



BIRZEIT UNIVERSITY

Electrical and Computer Engineering Department

ENCSENC5342: "Information Retrieval, Web Search and NLP"

Assignment #1: Instructor: Dr. Adnan H. Yahya,

Leyan Burait 1211439

Due date: May 25, 2025, strongly encouraged to submit before the exam (5/5/2025).

Problem1. Create the term frequency matrix for the four document titles below:

Doc_1 = "breakthrough drug for Schizophrenia"
Doc_2 = "new schizophrenia drug"
Doc_3 = "new approach for treatment of schizophrenia"
Doc_4 = "new hopes for Schizophrenia patients"

Retain every word (no stopword removal). The indexed units should be words converted to lower case, no stemming. The entry in each cell should be the number of occurrences of the term in the associated document.

(problem 1)

tf Matrix

ads	D ₁	D ₂	D ₃	D ₄
approach	0	0	1	0
breakthrough	1	0	0	0
drug	1	1	0	0
for	1	0	1	1
hopes	0	0	0	1
new	0	1	1	1
of	0	0	1	0
Patients	0	0	0	1
Schizophrenia	1	1	1	1
treatment	0	0	1	0

Problem2. Create a new incidence matrix from the matrix you prepared in problem 1 (1s and 0's only).

ads	D ₁	D ₂	D ₃	D ₄
approach	0	0	1	0
breakthrough	1	0	0	0
drug	1	1	0	0
for	1	0	1	1
hopes	0	0	0	1
new	0	1	1	1
of	0	0	1	0
Patients	0	0	0	1
Schizophrenia	1	1	1	1
treatment	0	0	1	0

Problem3. Extract the rows for the terms "schizophrenia" and "patients" from the matrix and write them down, one above the other. Each should have four binary digits (0 or 1) in them. Then draw a horizontal line under them and compute the AND of the two rows. This should produce a single row in which a 1 appears for each document that has both "schizophrenia" AND "patients". What are these documents? Show your work.

Problem 3:

① row for schizophrenia & patients

	D ₁	D ₂	D ₃	D ₄
schizophrenia:	1	1	1	1
&				
patient:	0	0	0	1

* bitwise & operation

	1	1	1	1
&	0	0	0	1
Result:	0	0	0	1

↳ D₄ to proof that use Cosine similarity

04/05/2025 01:25

Q: schizophrenia patients

$$Q: [0, 0, 0, 0, 0, 0, 0, 1, 1, 0]$$

$$D_4: [0, 0, 0, 1, 1, 1, 0, 1, 1, 0]$$

$$\text{Cos similarity} = \frac{Q \cdot D_4}{|Q| |D_4|}$$

$$Q \cdot D_4 = 1 \cdot 1 + 1 \cdot 1 + 0 = 2$$

$$|Q| = \sqrt{(1)^2 + (1)^2} = \sqrt{2} = 1.4$$

$$|D_4| = \sqrt{(1)^2 + (1)^2 + (1)^2 + (1)^2 + (1)^2} = \sqrt{5} = 2.236$$

$$\text{Cos similarity}(Q, D_4) = \frac{2}{3.13} = \boxed{0.6388}$$

Q with D_3

$$D_3 = [1, 0, 0, 1, 0, 1, 1, 0, 1, 1]$$

$$\frac{Q \cdot D_3}{|Q| |D_3|} = \frac{1}{1.4 \cdot 2.449} = \frac{1}{3.4292} = \boxed{0.2916}$$

$$|D_3| = \sqrt{6 \cdot (1)} = 2.449$$

Q with D_2

$$D_2 = [0, 0, 1, 0, 0, 1, 0, 0, 1, 0]$$

$$|D_2| = \sqrt{3} = 1.73$$

$$\text{Cos similarity}(Q, D_2) = \frac{Q \cdot D_2}{|Q| |D_2|}$$

$$\left. \begin{array}{l} Q \cdot D_2 = 1 \\ |Q| = 1.4 \\ |D_2| = 1.73 \end{array} \right\}$$

