3}{1} \right) = 0.477 \\ \text{brown} &= \log_{10} \left( \frac{3}{1} \right) = 0.477 \\ \text{fox} &= \log_{10} \left( \frac{3}{2} \right) = 0.1762 \\ \text{jump} &= \log_{10} \left( \frac{3}{1} \right) = 0.477 \\ \text{high} &= \log_{10} \left( \frac{3}{1} \right) = 0.477 \\ \text{run} &= \log_{10} \left( \frac{3}{2} \right) = 0.1762 \\ \text{fast} &= \log_{10} \left( \frac{3}{1} \right) = 0.477\end{aligned}$$

**Term Frequency \* Inverse Document Frequency**

**S<sub>1</sub> :-**

$$\text{Sunshine} = \frac{2}{4} * 0.1762 = 0.0881$$

$$\text{state} = \frac{1}{4} * 0.176 = 0.044$$

$$\text{enjoy} = \frac{1}{4} * 0.477 = 0.11925$$

**S<sub>2</sub> :-**



$$\text{brown} = \frac{2}{7} * 0.477 = 0.136$$

$$\text{fox} = \frac{2}{7} * 0.176 = 0.051$$

$$\text{jump} = \frac{1}{7} * 0.477 = 0.068$$

$$\text{high} = \frac{1}{7} * 0.477 = 0.068$$

$$\text{run} = \frac{1}{7} * 0.176 = 0.025$$

**S<sub>3</sub> :-**

$$\text{sunshine} = \frac{1}{5} * 0.176 = 0.0352$$

$$\text{state} = \frac{1}{5} * 0.176 = 0.0352$$

$$\text{fox} = \frac{1}{5} * 0.176 = 0.0352$$

$$\text{run} = \frac{1}{5} * 0.176 = 0.0352$$

$$\text{fast} = \frac{1}{5} * 0.477 = 0.0954$$

### TF \* IDF Table

Vocabulary	TF*IDF(S1)	TF*IDF(S2)	TF*IDF(S3)
sunshine	0.088	0	0.0352
state	0.044	0	0.0352
enjoy	0.11925	0	0
brown	0	0.136	0
fox	0	0.051	0
jump	0	0.068	0.0352



high	0	0.068	0
own	0	0.025	0.0352
fast	0	0	0.0954

**QUESTION-NO-02**  
**Cosine Similarity between S1 and S3 :-**  

$$\text{Cosine Similarity} = \frac{\vec{S_1} \cdot \vec{S_3}}{|\vec{S_1}| |\vec{S_3}|}$$

Document Vectors :

$$S_1 = [2, 1, 1, 0, 0, 0, 0, 0, 0, 0]$$

$$S_3 = [1, 1, 0, 0, 1, 0, 0, 1, 1]$$

$$\begin{aligned}\vec{S_1} \cdot \vec{S_3} &= (2 \times 1) + (1 \times 1) + (1 \times 0) + (0 \times 0) + (0 \times 1) + (0 \times 0) + (0 \times 0) + \\ &\quad (0 \times 1) + (0 \times 1) \\ &= 2 + 1\end{aligned}$$

$$\vec{S_1} \cdot \vec{S_3} = 3$$

$$\begin{aligned}|\vec{S_1}| &= \sqrt{(2 \times 2) + (1 \times 1) + (1 \times 1)} \\ &= \sqrt{4 + 1 + 1}\end{aligned}$$

$$|\vec{S_1}| = 2.45$$

$$\begin{aligned}|\vec{S_3}| &= \sqrt{1 \times 1 + 1 \times 1 + 1 \times 1 + 1 \times 1 + 1 \times 1} \\ &= \sqrt{5} \\ &= 2.24\end{aligned}$$

putting values in formula :-

$$\begin{aligned}\cos \theta &= \frac{\vec{S_1} \cdot \vec{S_3}}{|\vec{S_1}| |\vec{S_3}|} \\ &= \frac{3}{2.45 \times 2.24}\end{aligned}$$

$$\cos(S_1, S_3) = 0.547$$