



Tell Me This 20 hours ago (edited)

Human: What do we want!?

Computer: Natural language processing!

Human: When do we want it!?

Computer: When do we want what?

Reply • 203

[View reply](#) ▾



↪ Gretchen McCulloch Retweeted

Emma Manning (they/t... · 8h

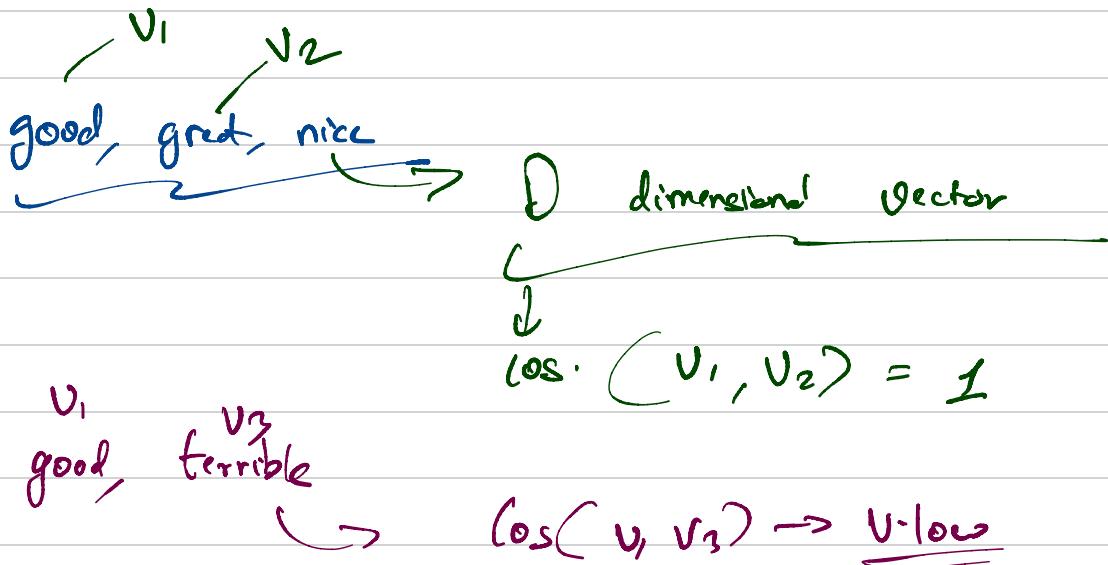
I like how the difference between "language processing" and "natural language processing" is that the one with "natural" is done by computers

2 77 361

Agenda

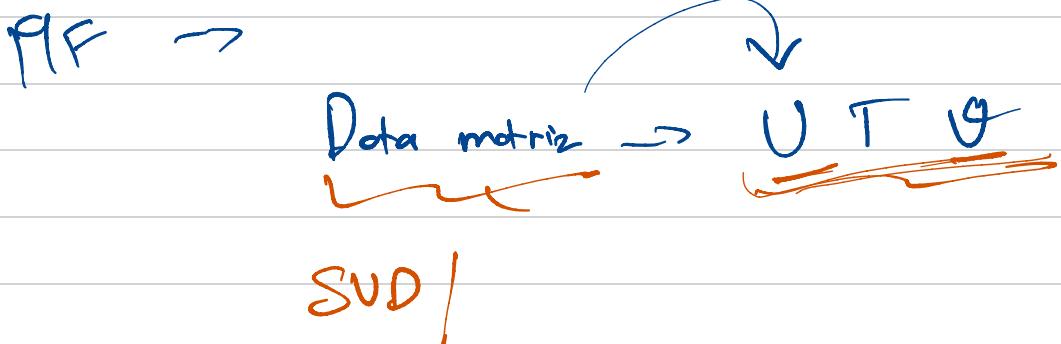
1. Matrix Factorization
2. Word2Vec model
3. CBOW, skip-gram v1/v2

Problem statement



Task:

