

National University of Computer and Emerging Sciences, Lahore Campus

Course: Information Retrieval
Program: BS(Data Science)
Duration: 25 Minutes
Paper Date: 22 Feb 2023
Section: BDS-6A
Exam: Quiz 2

Course Code: CS4051
Semester: Spring 2023
Total Marks: 10
Weight 0
Page(s): 2
Roll No:

Question 1

a) What is size of vocabulary (number of unique words) if total words in a collection are 100,000 Use Heap's law with $k = 10$ and $\beta = 0.5$ [2 Marks]

$$V = k(N)^\beta$$

$$V = 10(100,000)^{0.5}$$

$$V = 3160$$

(b) What is size of vocabulary if 500 new documents are added to the above collection? Assume each document has 1000 words on average. [2 Marks]

$$N = 100,000 + 500 * 1000 = 600,000$$

$$V = k(N)^\beta$$

$$V = 10(600,000)^{0.5}$$

$$V = 7740$$

Question 2

Represent following 4 documents as vectors. Use TF.IDF weights. [3 Marks]

Document 1: Omar drove the car.

Document 2: Nadia and Omar ate oranges and snacks.

Document 3: Nadia ate oranges not snacks.

Document 4: Nadia and Omar ate snacks.

	Omar	drove	the	car	Nadia	and	ate	oranges	snacks	not
<D1>	1	1	1	1	0	0	0	0	0	0
<D2>	1	0	0	0	1	1	1	1	1	0
<D3>	0	0	0	0	1	0	0	1	1	1
<D4>	1	0	0	0	1	1	1	0	1	0

(a) Calculate cosine similarity between vectors of document 2 and document 3. [3 Marks]

$$\text{Cosine sim (D2,D3)} = \mathbf{D2.D3} / (|\mathbf{D2}|*|\mathbf{D3}|)$$

$$\mathbf{D2.D3} = 1*0 + 0*0 + 0*0 + 0*0 + 1*1 + 1*0 + 1*0 + 1*1 + 1*1 + 0*1 = 1+1+1 = 3$$

$$|\mathbf{D2}| = \sqrt{(1)^2 + (0)^2 + (0)^2 + (0)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2 + (0)^2} = \sqrt{6} = 2.45$$

$$|\mathbf{D3}| = \sqrt{(0)^2 + (0)^2 + (0)^2 + (0)^2 + (1)^2 + (0)^2 + (0)^2 + (1)^2 + (1)^2 + (1)^2} = \sqrt{4} = 2$$

$$\text{Cosine sim (D2,D3)} = \mathbf{D2.D3} / (|\mathbf{D2}|*|\mathbf{D3}|) = 3 / (2.45*2) = 0.61$$