$$= \frac{1}{1.4 \times 1.73} = \frac{1}{2.424}$$

$$\text{Cos}(Q, D_2) = \boxed{0.412}$$

$$\text{Cos similarity}(Q, D_1) = \frac{Q \cdot D_1}{|Q| |D_1|}$$

$$D_1 = [0, 1, 1, 1, 0, 0, 0, 0, 1, 0]$$

$$|D_1| = \sqrt{4} = 2$$

$$|Q| = 1.4$$

$$\Rightarrow \text{Cos}(Q, D_1) = \frac{1}{2 \times 1.4}$$

$$Q \cdot D_1 = 1$$

$$\boxed{\text{Cos}(Q, D_1) = 0.357}$$

$$\text{Most similarity } \boxed{Q \text{ with } D_4 = 0.6388}$$

$$\text{then } Q \text{ with } D_2 = 0.412$$

$$\text{then } Q \text{ with } D_1 = 0.357$$

$$\text{then less similarity } Q \text{ with } D_3 = 0.2916$$

4. Perform the same computation for the following queries:

- (new AND for) OR (breakthrough AND drug)

Problem 3:

$Q_1: (\text{new AND for}) \text{ OR } (\text{breakthrough AND drug})$

doc	D_1	D_2	D_3	D_4
new	0	1	1	1
for	1	0	1	1
breakthrough	1	0	0	0
drug	1	1	0	0

new	0	1	1	1	breakthrough	1	0	0	0
& for	1	0	1	1	& drug	1	1	0	0
new & for :	0	0	1	1	breakthrough & drug :	1	0	0	0

OR

OR	0	0	1	1
OR	1	0	0	0
	1	0	1	1

$\rightarrow D_1, D_3, D_4$

$Q_1: [0, 1, 1, 1, 0, 1, 0, 0, 0, 0]$

$D_1: [0, 1, 1, 1, 0, 0, 0, 0, 0, 0]$

By Cos similarity :- $\cos(Q_1, D_1) = \frac{Q_1 \cdot D_1}{|Q_1| |D_1|}$

$\cos(Q_1, D_1) = \frac{1+1+1}{2 \times 2} = \frac{3}{4} = 0.75$

$|Q_1| = \sqrt{4 \times (1)^2} = 2$

$|D_1| = \sqrt{4} = 2$

$\cos(Q_1, D_2) = \frac{2}{2 \times 1.73} = 0.578$

$$\cos(Q_1, D_3) = \frac{1+1}{2.449 \cdot 2} = 0.408$$

$$\cos(Q_1, D_4) = \frac{2}{2.23 \cdot 2} = 0.4484$$

most similar $D_1 = 0.75$

then $D_2 = 0.578$
 then $D_4 = 0.4484$
 then $D_3 = 0.408$

- (new NOT schizophrenia)

Q_2 : new NOT schizophrenia

	D_1	D_2	D_3	D_4
new	0	1	1	1
schizophrenia	1	1	1	1

$Q_2 = [0, 0, 0, 0, 0, 1, 0, 0, 1, 0]$
 $|Q_2| = \sqrt{2}$

$$\cos(Q_2, D_1) = \frac{1}{2 \cdot \sqrt{2}} = 0.353$$

$$\cos(Q_2, D_2) = \frac{2}{\sqrt{2} \cdot 1.73} = 0.8174$$

$$\cos(Q_2, D_3) = \frac{2}{\sqrt{2} \cdot 2.449} = 0.57746$$

04/05/2025 01:57

$$\cos(Q_2, D_4) = \frac{2}{\sqrt{2} * 2.236} = 0.632$$

Q_2 similarity most with $D_2 = 0.8174$

then Q_2 with 0.632 D_4

then Q_2 with $D_3 = 0.5746$

then Q_2 with $D_1 = 0.353$

Matrix with Posting List

IDS	Posting List (DocIDs)
approach	3
breakthrough	1
drug	1, 2
for	1, 3, 4
hopes	4
new	2, 3, 4
of	3
patients	4
schizophrenia	1, 2, 3, 4
treatment	3

