

#Date: 11-12-23

#CSC 461 - Assignment 4 - NLP

#Nosheen Azhar

#FA20-BSE-061

#Description: In the below task we will calculate TF, WOF, IDF and TF-IDF in question 1 and then the similarity between S1, S2 and S3 using cosine, manhattan and euclidean distances.

Question 1.

S1: The "Data is one of the most important courses in computers science"

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

S3: The data scientists perform data analysis

BOW = 'data', 'science', 'is', 'one', 'of', 'the', 'most', 'important', 'courses', 'in', 'computer', 'this', 'best', 'scientists', 'perform', 'analysis'

Vector R1: [1, 2, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0]

Vector R2: [1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0]

Vector R3: [2, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1]

Length of Vector R1 = 12

Length of Vector R2 = 9

Length of Vector R3 = 6

TF of Vector S_1

$$\begin{aligned}
 tf('data') &= \frac{1}{12} \\
 tf('science') &= \frac{2}{12} = \frac{1}{6} \\
 tf('is') &= \frac{1}{12} \\
 tf('one') &= \frac{1}{12} \leftarrow tf('of') = \frac{1}{12} \\
 tf('the') &= \frac{1}{12} \\
 tf('most') &= \frac{1}{12} \\
 tf('important') &= \frac{1}{12} \\
 tf('courses') &= \frac{1}{12} \\
 tf('in') &= \frac{1}{12} \\
 tf('computer') &= \frac{1}{12} \\
 tf('this') &= \frac{0}{12} = 0 \\
 tf('best') &= \frac{0}{12} = 0 \\
 tf('scientist') &= \frac{0}{12} = 0 \\
 tf('perform') &= \frac{0}{12} = 0 \\
 tf('analysis') &= \frac{0}{12} = 0
 \end{aligned}$$

TF of Vector S_2

$$\begin{aligned}
 tf('data') &= \frac{1}{9} \\
 tf('science') &= \frac{1}{9} \\
 tf('is') &= \frac{1}{9} \\
 tf('one') &= \frac{1}{9} \leftarrow tf('of') = \frac{1}{9} \\
 tf('the') &= \frac{1}{9} \\
 tf('most') &= \frac{0}{9} = 0 \\
 tf('important') &= \frac{0}{9} = 0 \\
 tf('courses') &= \frac{1}{9} \\
 tf('in') &= \frac{0}{9} = 0 \\
 tf('computer') &= \frac{0}{9} = 0 \\
 tf('this') &= \frac{0}{9} = 0 \\
 tf('best') &= \frac{1}{9} \\
 tf('scientist') &= \frac{0}{9} = 0 \\
 tf('perform') &= \frac{0}{9} = 0 \\
 tf('analysis') &= \frac{0}{9} = 0
 \end{aligned}$$

TF of Vector S_3

$$\begin{aligned}
 tf('data') &= \frac{2}{6} = \frac{1}{3} \\
 tf('science') &= \frac{0}{6} = 0 \\
 tf('is') &= \frac{0}{6} = 0 \\
 tf('one') &= \frac{0}{6} = 0 \\
 tf('of') &= \frac{0}{6} = 0 \\
 tf('the') &= \frac{1}{6} \\
 tf('most') &= \frac{0}{6} = 0 \\
 tf('important') &= \frac{0}{6} = 0 \\
 tf('courses') &= \frac{0}{6} = 0 \\
 tf('in') &= \frac{0}{6} = 0 \\
 tf('computer') &= \frac{0}{6} = 0 \\
 tf('this') &= \frac{0}{6} = 0 \\
 tf('best') &= \frac{0}{6} = 0 \\
 tf('scientist') &= \frac{1}{6} \\
 tf('perform') &= \frac{1}{6} \\
 tf('analysis') &= \frac{1}{6}
 \end{aligned}$$

idf()

$$\text{idf}(\text{'data'}) = \log(3/3) = 0$$

$$\text{idf}(\text{'Science'}) = \log(3/2) = 0.18$$

$$\text{idf}(\text{'is'}) = \log(3/2) = 0.18$$

$$\text{idf}(\text{'one'}) = \log(3/2) = 0.18$$

$$\text{idf}(\text{'of'}) = \log(3/2) = 0.18$$

$$\text{idf}(\text{'The'}) = \log(3/3) = 0$$

$$\text{idf}(\text{'most'}) = \log(3/1) = 0.48$$

$$\text{idf}(\text{'important'}) = \log(3/1) = 0.48$$

$$\text{idf}(\text{'courses'}) = \log(3/2) = 0.18$$

$$\text{idf}(\text{'in'}) = \log(3/1) = 0.48$$

$$\text{idf}(\text{'computes'}) = \log(3/1) = 0.48$$

$$\text{idf}(\text{'this'}) = \log(3/1) = 0.48$$

$$\text{idf}(\text{'best'}) = \log(3/1) = 0.48$$

$$\text{idf}(\text{'scientist'}) = \log(3/1) = 0.48$$

$$\text{idf}(\text{'perform'}) = \log(3/1) = 0.48$$

$$\text{idf}(\text{'analysis'}) = \log(3/1) = 0.48$$

data	$tf \times idf(R_1)$	$tf \times idf(R_2)$	$tf \times idf(R_3)$
Science	0	0	0
is	0.03	0.02	0
one	0.015	0.02	0
of	0.015	0.02	0
the	0	0	0
most	0.04	0	0
important	0.04	0	0
courses	0.015	0.02	0
in	0.04	0	0
computer	0.04	0	0
this	0	0.053	0
best	0	0.053	0
2) scientist	0	0	0.08
perform	0	0	0.08
analysis	0	0	0.08

