# Introduction To Data Science

NAME: MUHAMMAD HUZAIFA JAWAD
Roll No: SP20-BCS-144
Group: 4 (IV)
Submitted To: SIR Muhammad Sharjeel

## ASSIGNMENT # 05

### QUESTION 01:

S1: "sunshine state enjoy sunshine"
S2: "brown fox jump high, brown fox run"
S3: "sunshine state fox run fast"

### BoW Model

| | Sunshine | state | enjoy | brown | fox | jump | high | run | fast | Total length |
|---|---|---|---|---|---|---|---|---|---|---|
| S1 | 2 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 4 |
| S2 | 0 | 0 | 0 | 2 | 2 | 1 | 1 | 1 | 0 | 7 |
| S3 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 5 |

### TF Model

| | sunshine | state | enjoy | brown | fox | jump | high | run | fast | Total length |
|---|---|---|---|---|---|---|---|---|---|---|
| S1 | 2/4 | 1/4 | 1/4 | 0 | 0 | 0 | 0 | 0 | 0 | 4 |
| S2 | 0 | 0 | 0 | 2/7 | 2/7 | 1/7 | 1/7 | 1/7 | 0 | 7 |
| S3 | 1/5 | 1/5 | 0 | 0 | 1/5 | 0 | 0 | 1/5 | 1/5 | 5 |

# IDF Model

formula :

$$IDF('word') = \log\left(\frac{\text{Total No. of Documents}}{\text{No. of Documents containing 'word'}}\right)$$

$$IDF(sunshine) = \log\left(\frac{3}{2}\right) = 0.1761$$

$$IDF(state) = \log\left(\frac{3}{2}\right) = 0.1761$$

$$IDF(enjoy) = \log\left(\frac{3}{1}\right) = 0.4771$$

$$IDF(brown) = \log\left(\frac{3}{1}\right) = 0.4771$$

$$IDF(fox) = \log\left(\frac{3}{2}\right) = 0.1761$$

$$IDF(jump) = \log\left(\frac{3}{1}\right) = 0.4771$$

$$IDF(high) = \log\left(\frac{3}{1}\right) = 0.4771$$

$$IDF(run) = \log\left(\frac{3}{2}\right) = 0.1761$$

$$IDF(fast) = \log\left(\frac{3}{1}\right) = 0.4771$$

|     | sunshine | state | enjoy | brown | fox | jump | high | run | fast |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| IDF | 0.1761 | 0.1761 | 0.4771 | 0.4771 | 0.1761 | 0.4771 | 0.4771 | 0.1761 | 0.4771 |

## TF.IDF values

for S1:

$tf.idf(sunshine) = 2/4 * 0.1761 = 0.0880$

$tf.idf(state) = 1/4 * 0.1761 = 0.0440$

$tf.idf(enjoy) = 1/4 * 0.4771 = 0.1192$

for S2:

$tf.idf(brown) = (2/7)(0.4771) = 0.1363$

$tf.idf(fox) = (2/7)(0.1761) = 0.0503$

$tf.idf(jump) = (1/7)(0.4771) = 0.0681$

$tf.idf(high) = (1/7)(0.4771) = 0.0681$

$tf.idf(run) = (1/7)(0.1761) = 0.0251$

for S3:

$tf.idf(sunshine) = (1/5)(0.1761) = 0.0352$

$tf.idf(state) = (1/5)(0.1761) = 0.0352$

$tf.idf(fox) = (1/5)(0.1761) = 0.0352$

$tf.idf(run) = (1/5)(0.1761) = 0.0352$

$tf.idf(fast) = (1/5)(0.4771) = 0.0954$

| | S1 | S2 | S3 |
|---|---|---|---|
| sunshine | 0.0880 | 0 | 0.0352 |
| state | 0.0440 | 0 | 0.0352 |
| enjoy | 0.1192 | 0 | 0 |
| brown | 0 | 0.1363 | 0 |
| fox | 0 | 0.0503 | 0.0352 |
| jump | 0 | 0.0681 | 0 |
| high | 0 | 0.0681 | 0 |
| run | 0 | 0.0251 | 0.0352 |
| fast | 0 | 0 | 0.0954 |

QUESTION # 02:

cosine similarity b/w S1 & S3.

using Bow model to generate vectors

$S1 = < 2 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 >$
$S3 = < 1 \quad 1 \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 1 \quad 1 >$

formula:

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

$S_1 \cdot S_3 = (2)(1) + (1)(1) + (1)(0) + (0)(0) + (0)(1) + (0)(0) + (0)(0) + (0)(1)$
$\qquad\qquad + (0)(1)$

$\qquad = 2+1$
$\qquad = 3$

$|S_1| = \sqrt{2^2 + 1^2 + 1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 - 0^2}$
$\qquad = \sqrt{4 + 1 + 1} \qquad = \sqrt{6}$
$|S_1| = 2.4494$

$|S_3| = \sqrt{1^2 + 1^2 + 0^2 + 0^2 + 1^2 + 0^2 + 0^2 + 1^2 + 1^2}$
$\qquad = \sqrt{5}$
$|S_3| = 2.2360$

$$\cos(S_1, S_3) = \frac{3}{(2.4494)(2.2360)}$$

$\boxed{\cos(S_1, S_3) = 0.5477}$

$S_1, S_3 = \cos^{-1}(0.5477)$
$\boxed{S_1, S_3 = 56.78}$