

① G1 Bag of Words.

	data	science	is	one	of	the	most	important	courses	in	computer	this	best	scientists	perform	analysis
S ₁	1	2	1	1	1	1	1	1	1	1	1	0	0	0	0	0
S ₂	1	1	1	1	1	1	0	0	1	0	0	1	1	0	0	0
S ₃	2	0	0	0	0	1	0	0	0	0	0	0	0	1	1	1

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

$$S_2 = [1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0]$$

$$S_3 = [2 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 1]$$

Term frequency.

	data	science	is	one	of	the	most	important	courses	in	computer	this	best	scientists	perform	analysis
S ₁	$\frac{1}{12}$	$\frac{2}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{12}$	0	0	0	0	0
S ₂	$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$	0	0	$\frac{1}{9}$	0	0	$\frac{1}{9}$	$\frac{1}{9}$	0	0	0
S ₃	$\frac{2}{6}$	0	0	0	0	$\frac{1}{6}$	0	0	0	0	0	0	0	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

Idf & tf ID F

IDF	idf	tf-IDF(S ₁)	tf-IDF(S ₂)	tf-idf(S ₃)
data	$\log(3/3) = 0$	0	$1/9 \times 0 = 0$	$2/6 \times 0 = 0$
Science	$\log(3/2) = 0.17609$	$2 \times 1/2 \times 0.17609 = 0.02935$	$1/9 \times 0.17609 = 0.0196$	0
is	$\log(3/2) = 0.17609$	$1/12 \times 0.17609 = 0.015$	$1/9 \times 0.17609 = 0.0196$	0
one	$\log(3/2) = 0.17609$	$1/12 \times 0.17609 = 0.015$	$1/9 \times 0.17609 = 0.0196$	0
of	$\log(3/2) = 0.17609$	$1/12 \times 0.17609 = 0.015$	$1/9 \times 0.17609 = 0.0196$	0
the	$\log(3/3) = 0$	$1/12 \times 0 = 0$	$1/9 \times 0 = 0$	$1/6 \times 0 = 0$
most	$\log(3/1) = 0.4771$	$1/12 \times 0.4771 = 0.03976$	$0 \times 0.4771 = 0$	0
important	$\log(3/1) = 0.4771$	$1/12 \times 0.4771 = 0.03976$	$0 \times 0.4771 = 0$	0
courses	$\log(3/2) = 0.17609$	$1/12 \times 0.17609 = 0.015$	$1/9 \times 0.17609 = 0.0196$	0
in	$\log(3/1) = 0.4771$	$1/12 \times 0.4771 = 0.03976$	$0 \times 0.4771 = 0$	0
computer	$\log(3/1) = 0.4771$	$1/12 \times 0.4771 = 0.03976$	$1/9 \times 0.4771 = 0.0530$	0
this	$\log(3/1) = 0.4771$	0	$1/9 \times 0.4771 = 0.053$	0
best	$\log(3/1) = 0.4771$	0	$0 \times 0.4771 = 0$	$1/6 \times 0.4771 = 0.07952$
scientist	$\log(3/1) = 0.4771$	0	$0 \times 0.4771 = 0$	$1/6 \times 0.4771 = 0.07952$
perform	$\log(3/1) = 0.4771$	0	$0 \times 0.4771 = 0$	$1/6 \times 0.4771 = 0.07952$
analysis	$\log(3/1) = 0.4771$	0		

Q2) Cosine:

