

Date: 11-12-2023  
Maryam Yousaf

Assignment 4  
FA21-GSE-077

## Question 1

S1: "data science is one of the most important courses in computer science".

S2: "This is one of the best data science courses"

S3: "The data scientists perform data analysis"

Compute Bag of words

vocabulary: 16 unique words

	data	science	is	one	of	the	most	important
S1	1	2	1	1	1	1	1	1
S2	1	1	1	1	1	1	0	0
S3	2	0	0	0	0	1	0	0

	courses	in	computer	This	best	scientists
S1	1	1	1	0	0	0
S2	1	0	0	1	1	0
S3	0	0	0	0	0	1

	perform	analysis
S1	0	0
S2	0	0
S3	1	1

vector  $S_1$  [121111111100000]  
 vector  $S_2$  [111111001001100]  
 vector  $S_3$  [200001000000011]

## Term frequency:

for  $S_1$  length: 12

$\text{data} = \frac{1}{12} = 0.0833$   
 $\text{science} = \frac{2}{12} = 0.1667$   
 $\text{is} = \frac{1}{12} = 0.0833$   
 $\text{one} = \frac{1}{12} = 0.0833$   
 $\text{of} = \frac{1}{12} = 0.0833$   
 $\text{the} = \frac{1}{12} = 0.0833$   
 $\text{most} = \frac{1}{12} = 0.0833$   
 $\text{important} = \frac{1}{12} = 0.0833$   
 $\text{courses} = \frac{1}{12} = 0.0833$   
 $\text{in} = \frac{1}{12} = 0.0833$   
 $\text{computer} = \frac{1}{12} = 0.0833$

for  $S_2$  length: 9

$\text{this: } \frac{1}{9} = 0.1111$   
 $\text{is: } \frac{1}{9} = 0.1111$   
 $\text{one: } \frac{1}{9} = 0.1111$   
 $\text{of: } \frac{1}{9} = 0.1111$   
 $\text{the: } \frac{1}{9} = 0.1111$   
 $\text{best: } \frac{1}{9} = 0.1111$   
 $\text{data: } \frac{1}{9} = 0.1111$   
 $\text{science: } \frac{1}{9} = 0.1111$

$\text{this} = 0$   
 $\text{best} = 0$   
 $\text{scientists} = 0$   
 $\text{perform} = 0$   
 $\text{analysis} = 0$

$\text{courses: } \frac{1}{9} = 0.1111$   
 $\text{scientists} = 0$   
 $\text{perform} = 0$   
 $\text{analysis} = 0$   
 $\text{important} = 0$   
 $\text{computer} = 0$   
 $\text{most} = 0$   
 $\text{in} = 0$

idf  
Total documents = 3

one of the

$$S_3 \text{ length: } 6 \\ \text{tf} = \frac{1}{6} = 0.167 \\ \text{tf}e = \frac{1}{6} = 0.167 \\ \text{data} = \frac{2}{6} = 0.333 \\ \text{scientists} = \frac{1}{6} = 0.167$$

$$\text{uniform} = \frac{1}{6} = 0.167 \\ \text{analysis} = \frac{1}{6} = 0.167$$

	data	science	is	one	of	the	most	imp	courses	in	computer	this
S <sub>1</sub>	0.833	0.1667	0.083	0.083	0.08	0.08	0.083	0.083	0.083	0.083	0.083	0
S <sub>2</sub>	0.11	0.111	0.11	0.11	0.11	0.11	0	0	0	0.11	0	0.11
S <sub>3</sub>	0.33	0	0	0	0	0.167	0	0	0	0	0	0

best scientist uniform analysis

	best	scientist	uniform	analysis
S <sub>1</sub>	0	0	0	0
S <sub>2</sub>	0.11	0	0	0
S <sub>3</sub>	0	0.1667	0.1667	0.1667

idf

Total documents = 3

	data	science	is	one	of	the
S1	$\log(\frac{3}{2}) 0$	$\log(\frac{3}{2}) 0$	$\log(\frac{3}{2}) 0.176$	0.176	0.176	$\log(\frac{3}{2}) 0$
S2	0	0	0.176	0.176	0.176	0
S3	0	0	0.176	0.176	0.176	0

	most	important	courses	in	computer
S1	$\log(\frac{3}{2}) 0.477$	0.477	0.176	0.477	0.477
S2	0.477	0.477	0.176	0.477	0.477
S3	0.477	0.477	0.176	0.477	0.477

	this	best	scientists	perform	analysis
S1	0.477	0.477	0.477	0.477	0.477
S2	0.477	0.477	0.477	0.477	0.477
S3	0.477	0.477	0.477	0.477	0.477