5. Draw the matrix with posting lists.

Problem4. For the following questions, use the term-by-document matrix provided below to perform vector space retrieval (blanks indicate zeros):

```

      d1 d2 d3
+---+---+---+
t1 |   | 3 | 5 |
+---+---+---+
t2 |   | 1 | 3 |
+---+---+---+ t3
| 5 | 4 | 4 | +-
--+---+---+ t4 |
6 | | 3 | |
+---+---+---+ t5
|   | 1 | 3 | +-
--+---+---+
t6 | 3 |   | 7 | +---+---+
--+
t7 |   | 6 | |
+---+---+---+
t8 | 2 |   | 2 |
+---+---+---+

```

1. Build the w matrix in which each element is computed as $TF * IDF$, where TF is what is specified in each element above and the IDF for term i is computed as:

2. Term-Document Matrix with TF-IDF weights:

Term	D1	D2	D3
T1	0	1.755	2.925
T2	0	0.585	1.755
T3	0	0	0
T4	3.51	1.755	0
T5	0	0.585	1.755
T6	1.755	0	4.095

Term	D1	D2	D3
T7	0	9.51	0
T8	1.17	0	1.17

- IDF Calculations (for $N=3$ documents):
 - Terms t_1, t_2, t_7 : $df=1 \rightarrow IDF = \log_2(3/1) \approx 1.585$
 - Terms t_4, t_5, t_6, t_8 : $df=2 \rightarrow IDF = \log_2(3/2) \approx 0.585$
 - Term t_3 : $df=3 \rightarrow IDF = \log_2(3/3) = 0$

Similarity Scores Calculation:

- For each document, sum the TF-IDF weights of the query terms t_2 and t_7 :
 - d_1 : $t_2=0+t_7=0=0$
 - d_2 : $t_2=0.585+t_7=9.51 \approx 10.095$
 - d_3 : $t_2=1.755+t_7=0 \approx 1.755$

$\log(\text{total_number_of_documents} / \text{number_of_documents_containing_term_i})$. Use **base 2 logarithms**.

Compute the rank order of the documents that would be found using the vector space method for the UNWEIGHTED query (terms are of equal importance): t_2 t_7

Ranking:

1. d_2 (Score ≈ 10.095)
2. d_3 (Score ≈ 1.755)
3. d_1 (Score $=0$)

Final Ranked List:

$d_2 > d_3 > d_1$ $d_2 > d_3 > d_1$

In your answer, give the similarity score that is computed for each document and give the ranked list that the system would return.

Problem 5.

Given a collection that contains 5 different words a, b, c, d, e only.

The frequency order is $f(a) > f(b) > f(c) > f(d) > f(e)$. The total number of tokens in the collection is 100000.

- 2-a. Assume that Zipf's law holds exactly for this collection. What are the collection frequencies of the 5 words?

Problem 5:-

unique word = 5 (a, b, c, d, e)

token word = 100,000 word

$f(a) > f(b) > f(c) > f(d) > f(e)$

① assume that Zipf's Law holds exactly for this collection what are the collection frequency of the 5 word?

rank of a = 1 $\Rightarrow f(a) = A/1 = 43795.620$

rank of b = 2 $\Rightarrow f(b) = A/2 = 21898$

rank of c = 3 $\Rightarrow f(c) = A/3 = 14599$

rank of d = 4 $\Rightarrow f(d) = A/4 = 10949$

rank of e = 5 $\Rightarrow f(e) = A/5 = 8759$

total frequency of 5 word = 100,000

Most frequency $\Rightarrow \text{freq.} = A$

$$\frac{A}{1} + \frac{A}{2} + \frac{A}{3} + \frac{A}{4} + \frac{A}{5} = 100,000$$

$$A * 2.283 = 100,000 \Rightarrow A = 43795.62044$$

2-b. Can you estimate the postings list size for each of these words: positional and nonpositional.

