

IDS-6II

Assignment no 5 (NLP)

FA19-BCS-131

Question no 1

Compute the Bow, TF, IDF and TFIDF.

S_1 "Sunshine state enjoy sunshine"

S_2 "Brown fox jump high, brown fox run"

S_3 "Sunshine State fox fun fast".

For Bow document term matrix

	Sunshine	state	enjoy	brown	fox	jump	high	run	fast
S_1	2	1	1	0	0	0	0	0	0
S_2	0	0	0	2	2	1	1	1	0
S_3	1	1	0	0	1	0	0	1	1

Vector S_1 : [2, 1, 1, 0, 0, 0, 0, 0, 0, 0]

Vector S_2 : [0, 0, 0, 2, 2, 1, 1, 1, 0]

Vector S_3 : [1, 1, 0, 0, 1, 0, 0, 1, 1]

Term - Frequency.

	Sunshine	state	enjoy	brown	fox	jump	high	run	fast
tf- S_1	2/4	1/4	1/4	0	0	0	0	0	0
tf- S_2	0	0	0	2/7	2/7	1/7	1/7	1/7	0
tf- S_3	1/5	1/5	0	0	1/5	0	0	1/5	1/5

: Inverse Document frequency (Idf)

$$\text{Idf}(\text{Sunshine}) = \log(3/2) = 0.18$$

$$\text{Idf}(\text{State}) = \log(3/2) = 0.18$$

$$\text{Idf}(\text{enjoy}) = \log(3/1) = 0.48$$

$$\text{Idf}(\text{brown}) = \log(3/1) = 0.48$$

$$\text{Idf}(\text{fox}) = \log(3/2) = 0.18$$

$$\text{Idf}(\text{jump}) = \log(3/1) = 0.48$$

$$\text{Idf}(\text{high}) = \log(3/1) = 0.48$$

$$\text{Idf}(\text{run}) = \log(3/2) = 0.18$$

$$\text{Idf}(\text{fast}) = \log(3/1) = 0.48$$

TF · IDF

Term Frequency · Inverse Document Frequency

S₁ →

$$\text{tf} \cdot \text{Idf}(\text{Sunshine}) = 2 \times 0.18 = 0.09$$

$$\text{tf} \cdot \text{Idf}(\text{state}) = 1/4 \times 0.18 = 0.045$$

$$\text{tf} \cdot \text{Idf}(\text{enjoy}) = 1/4 \times 0.48 = 0.12$$

S₂ →

$$\text{tf} \cdot \text{Idf}(\text{brown}) = 2/7 \times 0.48 = 0.14$$

$$\text{tf} \cdot \text{Idf}(\text{fox}) = 2/7 \times 0.18 = 0.05$$

$$\text{tf} \cdot \text{Idf}(\text{jump}) = 1/7 \times 0.48 = 0.07$$

$$\text{tf} \cdot \text{Idf}(\text{high}) = 1/7 \times 0.48 = 0.07$$

$$\text{tf} \cdot \text{Idf}(\text{run}) = 1/7 \times 0.18 = 0.03$$

$S_3 \rightarrow$

$$\text{tf.idf}(\text{Sunshine}) = 1/5 \times 0.18 = 0.04$$

$$\text{tf.idf}(\text{state}) = 1/5 \times 0.18 = 0.04$$

$$\text{tf.idf}(\text{fox}) = 1/5 \times 0.18 = 0.04$$

$$\text{tf.idf}(\text{run}) = 1/5 \times 0.18 = 0.04$$

$$\text{tf.idf}(\text{fast}) = 1/5 \times 0.48 = 0.096$$

Q2

Cosine Similarity.

S_1 & S_3 .

TF-IDF vectors.

$$S_1 = [0.09, 0.045, 0.12, 0, 0, 0, 0, 0, 0]$$

$$S_3 = [0.04, 0.04, 0, 0, 0.04, 0.04, 0, 0, 0.096]$$

$$\cos(S_1, S_3) = \frac{S_1 \cdot S_3}{|S_1| |S_3|}$$

$$S_1 \cdot S_3 = (0.09 \times 0.04) + (0.045 \times 0.04) + 0 + 0 + 0 + 0 + 0 + 0 + 0$$
$$+ 0 = 0.0054$$

$$|S_1| = \sqrt{(0.09 \times 0.09) + (0.045 \times 0.045) + 0 + 0 + (0.12)^2 + 0 + 0 + 0 + 0}$$
$$= 0.15660$$

$$|S_3| = \sqrt{0.0016 + 0.0016 + 0.0016 + 0.0016 + 0.0036 + 0.0036 + 0.008832}$$

$$|S_3| = \sqrt{(0.04)(0.04) + (0.04 \times 0.04) + (0.04 \times 0.04) + (0.04 \times 0.04) + (0.096 \times 0.096)}$$
$$= 0.125$$

$$\cos(s_1, s_3) = \frac{0.0054}{0.15660 \times 0.125}$$

$$\cos(s_1, s_3) = 0.2759$$