

③ (a) (Lexicographic order) $N=5$

Term	d_f	$idf = \ln\left(\frac{N}{d_f}\right)$	$tf(d_1)$	$tf(d_2)$	$tf(d_3)$	$tf(d_4)$	$tf(d_5)$	$tf \cdot idf(d_1)$	$tf \cdot idf(d_2)$	$tf \cdot idf(d_3)$	$tf \cdot idf(d_4)$	$tf \cdot idf(d_5)$
and	1	$\ln\left(\frac{5}{1}\right) = 1.61$	0	0	1	0	0	$0 * \ln\left(\frac{5}{1}\right) = 0$	$0 * \ln\left(\frac{5}{1}\right) = 0$	$1 * \ln\left(\frac{5}{1}\right) = 1.61$	$0 * \ln\left(\frac{5}{1}\right) = 0$	$0 * \ln\left(\frac{5}{1}\right) = 0$
cat	4	$\ln\left(\frac{5}{4}\right) = 0.22$	0	1	1	1	2	$0 * \ln\left(\frac{5}{4}\right) = 0$	$1 * \ln\left(\frac{5}{4}\right) = 0.22$	$1 * \ln\left(\frac{5}{4}\right) = 0.22$	$1 * \ln\left(\frac{5}{4}\right) = 0.22$	$2 * \ln\left(\frac{5}{4}\right) = 0.45$
dog	3	$\ln\left(\frac{5}{3}\right) = 0.51$	1	0	2	1	0	$1 * \ln\left(\frac{5}{3}\right) = 0.51$	$0 * \ln\left(\frac{5}{3}\right) = 0$	$2 * \ln\left(\frac{5}{3}\right) = 1.02$	$1 * \ln\left(\frac{5}{3}\right) = 0.51$	$0 * \ln\left(\frac{5}{3}\right) = 0$
eat	3	$\ln\left(\frac{5}{3}\right) = 0.51$	1	1	1	0	0	$1 * \ln\left(\frac{5}{3}\right) = 0.51$	$1 * \ln\left(\frac{5}{3}\right) = 0.51$	$1 * \ln\left(\frac{5}{3}\right) = 0.51$	$0 * \ln\left(\frac{5}{3}\right) = 0$	$0 * \ln\left(\frac{5}{3}\right) = 0$
homework	2	$\ln\left(\frac{5}{2}\right) = 0.92$	1	1	0	0	0	$1 * \ln\left(\frac{5}{2}\right) = 0.92$	$1 * \ln\left(\frac{5}{2}\right) = 0.92$	$0 * \ln\left(\frac{5}{2}\right) = 0$	$0 * \ln\left(\frac{5}{2}\right) = 0$	$0 * \ln\left(\frac{5}{2}\right) = 0$
host	1	$\ln\left(\frac{5}{1}\right) = 1.61$	0	0	1	0	0	$0 * \ln\left(\frac{5}{1}\right) = 0$	$0 * \ln\left(\frac{5}{1}\right) = 0$	$1 * \ln\left(\frac{5}{1}\right) = 1.61$	$0 * \ln\left(\frac{5}{1}\right) = 0$	$0 * \ln\left(\frac{5}{1}\right) = 0$
play	2	$\ln\left(\frac{5}{2}\right) = 0.92$	0	0	0	1	1	$0 * \ln\left(\frac{5}{2}\right) = 0$	$0 * \ln\left(\frac{5}{2}\right) = 0$	$0 * \ln\left(\frac{5}{2}\right) = 0$	$1 * \ln\left(\frac{5}{2}\right) = 0.92$	$1 * \ln\left(\frac{5}{2}\right) = 0.92$
the	5	$\ln\left(\frac{5}{5}\right) = 0$	1	1	3	2	2	$1 * \ln\left(\frac{5}{5}\right) = 0$	$1 * \ln\left(\frac{5}{5}\right) = 0$	$3 * \ln\left(\frac{5}{5}\right) = 0$	$2 * \ln\left(\frac{5}{5}\right) = 0$	$2 * \ln\left(\frac{5}{5}\right) = 0$
with	2	$\ln\left(\frac{5}{2}\right) = 0.92$	0	0	0	1	1	$0 * \ln\left(\frac{5}{2}\right) = 0$	$0 * \ln\left(\frac{5}{2}\right) = 0$	$0 * \ln\left(\frac{5}{2}\right) = 0$	$1 * \ln\left(\frac{5}{2}\right) = 0.92$	$1 * \ln\left(\frac{5}{2}\right) = 0.92$

* here $tf(d_n)$ signifies tf for a particular document n .
* \log has been written as \ln i.e. \log to the base e .

* All values rounded off to 2 decimals.

All terms (words) are written in the Lexicographic order in the first column.

df = It is the document frequency i.e. no. of documents in the collection containing the term.

tf = It is the frequency (count) of a term in document.

N = Total number of documents in the collection.