- 5 words: a, b, c, d, e
- Word frequencies according to Zipf's law:
 - a: 100,000
 - b: 50,000
 - c: 33,333
 - d: 25,000
 - e: 20,000
- Total tokens in collection: 100,000

Necessary Assumptions:

1. Number of documents (N): Assuming average document length of 100 words
 $N = 100,000/100 = 1,000$ documents

2. Document ID size: 10 bits (to represent 1,000 documents)
3. Position index size: 8 bits (to represent word positions within documents)

1. Non-positional Postings Lists:

Each entry contains:

- Only document ID (10 bits)

Estimating number of documents containing each term (df):

- Word a: Appears in nearly all documents ($df \approx 1,000$)
- Word b: Appears in ≈ 500 documents
- Word c: Appears in ≈ 333 documents
- Word d: Appears in ≈ 250 documents
- Word e: Appears in ≈ 200 documents

List sizes:

List size = $df \times 10$ bits

- a: $1,000 \times 10 = 10,000$ bits ≈ 1.25 KB
- b: $500 \times 10 = 5,000$ bits ≈ 0.625 KB
- c: $333 \times 10 \approx 3,330$ bits ≈ 0.416 KB
- d: $250 \times 10 = 2,500$ bits ≈ 0.312 KB
- e: $200 \times 10 = 2,000$ bits ≈ 0.25 KB

2. Positional Postings Lists:

Each entry contains:

- Document ID (10 bits)
- Number of positions (variable)
- Actual positions (8 bits per position)

Estimating positions per word:

- Assuming each word appears once on average in documents where it occurs

List sizes:

List size $\approx df \times (10 + 8 \times \text{average frequency per document})$

- a: $1,000 \times (10 + 8 \times 100) \approx 1,000 \times 810 = 810,000$ bits ≈ 101 KB
- b: $500 \times (10 + 8 \times 50) = 500 \times 410 = 205,000$ bits ≈ 25.6 KB
- c: $333 \times (10 + 8 \times 33) \approx 333 \times 274 \approx 91,242$ bits ≈ 11.4 KB
- d: $250 \times (10 + 8 \times 25) = 250 \times 210 = 52,500$ bits ≈ 6.56 KB
- e: $200 \times (10 + 8 \times 20) = 200 \times 170 = 34,000$ bits ≈ 4.25 KB

Problem 6.

2-a. A search engine has a collection of 16,000,000 pages (documents) with 200 tokens per page, on average.

What is the minimal length for document IDs for the postings? In bits and in full bytes.

2-b. If the vocabulary size is 400,000, and the average dictionary word length is 7 characters.

How many bits do you need for pointers if one is to store the dictionary as a single string with pointers to the start of each word (what is the length of each pointer).

2-c. What is the size of the incidence matrix? Can you estimate the number of nonzero elements in the incidence matrix? What about the number of nonpositional postings?

Problem 6:

(a) minimal length for document IDs for the posting?

$$\log_2(16000000) \approx 23.93 \approx 24 \text{ bit} \\ \approx 3 \text{ Byte}$$

(b) $V \text{ size} = 400000$, avg word length = 7 character

Length of pointer = $V \text{ size} * \text{word length}$

$$\text{Length of each pointer} = 400000 * 7$$

$$\text{Length of each pointer} = 2800000 \text{ bits}$$

$$\text{Length of each pointer} = 350000 \text{ Byte}$$

(c) What is the size of the index matrix?

$$\text{Size} = \# \text{document} * \# \text{Word} = 400000 * 1600000 = 640000000000$$

$$\text{Size of Byte} = 800000000000 \text{ Byte} = 800 \text{ GB}$$

* Can you estimate the number of nonzero element in the incidence matrix?

$$\# \text{token} = 200 \quad \text{so number of nonzero element} = \\ 1600000 * 200 = 320000000$$

What about the number of nonposting and Posting?

number of nonposting = number of nonzero

$$= 320000000 = 3.2 \text{ billion}$$

Problem 7.

3-a. Compute the γ -code for the sequence of decimals 1, 7, 14, 49, 127, 1023, 2047.

