

Week 7 NLP live session

(TF) - IDF
└→

TF: term frequency

IDF: Inverse Document frequency

TF-IDF (t, d) = TF(t, d) * IDF(t)

statistical measure

specific to single document

count of term
within a
document

depends on the entire
corpus

remains same
over the entire
corpus

↗ term count within the doc.

$$\underline{\text{TF}}(t, d) = \frac{\text{no. of times term 't' appears in doc. 'D'}}{\text{No. of terms in doc. 'D'}}$$

$$\underline{\text{IDF}}(t) = \log \left(\frac{\text{No. of documents in the corpus} \cancel{\text{that contain 't'}}}{\text{No. of doc in which term 't' appears}} \right)$$

Q1: consider a ^ddocument containing 100 words wherein the word 'cat' appears 3 times.

Now, suppose in the corpus there are 10 million documents Δ the word 'cat' appears in 1000 of these doc.

Calculate $TF - IDF (cat, d)$

$1 \text{ million} = 10^6$

$$tf(cat, 'd') = \frac{\text{no. of times 't' appears in d}}{\text{no. of terms in d}} = \frac{3}{100}$$

$$= 0.03$$

\log \downarrow

$$\text{IDF}(t) = \frac{\text{no. of doc. in the corpus}}{\text{no. of doc. in which term 'cat' appears}}$$

\log \downarrow

$$= \frac{10 \times 10^6}{1000} = 10^4$$

\log \downarrow

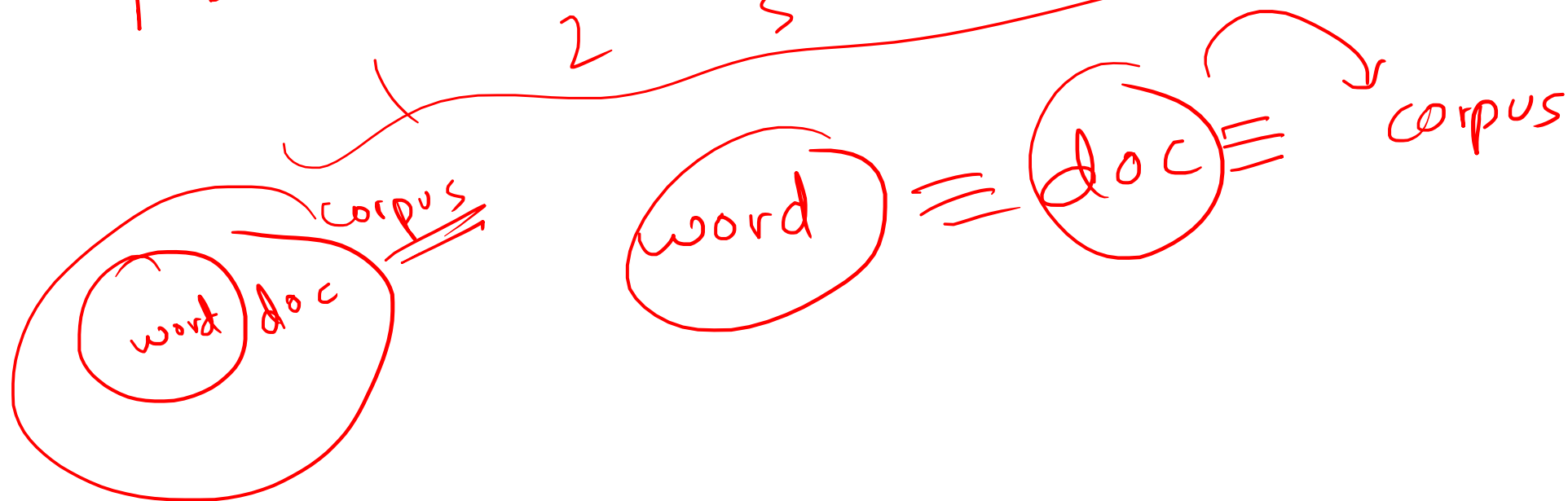
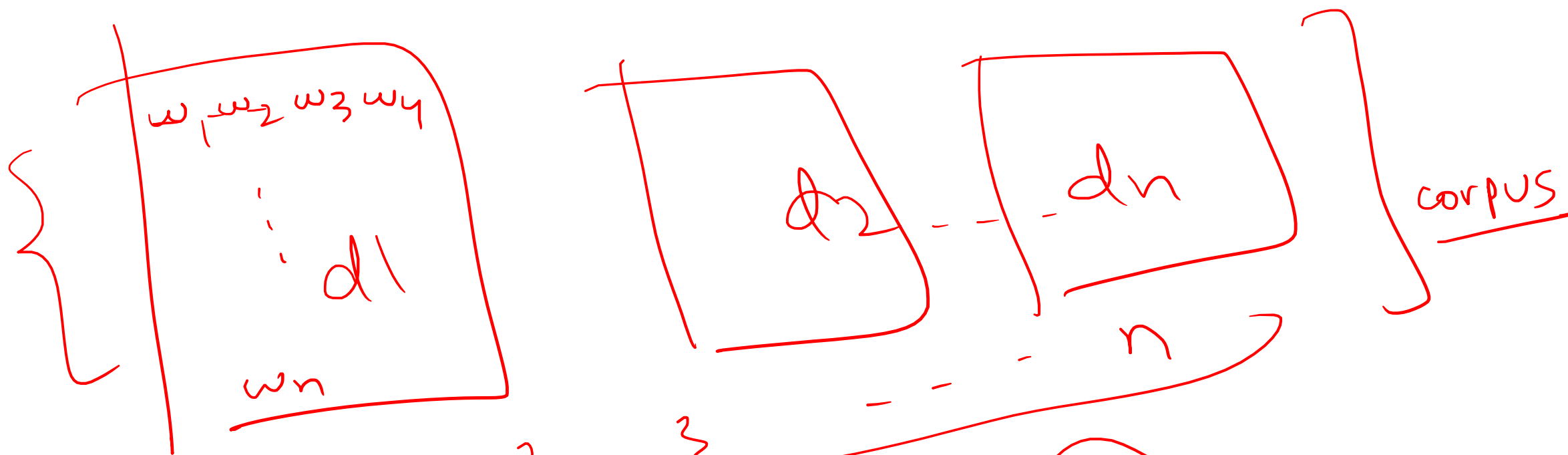
$$= \log_{10}(10^4) = 4$$