$$w_i \rightarrow v_i \in \mathbb{R}^d \text{ using my corpus}$$



Bow → It were → Santa

→ Doesn't care about content / sequence
→ Mr. Spock

→ Words Co-occurrence matrix

--- . . . $w_1, w_2, w_3, w_4, w_5, w_6, w_7, w_8, \dots$

Window-size-1

$N \times N$ matrix; $N \rightarrow$ Vocab of my corpus

$w_1, w_2, w_3, w_4, w_5, w_6, \dots$

w_3	w_4	w_1	w_1	w_1	w_1

Today I got to buy a GPU

Yesterday I got fired
content-window → 2

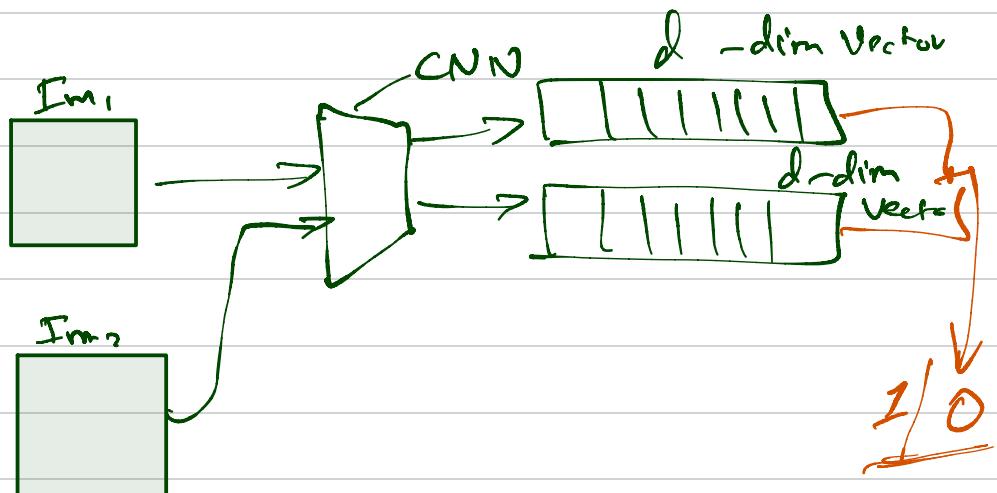
Today I got to

Today +1 +1
I M +1 +1

0 - occurrence → Pros → Partially Preserves
→ Doesn't ^{content} preserve
Sequence

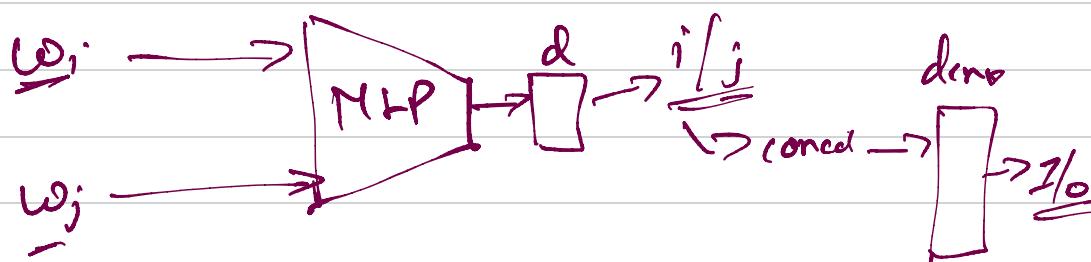
Word Vec

Siamese network



- ① Input
- ② Output
- ③ Loss
- ④ Architecture
- ⑤ Loss - differentiable

$$w_i \rightarrow \text{OHE} \rightarrow \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$



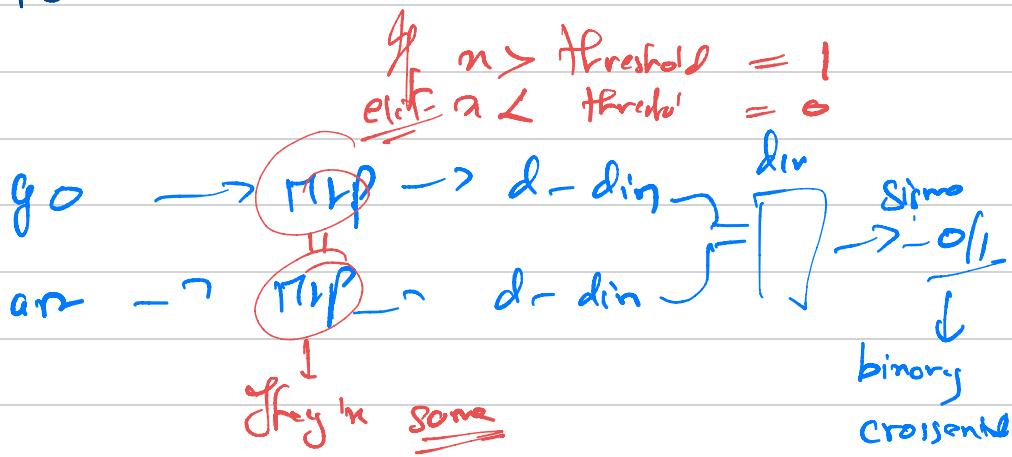
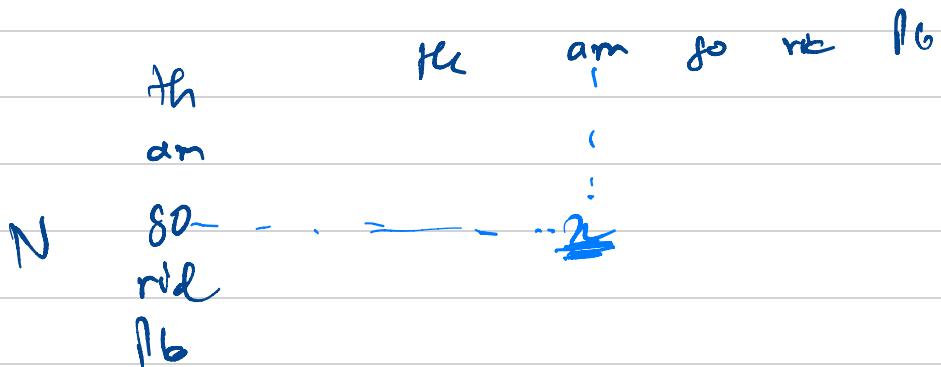
Vocab \rightarrow [was, the, some, writer, names]