If idf for S1

$$\text{data} : 0.0833 \times 0 = 0$$

$$\text{science} : 0.1667 \times 0 = 0$$

$$\text{is} : 0.0833 \times 0.176 = 0.0147$$

$$\text{one} : 0.0833 \times 0.176 = 0.0147$$

$$\text{of} : 0.0833 \times 0.176 = 0.0147$$

$$\text{the} : 0.0833 \times 0 = 0$$

$$\text{most} : 0.0833 \times 0.477 = 0.0391$$

$$\text{important} : 0.0833 \times 0.477 = 0.0391$$

$$\text{courses} : 0.0833 \times 0.176 = 0$$

$$\text{in} : 0.0833 \times 0.477 = 0$$

$$\text{computer} : 0.0833 \times 0.477 = 0$$

$$\text{this} : 0$$

$$\text{best} : 0$$

$$\text{scientists} : 0$$

$$\text{perform} : 0$$

$$\text{analysis} : 0$$

## Half for S<sub>2</sub>

This:  $0.11 \times 0.477 = 0.05247$   
is:  $0.11 \times 0.176 = 0.01936$   
One:  $0.11 \times 0.176 = 0.01936$   
of:  $0.11 \times 0.176 = 0.01936$   
the:  $0.11 \times 0 = 0$   
best:  $0.11 \times 0.477 = 0.05247$   
data:  $0.11 \times 0 = 0$   
Science:  $0.11 \times 0 = 0$   
analysis: 0

course:  $0.11 \times 0.176 = 0.01936$   
most: 0  
important: 0  
in: 0  
computer: 0  
scientists: 0  
perform: 0

## Half for S<sub>3</sub>

The:  $0.167 \times 0 = 0$   
data: 0  
Scientists:  $0.167 \times 0.477 = 0.079659$   
perform:  $0.167 \times 0.477 = 0.079659$   
data: 0  
analysis:  $0.167 \times 0.477 = 0.079659$

## High Table

Data Science is one of the most important courses in computer science

sh	0	0.047	0.047	0	0.0397	0.0397	0.0397
s <sub>1</sub>	0	0	0.019	0.019	0	0	0.019
s <sub>2</sub>	0	0	0.019	0.019	0	0	0
s <sub>3</sub>	0	0	0	0	0	0	0

this best student perform analysis

s <sub>1</sub>	0	0	0	0	0	0	0
s <sub>2</sub>	0.05247	0.05247	0	0	0	0	0
s <sub>3</sub>	0	0	0.0797	0.0797	0.0797	0.0797	0.0797

## Question 2

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

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

Compute Similarity b/w  $S_1, S_2, S_3$

using cosine

$$\cos(S_1, S_2) = \frac{S_1 \cdot S_2}{\sqrt{S_1^2} \sqrt{S_2^2}}$$

$$\begin{aligned} S_1 \cdot S_2 &= 1 \cdot 1 + 2 \cdot 1 + 1 \cdot 0 + 1 \cdot 0 + 1 \cdot 1 + 0 + 1 \cdot 0 + \\ &\quad 0 \cdot 1 + 0 \cdot 1 + 0 \cdot 0 + 0 \cdot 0 + 0 \cdot 0 \\ &= 1 + 2 + 1 + 1 + 1 + 1 + 1 \\ &= 8 \end{aligned}$$

$$\begin{aligned} |S_1| &= (1^*1 + 2^*2 + 1^*1 + 1^*1 + 1^*1 + 1^*1 + 1^*1 + 1^*1 + 1^*1 + 1^*1 + 1^*1 + \\ &\quad 1^*1 + 0^*0 + 0^*0 + 0^*0 + 0^*0 + 0^*0)^{0.5} \\ &= (1+4+1+1+1+1+1+1+1+1+1+1)^{0.5} \\ &= (14)^{0.5} = 7 \end{aligned}$$

$$\begin{aligned} |S_2| &= (1^*1 + 1^*1 + 1^*1 + 1^*1 + 1^*1 + 1^*1 + 0^*0 + 0^*0 + 1^*1 + 0^*0 \\ &\quad + 0^*0 + 1^*1 + 1^*1 + 0^*0 + 0^*0 + 0^*0)^{0.5} \\ &= (1+1+1+1+1+1+1+1+1+1)^{0.5} \\ &= (9)^{0.5} = 4.5 \end{aligned}$$