Since we have 5 documents d_1, d_2, d_3, d_4, d_5 . Hence $N = 5$

$$tfidf = tf \cdot \ln \left(\frac{N}{df} \right) \quad \text{Here } \ln() \text{ represent log to the base } e.$$

All calculations & values are mentioned in the Table.

Eg. For the term "and", $df = 1$, as there is only 1 document i.e. d_3 , in which this term is present. idf can be easily calculated using $\ln \left(\frac{N}{df} \right)$, if we know the value of df .

$tf(d_i)$ represents tf for document d_i . It is 0 as the term "and" is absent in document d_1 i.e. it is present 0 times in the document d_1 . Similarly all calculations are done, as shown in table.

So final document vectors using the $tfidf$ weights are -

$$\vec{d_1} = [0, 0, 0.51, 0.51, 0.92, 0, 0, 0, 0]$$

$$\vec{d_2} = [0, 0.22, 0, 0.51, 0.92, 0, 0, 0, 0]$$

$$\vec{d_3} = [1.61, 0.22, 1.02, 0.51, 0, 1.61, 0, 0, 0]$$

$$\vec{d_4} = [0, 0.22, 0.51, 0, 0, 0, 0.92, 0, 0.92]$$

$$\vec{d_5} = [0, 0.45, 0, 0, 0, 0, 0.92, 0, 0.92]$$

③ (b) Cosine Similarity between vectors d_x & $d_y \Rightarrow$

$$\text{sim}(d_x, d_y) = \frac{d_x \cdot d_y}{\|d_x\| \|d_y\|}$$

Here $d_x \cdot d_y$ represents dot product of d_x & d_y .

$\|d_x\|$ represents square root of sum of squares of

"terms" of different terms present in the document d_x .

Similar meaning is for $\|d_y\|$.

So cosine similarity between vectors for documents d_1 & $d_2 \Rightarrow$

$$\text{sim}(d_1, d_2) = \frac{d_1 \cdot d_2}{\|d_1\| \|d_2\|}$$

$d_1 \cdot d_2$ as per values taken Table

$$= (0 \times 0) + (0 \times 0.22) + (0.51 \times 0) + (0.51 \times 0.51) + (0.92 \times 0.92) + (0 \times 0) + (0 \times 0) + (0 \times 0) + (0 \times 0)$$

$$= 0 + 0 + 0 + 0.2601 + 0.8464 + 0 + 0 + 0 + 0$$

$$= 1.1065$$

$$\|d_1\| = \sqrt{0^2 + 0^2 + (0.51)^2 + (0.51)^2 + (0.92)^2 + 0^2 + 0^2 + 0^2 + 0^2}$$

$$= \sqrt{0 + 0 + 0.2601 + 0.2601 + 0.8464 + 0 + 0 + 0} = \sqrt{1.3666}$$

$$\|d_2\| = \sqrt{0^2 + (0.22)^2 + 0^2 + (0.51)^2 + (0.92)^2 + 0^2 + 0^2 + 0^2 + 0^2}$$

$$= \sqrt{0 + 0.0484 + 0 + 0.2601 + 0.8464 + 0 + 0 + 0 + 0} = \sqrt{1.1549}$$

$$\text{So sim}(d_1, d_2) = \frac{1.1065}{\sqrt{1.3666} \times \sqrt{1.1549}} = 0.88$$

Similarly Cosine Similarity between vectors for documents d_1 & $d_4 \Rightarrow$

$$\text{sim}(d_1, d_4) = \frac{d_1 \cdot d_4}{\|d_1\| \cdot \|d_4\|}$$

$d_1 \cdot d_4$ as per values taken from table \Rightarrow

$$= (0 \times 0) + (0 \times 0.22) + (0.51 \times 0.51) + (0.51 \times 0) + (0.92 \times 0) \\ + (0 \times 0) + (0 \times 0.92) + (0 \times 0) + (0 \times 0.92)$$

$$= 0 + 0 + 0.2601 + 0 + 0 + 0 + 0 + 0 + 0$$

$$= 0.2601$$

$\|d_1\|$ as already calculated is $\sqrt{1.3666}$

$$\|d_4\| = \sqrt{0^2 + (0.22)^2 + (0.51)^2 + 0^2 + 0^2 + 0^2 + (0.92)^2 + 0^2 + (0.92)^2}$$

$$= \sqrt{0 + 0.0484 + 0.2601 + 0 + 0 + 0 + 0.8464 + 0 + 0.8464}$$

$$= \sqrt{2.0013}$$

$$\text{So } \text{sim}(d_1, d_4) = \frac{0.2601}{\sqrt{1.3666} \times \sqrt{2.0013}} = 0.16$$

More the value of similarity^(Sim), lesser the angle between the two vectors & more is the similarity between them.

$$\text{Since } \text{sim}(d_1, d_2) > \text{sim}(d_1, d_4)$$

Hence d_1 is more similar to d_2 .