[0, 0, 1, 0, 0]

\rightarrow In slopes \rightarrow we measure similarity

\rightarrow What we will do is train the model to learn whether w_i is in contrast of w_j or not.

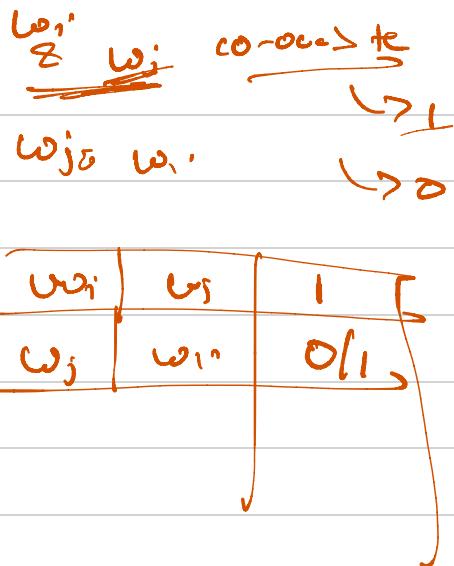
$N \rightarrow$ words \rightarrow entire corp



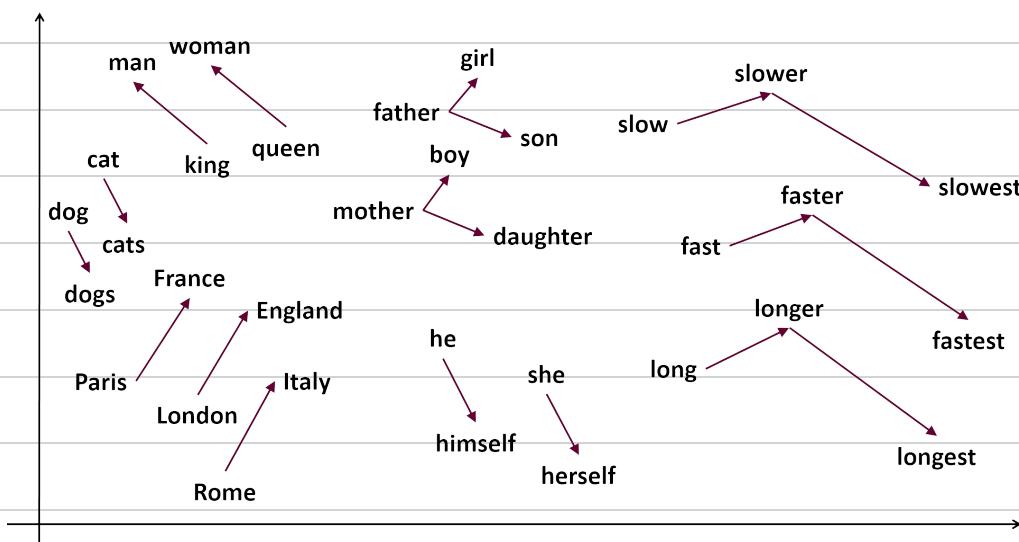
\Rightarrow When training the embedding \rightarrow to learn
 the contrast in which the word
is present

only 1 MLP (Neural network)



