

ANS: 3

$$q) \text{tfidf} = \text{tf} \cdot \log\left(\frac{N}{\text{df}}\right)$$

Now, removing stop words such as

$d_1 = \text{dog ate homework}$

$d_2 = \text{cat ate homework}$

$d_3 = \text{dog cat ate hot dog}$

$d_4 = \text{dog cat wrote homework}$

$d_5 = \text{yesterday hot day}$

$$\rightarrow \text{idf} = \log(N/\text{df})$$

here  $N=5$  where  $n = \text{number of documents}$

$N/\text{df}$ (for all words)	$\text{idf} = \log(N/\text{df})$
ate - 5/3	0.51
cat - 5/2	0.92
day - 5/1	1.61
dog - 5/4	0.22
homework - 5/3	0.51
hot - 5/1	1.61
wrote - 5/1	1.61
yesterday - 5/1	1.61

TF:

	ate	cat	day	dog	homework	hot	wrote	yesterday
d1	1	0	0	1	1	0	0	0
d2	1	1	0	0	1	0	0	0
d3	1	1	0	2	0	0	0	0
d4	0	1	0	1	1	0	1	0
d5	0	0	1	0	0	1	0	1

- Now we multiply TF with IDF

	ate	cat	day	dog	homework	hot	wrote	yesterday
d1	0.51	0	0	0.22	0.51	0	0	0
d2	0.51	0.92	0	0	0.51	0	0	0
d3	0.51	0.92	0	0.44	0	0	0	0
d4	0	0.92	0	0.22	0.51	0	1.61	0
d5	0	0	1.61	0	0	1.61	0	1.61



NOW we can rewrite our results (document terms vs documents)

terms	documents				
	0.51	0.51	0.51	0	0
	0	0.92	0.92	0.92	0
	0	0	0	0	1.61
	0.22	0	0.44	0.22	0
	0.51	0.51	0	0.51	0
	0	0	0	0	1.61
	0	0	0	1.61	0
	0	0	0	0	1.61

b)

1. cosine similarity bet<sup>n</sup>  
D<sub>1</sub> & D<sub>2</sub>

$$\text{Sim}(d_1, d_2) = \frac{\vec{d}_1 \cdot \vec{d}_2}{|\vec{d}_1| |\vec{d}_2|}$$

$$= 2/3$$

$$= 0.67$$

here D<sub>1</sub> is more similar  
to D<sub>2</sub>

2. cosine similarity (D<sub>1</sub> & D<sub>5</sub>)

$$\text{Sim}(d_1, d_5) = \frac{\vec{d}_1 \cdot \vec{d}_5}{|\vec{d}_1| |\vec{d}_5|}$$

$$= 0/3$$

$$= 0$$

here we can say that  
D<sub>1</sub> & D<sub>5</sub> are totally different.