

NAME: ALEENA KHAN

ROLL NO: SP20-BCS-112

Group #4:

- Q: S1 : sunshine state enjoy sunshine
 S2 : brown fox jump high, brown fox run
 S3 : sunshine state fox run fast .

	BoW1	BoW2	BoW3	TF-1	TF-2	TF-3	IDF
sunshine	2	0	1	0.5	0	0.2	0.176
state	1	0	1	0.25	0	0.2	0.176
enjoy	1	0	0	0.25	0	0	0.4471
brown	0	2	0	0	0.285	0	0.4471
fox	0	2	1	0	0.285	0.2	0.176
jump	0	1	0	0	0.142	0	0.4771
high	0	1	0	0	0.142	0	0.4771
run	0	1	1	0	0.142	0.2	0.176
fast.	0	0	1	0	0	0.2	0.4771
length:	4	7	5				

Formulas:

• BoW = count in single document

• TF = $\frac{\text{count in single doc}}{\text{length of document}}$

تاریخ: _____

$$\text{IDF} = \log \left(\frac{\text{no of docu}}{\text{no of docu in which it occurred}} \right)$$

$$\bullet \text{ sunshine} = \log \left(\frac{3}{2} \right) = 0.176$$

$$\text{state} = \log (3/2) = 0.176$$

$$\text{enjoy} = \log (3/1) = 0.4771$$

$$\text{brown} = \log (3/1) = 0.4771$$

$$\text{fox} = \log (3/2) = 0.176$$

$$\text{jump} = \log (3/1) = 0.4771$$

$$\text{high} = \log (3/1) = 0.4771$$

$$\text{run} = \log (3/2) = 0.176$$

$$\text{fast} = \log (3/1) = 0.4771$$

For $S1 \cdot S3$:

$$\begin{aligned} &= (2 \times 1) + (1 \times 1) + (1 \times 0) + (0 \times 0) + (0 \times 1) + \\ &\quad (0 \times 0) + (0 \times 0) + (0 \times 1) + (0 \times 1) \\ &= 2 + 1 = 3 \end{aligned}$$

$$\begin{aligned} |S1| &= \sqrt{2^2 + 1^2 + 1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2} \\ &= \sqrt{4+1+1} = \sqrt{6} = 2.45 \end{aligned}$$

$$\begin{aligned} |S3| &= \sqrt{1^2 + 1^2 + 0^2 + 0^2 + 1^2 + 0^2 + 0^2 + 1^2 + 1^2} \\ &= \sqrt{1+1+1+1+1} = \sqrt{5} = 2.24 \end{aligned}$$

$$\begin{aligned} \text{So } \cos(S1, S3) &= \frac{S1 \cdot S3}{|S1||S3|} \\ &= \frac{3}{(2.45)(2.24)} = 0.55 \end{aligned}$$

When I was trying to compute the TF-IDF matrix by hand using the standard formula, after computing TF and IDF first and then multiplying the two I did realize, there was something different compared to the one I obtained with Scikit-learn on my sample corpus, there I realized the difference between the Scikit-learn version and most standard and traditional version.

Formulas:

Scikit-Learn

- $\text{IDF}(t) = \log \frac{1+n}{1+\text{df}(t)} + 1$

Standard notation

- $\text{IDF}(t) = \log \frac{n}{\text{df}(t)}$