Bag of words:

$$\cos \theta = \frac{\bar{s}_1 \cdot \bar{s}_2 \cdot \bar{s}_3}{|\bar{s}_1| |\bar{s}_2| |\bar{s}_3|}$$

$$\cos(s_1, s_2) = \frac{s_1 \cdot s_2}{|s_1| |s_2|}$$

$$s_1 \cdot s_2 = 1 \cdot 1 + 2 \cdot 1 + 1 \cdot 1 + 1 \cdot 1 + 1 \cdot 1 + 1 \cdot 1 + 1 \cdot 0 + 1 \cdot 0 + 1 \cdot 1 + 1 \cdot 0 + 1 \cdot 0 \\ + 1 \cdot 1 + 1 \cdot 0 + 1 \cdot 0 + 0 \cdot 1 + 0 \cdot 1 + 0 \cdot 0 + 0 \cdot 0 + 0 \cdot 0$$

$$s_1 \cdot s_2 = 9$$

$$|s_1| = (1 + 4 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 0 + 0 + 0 + 0)^{10.5} \\ = 14^{10.5} = 3.7427$$

$$|s_2| = 1 + 1 + 1 + 1 + 1 + 1 + 0 + 0 + 1 + 0 + 0 + 1 + 1 + 0 + 0 + 0 \\ = \sqrt{9} = 3$$

$$|s_1| |s_2| = 11.2251$$

$$\cos(s_1, s_2) = 0.8017$$

$$\cos(s_2, s_3) = \frac{s_2 \cdot s_3}{|s_2| |s_3|}$$

$$s_2 \cdot s_3 = 2 + 1 + 0 = 3$$

$$|s_2| = \sqrt{9} = 3$$

$$|s_3| = \sqrt{6} = 2.4495$$

$$\cos(s_2, s_3) = 0.4082$$

$$\bullet \cos(s_1, s_3) = \frac{s_1 \cdot s_3}{|s_1| |s_3|}$$

$$\times s_1 \cdot s_2 = 1 + 2 + 1 + 1 + 1 + 1 + 1 + 0 = 8$$

$$|s_1| = s_1 \cdot s_1 = 2 + 1 + 0 = 3$$

$$|s_1| = 3.7417$$

$$|s_3| = \sqrt{6}$$

$$= 2.4495$$

$$\cos(s_1, s_2) = 0.327$$

Euclidean (Bag of words)
distance.

$$e-d(s_1, s_2) = \sqrt{\sum_{i=1}^n (s_{1i} - s_{2i})^2}$$

$$= (0^2 + 1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1^2 + 1^2 + 0^2 + 1^2 + 1^2 + (-1)^2 + (-1)^2 + 0^2 + 0^2 + 0^2)^{1/2}$$

$$= \sqrt{7} = 2.646$$

$$e-d(s_1, s_3) = \sqrt{\sum_{i=1}^n (s_{1i} - s_{3i})^2}$$

$$= [(-1)^2 + (2)^2 + (1)^2 + (1)^2 + (1)^2 + (0)^2 + 1^2 + 1^2 + 1^2 + 1^2 + 1^2 + 0^2 + 0^2 + (-1)^2 + (-1)^2 + (-1)^2]^{1/2}$$

$$= \sqrt{16} = 4$$

$$e-d(s_2, s_3) = \sqrt{\sum_{i=1}^n (s_{2i} - s_{3i})^2}$$

$$[(-1)^2 + 1^2 + 1^2 + 1^2 + 1^2 + 0^2 + 0^2 + 0^2 + 1^2 + 1^2 + 1^2 + (-1)^2 + (-1)^2 + (-1)^2]^{1/2}$$

$$= \sqrt{11} = 3.32$$

Manhattan distances (Bag of words)

$$m-d(s_1, s_2) = \sum_{i=1}^n |s_{1i} - s_{2i}|$$

$$= [1 + 1 + 0 + 0 + 0 + 0 + 1 + 1 + 0 + 1 + 1 + 1 + 1 + 0 + 0 + 0] = 8$$

$$m-d(s_1, s_3) = \sum_{i=1}^n |s_{1i} - s_{3i}|$$

$$[1 + 2 + 1 + 1 + 1 + 0 + 1 + 1 + 1 + 1 + 1 + 0 + 0 + 1 + 1 + 1] = 14$$

$$m-d(s_2, s_3) = \sum_{i=1}^n |s_{2i} - s_{3i}|$$

$$[1 + 1 + 1 + 1 + 1 + 0 + 0 + 0 + 1 + 0 + 0 + 1 + 1 + 1 + 1 + 1] = 11$$