(Problem 7)

(a) 1, 7, 14, 49, 127, 1023, 2047

① 1:

$$\begin{aligned} \log_2 1 &= 0 \rightsquigarrow 0 \\ 1 - 2^0 &= 0 \Rightarrow 0 \end{aligned} \quad \left. \vphantom{\begin{aligned} \log_2 1 &= 0 \rightsquigarrow 0 \\ 1 - 2^0 &= 0 \Rightarrow 0 \end{aligned}} \right\} 0$$

② 7:

$$\begin{aligned} \log_2 7 &= 2 \rightsquigarrow 110 \\ 7 - 2^2 &= 3 \xrightarrow{\text{binary}} 11 \end{aligned} \quad \left. \vphantom{\begin{aligned} \log_2 7 &= 2 \rightsquigarrow 110 \\ 7 - 2^2 &= 3 \xrightarrow{\text{binary}} 11 \end{aligned}} \right\} 11011$$

③ 14

$$\begin{aligned} \log_2 14 &= 3 \rightsquigarrow 1110 \\ 14 - 2^3 &= 6 \xrightarrow{\text{binary}} 110 \end{aligned} \quad \left. \vphantom{\begin{aligned} \log_2 14 &= 3 \rightsquigarrow 1110 \\ 14 - 2^3 &= 6 \xrightarrow{\text{binary}} 110 \end{aligned}} \right\} 1110110$$

④ 49:

$$\begin{aligned} \log_2 49 &= 5 \rightsquigarrow 111110 \\ 49 - 2^5 &= 17 \rightsquigarrow 10001 \end{aligned} \quad \left. \vphantom{\begin{aligned} \log_2 49 &= 5 \rightsquigarrow 111110 \\ 49 - 2^5 &= 17 \rightsquigarrow 10001 \end{aligned}} \right\} 11111010001$$

⑤ 127:

$$\begin{aligned} \log_2 127 &= 6 \rightsquigarrow 1111110 \\ 127 - 2^6 &= 63 \rightsquigarrow 111111 \end{aligned} \quad \left. \vphantom{\begin{aligned} \log_2 127 &= 6 \rightsquigarrow 1111110 \\ 127 - 2^6 &= 63 \rightsquigarrow 111111 \end{aligned}} \right\} 111111011111$$

$$\begin{aligned}
 &1023 \\
 &\log_2 1023 \approx 9 \rightarrow 111111110 \\
 &1023 - 2^9 = 511 \rightarrow 110111111 \\
 &1023 : 111111101111111 \\
 &2047 : 1111111101111111 \\
 &2^{11} = 2048 \text{ so } 10 \text{ bits}
 \end{aligned}$$

3-b. Represent the least significant 5 decimal digits of your ID number in γ -code.

(b) Last 5 decimal in ID student 1211439

$$\begin{aligned}
 &1 : 0 \\
 &1 : 0 \\
 &4 : \log_2 4 = 2 \rightarrow 110 \quad \left. \begin{array}{l} 110 \\ 0 \end{array} \right\} 1100 \\
 &\quad 4 - 2^2 = 0 \rightarrow 0 \quad \text{offset: Binary} \\
 &3 : \log_2 3 = 1 \rightarrow 10 \quad \left. \begin{array}{l} 10 \\ 01 \end{array} \right\} 1001 \\
 &\quad 3 - 2^1 = 1 \rightarrow 01 \quad \text{Binary} \\
 &9 : \log_2 9 = 3 \rightarrow 1110 \quad \left. \begin{array}{l} 1110 \\ 01 \end{array} \right\} 111001 \\
 &\quad 9 - 2^3 = 1 \rightarrow 01 \quad \text{Binary} \\
 &\gamma\text{-code: } 001100100111001
 \end{aligned}$$

04/05/2025 02:33

3-c. Recover the gap value in decimal for the following string representing a sequence of gaps in a posting list, if possible.

1111011111110101111110101011110101111111011101010

③

$\overbrace{1111}^4$ $\overbrace{011111}^{63}$ $\overbrace{110}^2$ $\overbrace{10}^2$ $\overbrace{1111110}^7$ $\overbrace{1010111}^8$

$2^4 + 63$ $2^2 + 2$ $2^7 + 87$
 $= 16 + 63$ $= 6$ $= 215$
 $= 79$

$\overbrace{10}^1$ $\overbrace{1011111110}^8$ $\overbrace{11101010}^8$

$2^1 + 1$ $2^8 + 234$
 $= 3$ $= 490$

So Result : 79, 6, 215, 3, 490

Good Luck