Question 2: Cosine Similarity

Through BotW

$$\cos \theta = \frac{\bar{S}_1 \cdot \bar{S}_2 \cdot \bar{S}_3}{\|S_1\| \|S_2\| \|S_3\|}$$

$$\cos(S_1, S_2) = \frac{S_1 \cdot S_2}{\|S_1\| \|S_2\|} = \frac{9}{(3.74)(3)} = 0.3017$$

$$\cos(S_2, S_3) = \frac{S_2 \cdot S_3}{\|S_2\| \|S_3\|} = \frac{3}{(3)(2.4495)} = 0.4082$$

$$\cos(S_1, S_3) = \frac{S_1 \cdot S_3}{\|S_1\| \|S_3\|} = \frac{8}{(3.74)(2.4495)} = 0.327$$

Through Rfid

$$\begin{aligned} \cos(S_1, S_2) &= \frac{S_1 \cdot S_2}{\|S_1\| \|S_2\|} = \frac{(0 \times 0) + (0.03 \times 0.02) + \dots + (0 \times 0)}{(0.009)(0.087)} = \frac{1.8 \times 10^{-3}}{7.85 \times 10^{-3}} \\ &= \frac{36}{157} = 0.229 \end{aligned}$$

$$\cos(S_2, S_3) = \frac{S_2 \cdot S_3}{\|S_2\| \|S_3\|} = \frac{0}{(0.087)(0.138)} = 0$$

$$\cos(S_1, S_3) = \frac{S_1 \cdot S_3}{\|S_1\| \|S_3\|} = \frac{0}{(0.09)(0.138)} = 0$$

Manhattan Distance through BOW.

$$\begin{aligned}\text{Manhattan Distance}(S_1, S_2) &= \sum_i |S_{1i} - S_{2i}| \\ &= |0-1| + (1-1) + (2-1) + \dots + (0-0)| \\ &= 7\end{aligned}$$

$$\begin{aligned}\text{Manhattan Distance}(S_2, S_3) &= \sum_i |S_{2i} - S_{3i}| \\ &= |(1-2) + (1-0) + \dots + (0-1)| \\ &= 14\end{aligned}$$

$$\begin{aligned}\text{Manhattan Distance}(S_1, S_3) &= \sum_i |S_{1i} - S_{3i}| \\ &= |(1-2) + (2-0) + (\dots + (0-1))| \\ &= 11\end{aligned}$$

Euclidean Distance Using BOW.

$$\begin{aligned}\sqrt{(S_1, S_2)^2} &= \sqrt{(1-1)^2 + (2-1)^2 + \dots + (0-0)^2} \\ &= 2.6458\end{aligned}$$

$$\begin{aligned}\sqrt{(S_2, S_3)^2} &= \sqrt{(1-2)^2 + (1-0)^2 + \dots + (0-1)^2} \\ &= 14\end{aligned}$$

$$\begin{aligned}\sqrt{(S_1, S_3)^2} &= \sqrt{(1-2)^2 + (2-0)^2 + \dots + (0-1)^2} \\ &= 11.\end{aligned}$$

Manhattan Distance through fidf

$$\begin{aligned}\text{Manhattan Distance}(S_1, S_2) &= \sum_i |S_{1i} - S_{2i}| \\ &= |(0-0) + (0.03-0.02) + \dots + (0-0)| \\ &= 0.054\end{aligned}$$

$$\begin{aligned}\text{Manhattan Distance}(S_2, S_3) &= \sum_i |S_{2i} - S_{3i}| \\ &= |(0-0) + (0.02-0) + \dots + (0-0.08)| \\ &= |-0.034| = 0.034\end{aligned}$$

$$\begin{aligned}\text{Manhattan Distance}(S_1, S_3) &= \sum_i |S_{1i} - S_{3i}| \\ &= |(0-0) + (0.03-0) + \dots + (0-0.08)| \\ &= |0.01| = 0.01\end{aligned}$$

Euclidean Distance through fidf

$$\begin{aligned}\text{Euclidean } D(S_1, S_2) &= \sqrt{\sum_i (S_{1i} - S_{2i})^2} = \sqrt{(0-0)^2 + (0.03-0.02)^2 + \dots + (0-0)^2} \\ &= 0.1105\end{aligned}$$

$$\begin{aligned}\text{Euclidean } D(S_2, S_3) &= \sqrt{\sum_i (S_{2i} - S_{3i})^2} = \sqrt{(0-0)^2 + (0.02-0)^2 + \dots + (0-0.08)^2} \\ &= 0.1637\end{aligned}$$

$$\begin{aligned}\text{Euclidean } D(S_1, S_3) &= \sqrt{\sum_i (S_{1i} - S_{3i})^2} \\ &= \sqrt{(0-0)^2 + (0.03-0)^2 + \dots + (0-0.08)^2} \\ &= 0.1655\end{aligned}$$