\log \downarrow

$$\log\left(\frac{10^7}{10^3}\right) = \log_{10}(10^4) = 4$$

$$T_{F-JDF}(t, d) = \frac{T_F(t, d) * JDF(t)}{JDF(t)}$$

$$= \left(\frac{3}{100} \right) * 4 = 0.12$$



Q2) D1: A quick brown fox jumps over the lazy dog.
what a fox! ✓ $|D1| = 12$

D2: A quick brown fox jumps over the lazy fox.
What a fox! ✓ $|D2| = 12$

~~D3~~: How word fox is relevant to corpus D doc?

Hint: $TF - IDF(\text{fox}, \underline{D})$

$$TF(fox, D1) = \frac{\text{No. of times fox appears in } D1}{\text{No. of terms in } D1 (|D1|)}$$

$$= \frac{2}{12} = \frac{1}{6} = 0.17$$

$$TF(fox, D2) = \frac{3}{|D2|} = \frac{3}{12} = \frac{1}{4} = \underline{0.25}$$

$$D = \{D_1, D_2\}$$

$$IDF(\text{fox}, 'D') = \frac{\log(\text{No. of doc. in corpus 'D'})}{\log(\text{No. of doc in which 'fox' appears})}$$

$$= \log \left(\frac{2}{2} \right)$$

$$= 0$$

$$\text{TF-IDF}(\underline{\text{fox}}, \underline{D1}) = \text{TF}(\underline{\text{fox}}, D1) * \text{IDF}(\underline{\text{fox}})$$

$$= 0.17 * 0 = 0 \quad \checkmark$$

$$\text{TF-IDF}(\underline{\text{fox}}, \underline{D2}) = \text{TF}(\underline{\text{fox}}, D2) * \text{IDF}(\underline{\text{fox}})$$

$$= 0.25 * 0 = 0 \quad \checkmark$$

~~Ans:~~ The word 'fox' is equally relevant to corpus D
 $\therefore \text{TF-IDF}(\text{fox}, D1) = \text{TF-IDF}(\text{fox}, D2)$ is same

tf: simple choice (raw count of a term in a doc.)

idf: how much infoⁿ the word provides in
our corpus

TF-IDF : statistical measurement
↳ order / sequence of words (~~semantic~~)

Q3: d1: the man went out for a walk.

d2: the children sat around the fire

✓ ~~a) TF IDF (the, D)~~ = 0

✓ b) TF IDF (fire, D) $\frac{1}{6}(\log 2)$

✓ c) TF IDF (children, D) $\frac{1}{6}(\log 2)$

$$TF(the, d_1) = \frac{1}{7}, \quad TF(the, d_2) = \frac{2}{6} = \frac{1}{3}$$

$$IDF(the, D) = \log\left(\frac{2}{2}\right) = 0$$

$$TF(\text{fire}, D1) = \underline{1} \checkmark$$

$$TF(\text{fire}, D2) = \frac{1}{6}$$

$$IDF(\text{fire}, D) = \log\left(\frac{2}{1}\right) = \log_{10} 2 = \underline{0.3010}$$

$$TFIDF(\text{fire}, D1) = 0$$

$$TFIDF(\text{fire}, D2) = \frac{1}{6} * \log 2$$

R1: This movie is very scary and long

R2: This movie isn't scary and is slow.

R3: This movie is spooky and good.

Bag of words:

Build a vocabulary \equiv unique words from all doc.

11

	this	movie	is	very	scary	and	long	not	slow	spooky	good
R1	1	1	1	1	1	1	1	0	0	0	0
R2	1	1	2	0	1	1	0	1	1	0	0
R3	1	1	1	0	0	1	0	0	0	1	1

this

✓

Bow

Bag of words Vector Representation

R1: [1 1 1 1 1 1 1 1 0 0 0 0] sparse

R2: [1 1 2 0 1 1 0 1 1 0 0]

R3: [1 1 1 0 0 1 0 0 0 0 1 1]

$$\textcircled{*} \text{Jaccard similarity } (A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$$** \text{Dice coefficient} = \frac{2|A \cap B|}{|A| + |B|}$$

$$|A| =$$

(A) S1: the cat sat on the mat

(B) S2: the kitten rested on the rug

the, cat, sat, on, ^{the} mat, kitten, rested, rug

$$|B| =$$

$$IS(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{2}{8} = \frac{1}{4} = 0.25$$

$$\frac{3}{4} = 0.75$$

