

Introduction to Data Science**Assignment#5****Sentences:**

- sunshine state enjoy sunshine
- brown fox jump high, brown fox run
- sunshine state fox run fast

Bag Of Words:

	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

Term Frequencies:

	Sunshine	State	Enjoy	Brown	Fox	Jump	High	Run	fast
S1	1/2	1/4	1/4	0	0	0	0	0	0
S2	0	0	0	2/7	2/7	1/7	1/7	1/7	0
S3	1/5	1/5	0	0	1/5	0	0	1/5	1/5

Inverse Document Frequencies:

$$\text{Idf('sunshine')} = \log(3/2) = 0.1760$$

$$\text{Idf('state')} = \log(3/2) = 0.1760$$

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

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

$$\text{Idf('fox')} = \log(3/2) = 0.1760$$

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

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

$$\text{Idf('run')} = \log(3/2) = 0.1760$$

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

	Sunshine	State	Enjoy	Brown	Fox	Jump	High	Run	fast
IDF	0.1760	0.1760	0.4771	0.4771	0.1760	0.4771	0.4771	0.1760	0.4771

TF-IDF:

S1:

$\text{Idf}(\text{'sunshine'}) = 0.1760 \times 1/2 = 0.088$

$\text{Idf}(\text{'state'}) = 0.1760 \times 1/4 = 0.044$

$\text{Idf}(\text{'enjoy'}) = 0.4771 \times 1/4 = 0.1192$

$\text{Idf}(\text{'brown'}) = 0.4771 \times 0 = 0$

$\text{Idf}(\text{'fox'}) = 0.1760 \times 0 = 0$

$\text{Idf}(\text{'jump'}) = 0.4771 \times 0 = 0$

$\text{Idf}(\text{'high'}) = 0.4771 \times 0 = 0$

$\text{Idf}(\text{'run'}) = 0.1760 \times 0 = 0$

$\text{Idf}(\text{'fast'}) = 0.4771 \times 0 = 0$

S2:

$\text{Idf}(\text{'sunshine'}) = 0.1760 \times 0 = 0$

$\text{Idf}(\text{'state'}) = 0.1760 \times 0 = 0$

$\text{Idf}(\text{'enjoy'}) = 0.4771 \times 0 = 0$

$\text{Idf}(\text{'brown'}) = 0.4771 \times 2/7 = 0.1363$

$\text{Idf}(\text{'fox'}) = 0.1760 \times 2/7 = 0.0502$

$\text{Idf}(\text{'jump'}) = 0.4771 \times 1/7 = 0.0681$

$\text{Idf}(\text{'high'}) = 0.4771 \times 1/7 = 0.0681$

$\text{Idf}(\text{'run'}) = 0.1760 \times 1/7 = 0.0251$

$\text{Idf}(\text{'fast'}) = 0.4771 \times 0 = 0$

S3:

$\text{Idf}(\text{'sunshine'}) = 0.1760 \times 1/5 = 0.0352$

$\text{Idf}(\text{'state'}) = 0.1760 \times 1/5 = 0.0352$

$\text{Idf}(\text{'enjoy'}) = 0.4771 \times 0 = 0$

$\text{Idf}(\text{'brown'}) = 0.4771 \times 0 = 0$

$\text{Idf}(\text{'fox'}) = 0.1760 \times 1/5 = 0.0352$

$\text{Idf}(\text{'jump'}) = 0.4771 \times 0 = 0$

$\text{Idf}(\text{'high'}) = 0.4771 \times 0 = 0$

$\text{Idf}(\text{'run'}) = 0.1760 \times 1/5 = 0.0352$

$\text{Idf}(\text{'fast'}) = 0.4771 \times 1/5 = 0.0954$

	Sunshine	State	Enjoy	Brown	Fox	Jump	High	Run	fast
$\text{tfidf}(S1)$	0.088	0.044	0.11925	0	0	0	0	0	0
$\text{tfidf}(S2)$	0	0	0	0.1363	0.0502	0.0681	0.0681	0.0251	0
$\text{tfidf}(S3)$	0.0352	0.0352	0	0	0.0352	0	0	0.0352	0.0954

Cosine Similarity Between S1 and S3:

$\text{Cos}(S1, S3) = S1 \cdot S3 / |S1| |S3|$

Taking Bag of Words Vector:

$S1 = [2, 1, 1, 0, 0, 0, 0, 0, 0, 0]$

$S3 = [1, 1, 0, 0, 1, 0, 0, 0, 1, 1]$

$S1 \cdot S3 = 2 \cdot 1 + 1 \cdot 1 + 1 \cdot 0 + 0 \cdot 0 + 0 \cdot 1 + 0 \cdot 0 + 0 \cdot 0 + 0 \cdot 0 + 0 \cdot 1 + 0 \cdot 1 = 3$

$|S1| = (2^2 + 1^2 + 1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2)^{0.5} = 2.4494$

$|S2| = (1^2 + 1^2 + 0^2 + 0^2 + 1^2 + 0^2 + 0^2 + 0^2 + 1^2 + 1^2)^{0.5} = 2.2360$

$\text{Cos}(S1, S3) = 3 / (2.4494 \cdot 2.2360)$

$\text{Cos}(S1, S3) = 3 / 5.4768$

$\text{Cos}(S1, S3) = 0.5477$