$$\cos(S_1, S_2) = \frac{8}{4.5 \times 7} = \frac{8}{31.5} = 0.25 \text{ Ans}$$

$S_1 [12111111100000]$

$S_3 [2000010000000111]$

$$\cos(S_1, S_3) = \frac{S_1 \cdot S_3}{\sqrt{S_1^2} \sqrt{S_3^2}}$$

$$S_1 \cdot S_3 = 2 + 1 + 0$$

$$= 3$$

$$\begin{aligned}\sqrt{S_1^2} &= 1+4+1+1+1+1+1+1+1+1 \\ &= 14^{0.5} \\ &= 7\end{aligned}$$

$$\begin{aligned}\sqrt{S_3^2} &= 4+1+1+1+1 \\ &= 8^{0.5} \\ &= 4\end{aligned}$$

$$\cos(S_1, S_3) = \frac{3}{7 \times 4} = \frac{3}{28} = 0.107 \text{ Ans}$$

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

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

$$\cos(S_2, S_3) = \frac{S_2 \cdot S_3}{\sqrt{S_2^2} \sqrt{S_3^2}}$$

$$S_2 \cdot S_3 = 2 + 1 = 3$$

$$\begin{aligned}\sqrt{S_2^2} &= (1+1+1+1+1+1+1+1+1)^{0.5} \\ &= 9^{0.5} \\ &= 3\end{aligned}$$

$$\begin{aligned}\sqrt{s_3^2} &= (4+1+1+1+1)^{0.5} \\ &= 8^{0.5} \\ &= 4\end{aligned}$$

$$\cos(S_2, S_3) = \frac{3}{4 \times 4.5} = \frac{3}{18} = 0.167$$

## Using Manhattan

$$\begin{aligned}
 S_1, S_2 &= |1-1| + |2-1| + |1-1| + |1-1| + |1-1| + |1-1| + |1-0| + |1-0| + |1-1| + |1-0| + \\
 &\quad |1-0| + |0-1| + |0-1| + |0-0| + |0-0| + |0-0| \\
 &= 1 + 1 + 1 + 1 + 1 + 1 \\
 &= 7
 \end{aligned}$$

$$\begin{aligned}
 S_2, S_3 &= |1-2| + |1-0| + |1-0| + |1-0| + |1-0| + |1-1| + |0-0| + |0-0| + |1-0| \\
 &\quad + |0-0| + |0-0| + |1-0| + |1-0| + |0-1| + |0-1| + |0-1| \\
 &= 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 \\
 &= 11
 \end{aligned}$$

Using Euclidean Distances

formula  $P_3 - P_2 = \sqrt{\sum(P_3 - P_2)^2}$

$$\begin{aligned}
 S_1, S_2 &= [(1-1)^2 + (2-1)^2 + (1-1)^2 + (1-1)^2 + (1-1)^2 + (1-1)^2 + (1-0)^2 + (1-0)^2 + (1-1)^2 + \\
 &\quad (1-0)^2 + (1-0)^2 + (0-1)^2 + (0-1)^2 + (0-0)^2 + (0-0)^2 + (0-0)^2]^{0.5} \\
 &= [(1)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2 + (-1)^2 + (-1)^2]^{0.5} \\
 &= 7^{0.5} \\
 &= 3.5
 \end{aligned}$$

$$S_2, S_3 = \left[ (-2)^2 + (-0)^2 + (-0)^2 + (-0)^2 + (-1)^2 + (0-0)^2 + (0-0)^2 + (-0)^2 + (0-0)^2 + (0-0)^2 + (-0)^2 + (-1)^2 + (0-1)^2 + (0-1)^2 \right] \times 0.5$$

$$= (-1)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2 + (-1)^2 + (-1)^2 + (-1)^2$$

$$= 11^{\wedge} 0.5$$

$$= 5.5$$

$$S_1, S_3 = (-2)^2 + (2-0)^2 + (-0)^2 + (-0)^2 + (-0)^2 + (-1)^2 + (1-0)^2 + (1-0)^2 + (0-0)^2 + (0-0)^2 + (0-1)^2 + (0-1)^2 + (0-1)^2 + (0-1)^2$$

$$= (-1)^2 + (2)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2 + (-1)^2 + (-1)^2 + (-1)^2$$

$$= 12 + 4$$

$$= 16^{\wedge} 0.5$$

$$= 8$$