

NLP Sheet 4

- 1) . $N = 10,000$ documents
• One of them has 250 words where 'sheet' occurs 20 times (call it D)
• 'Sheet' also occurs in 2500 documents (overall)
→ What is tf-idf for 'sheet' in a bag of words representation of D .

Recall

$$tf_{t,d} = \frac{\text{Count}(t,d)}{\text{len}(d)}$$

$$\rightarrow tf_{\text{sheet},D} = \frac{20}{250} = \frac{2}{25}$$

$$idf_t = \log_{10}(N / df_t)$$

↓
No. of docs where t appears

$$\rightarrow idf_{\text{sheet}} = \log_{10}\left(\frac{10,000}{2500}\right) = 0.6$$

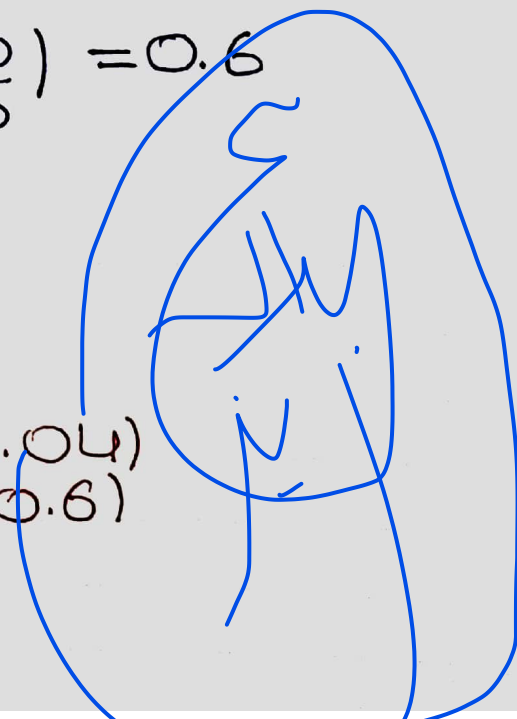
$$\bullet \text{tf-idf} = \text{tf} * \text{idf} = 0.048$$

2) Given word vectors

$$w_1 = (1, 0.4, 0.3, 0.8, 0.04)$$

$$w_2 = (1, 0.2, 0.5, 0.7, 0.6)$$

• Are they similar or not?



$$\cos(w_1, w_2) = \frac{1.1 + 0.4 \cdot 0.2 + 0.3 \cdot 0.5 + 0.8 \cdot 0.7 + 0.04 \cdot 0.6}{\sqrt{1^2 + 0.4^2 + 0.3^2 + 0.8^2 + 0.04^2} \sqrt{1 + 0.2^2 + 0.5^2 + 0.7^2 + 0.6^2}}$$

$$= 0.96$$

• They are similar (slightly below max ^{1.0} similarity)

• Make one change to reduce similarity
 → let's try inverting the biggest term
 (causes largest reduction in $\cos(\theta)$)

$$\rightarrow \text{Now } \cos(w_1, w_2) = -0.09$$

• $-0.09 \ll 1$ and we can no longer say the two words are similar.

3. Given are 3 sentences

→ Perform text normalization (Stop word removal & lemmatization)

D1: natural language ^{PROCESS} ~~processing~~ is ^{become} ~~becoming~~ important since
 soon we will begin ^{talk} ~~talking~~ to our ^{computer} ~~computers~~.

D2: If ^{computer} ~~computers~~ understand natural language they will become
 much ^{simple} ~~simpler~~ to use

D3: Speech recognition is the first step to build ^{computer} ~~computers~~
 like us

Vocab: natural, language, process, become, important, begin, talk, computer,
 understand, simple, use, speech, recognition, first, step, build, like

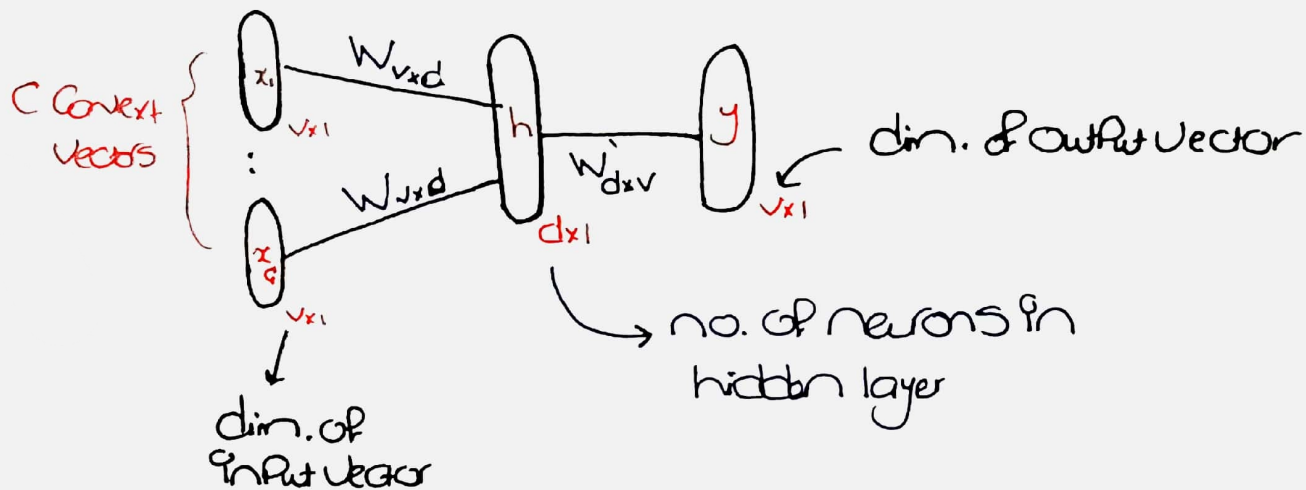
b) what is document vector for $D_3 \rightarrow$ use freq. counts ^①
 \rightarrow then tf-idf ^②
 \rightarrow has $|V|$ elements (for each word in vocab.)

vocab.	D_3	tf_{t,D_3}	idf_t	tf-idf
natural	○			○
language	○			○
Process	○			○
become	○			○
important	○			○
begin	○			○
talk	○			○
Computer	1	1/7	$\log_{10}(\frac{3}{3})$	○
understand	○		$\nwarrow N_{docs}$	○
Simple	○		$\nwarrow docs \text{ with Computer}$	○
use	0			○
Speech	1	1/7	$\log_{10}(3/1)$	0.068
Recognition	1	1/7	$\log_{10}(3/1)$	0.068
First	1	1/7	$\log_{10}(3/1)$	0.068
Step	1	1/7	$\log_{10}(3/1)$	0.068
build	1	1/7	$\log_{10}(3/1)$	0.068
like	1	1/7	$\log_{10}(3/1)$	0.068
	• Document vector (frequency count)	• normalized frequency counts (D_3 has 7 words)	• tf-idf vector	○

4. Consider CBOW model

- V words in the vocab.
- d is embedding dimension

• Draw the architecture



• For Skipgram each of the dimensions is the same ($v, v \times d, d, d \times v, v$)

→ but only 1 input vector (target word)

→ 1 or 2 output vectors depending on formulation

5. Find big-O runtime of computing a single Prob.

$P(\text{context} = c | \text{word} = w)$ for Skipgram

→ Write it in terms of embedding dim d and vocab size V

• To multiply $A_{k \times m} \times B_{m \times n}$ it takes

→ we do kn dot products (for each elem. in the result)

→ each dot product involves two m vectors

- m multiplications
- $m-1$ additions

- hence, overall the matrix multiplication corresponds to Knm multiplications
 → will disregard the $Kn(m-1)$ additions since the Knm multiplications dominate the complexity.

• SKIP-gram:

Big-oh

$$h = x_{1 \times V} W_{V \times d} \rightarrow d(V) \text{ multiplications}$$

$$u_c = h_{1 \times d} W'_{d \times V} \rightarrow V(d) \text{ multiplications}$$

$$j_c = \text{softmax}(u_c) \rightarrow V \text{ multiplications}$$

• Hence, Overall Complexity is $O(dV)$

6. users A & B have used word2vec on a specific vocabulary

- Each of them has obtained two word vectors u_w, v_w for each word in the vocabulary.
 → For w as Context → For w as Center (target)

- If for every two words w, w' in V it holds that

$$u_w^A \cdot v_{w'}^A = u_w^B \cdot v_{w'}^B \quad (1)$$

- Can we claim $v_w^A = v_w^B \quad \forall w \in V$

No.

• It's obvious that if $a.b = c.d$ then that doesn't imply $c=d$

• hence, can make the following counterexample

let vocab:

	by A		by B	
	U_w^A	V_w^A	U_w^B	V_w^B
Maltahar	(1 1)	(0.5 0.5)	(2 2)	($\frac{1}{3}$ $\frac{1}{3}$)
Thresh	(2 2)	(1 0)	(3 3)	($\frac{1}{2}$ 0)

• There's only two Pairs of words

(Malt, thresh):

$$\begin{matrix} U^A & V^A \\ (1 & 1) \end{matrix} \cdot \begin{matrix} U^A & V^A \\ (1 & 0) \end{matrix} = \begin{matrix} U^B & V^B \\ (2 & 2) \end{matrix} \cdot \begin{matrix} U^B & V^B \\ (\frac{1}{2} & 0) \end{matrix}$$

(thresh, malt):

$$\begin{matrix} U^A & V^A \\ (2 & 2) \end{matrix} \cdot \begin{matrix} U^A & V^A \\ (\frac{1}{2} & \frac{1}{2}) \end{matrix} = \begin{matrix} U^B & V^B \\ (3 & 3) \end{matrix} \cdot \begin{matrix} U^B & V^B \\ (\frac{1}{3} & \frac{1}{3}) \end{matrix}$$

* hence, ① from last page holds for every pair of words

• however, clearly for each word

$$\rightarrow V_w^A \neq V_w^B$$

↓
Thresh Malt

Q.E.D

7. Compare word2vec to co-occurrence based count methods

Word2vec

co-occurrence

disadv. { • Only takes local statistics into account

• Global & local* statistics (e.g. word-word matrix)

adv. { • Dense, short vector

• Sparse, long vector