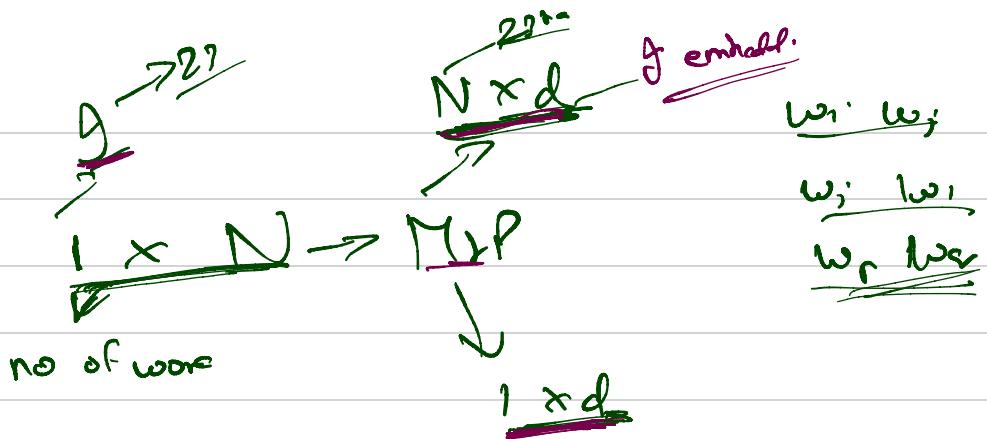


Word2vec in 2d representation



man
King
woman
Queen

(King - man) + Woman
 $\underbrace{\qquad\qquad}_{\text{Royalty}}$ = Queen



$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & x & 2 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}_{3 \times 3} = \boxed{\begin{bmatrix} 1 & 7 & 8 & 9 \\ \downarrow \end{bmatrix}}$$

~~we~~
 $N = 2$

all the words in
while we

is content with

hit losing
stim:

$$d = 100$$

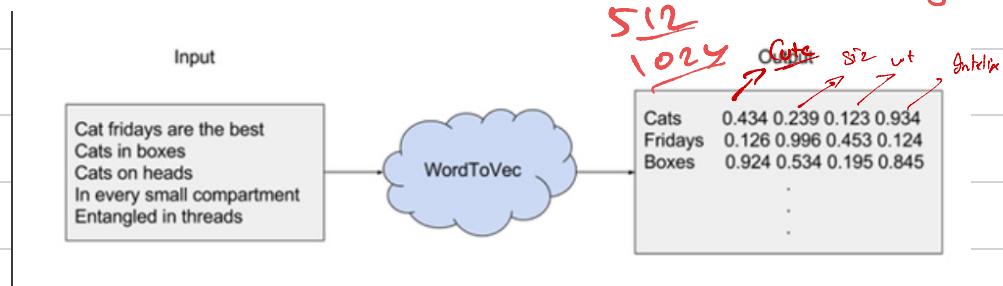
$$512$$

$$1024$$

Cats size ut

Intellix

	Cats	Fridays	Boxes	Cute	size	ut	Intellix
	0.434	0.239	0.123	0.934			
	0.126	0.996	0.453	0.124			
	0.924	0.534	0.195	0.845			
			



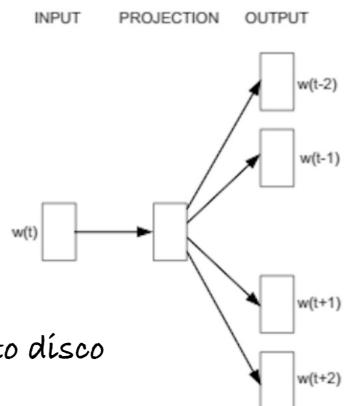
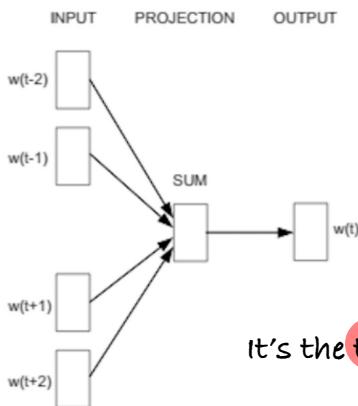
Word2Vec \rightarrow 2013

Glove \rightarrow 2014 \leftarrow (co-occurrence matrix)

FastText \rightarrow 2016-17 \uparrow NLP
 \searrow Phoenetic

Skip-gram with noise sampling

Voldemort



It's the time to disco

$w_1 \cdots w_d$ - - $w_{i,j} \cdots w_d$

loop

}

$S_1 \rightarrow 10$ word \rightarrow Entire Sentence
 $\hookrightarrow \underline{[10 \times 100]}$

$S_2 \rightarrow 20$ word \rightarrow (20×100)

① Average / centroid

② Average with PDF $\rightarrow S_{1,2}$
rare

multiply embedding of word
with rarity factor
before averaging

$S = R_c$

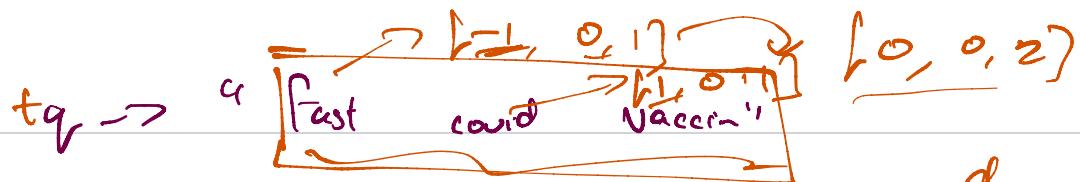
df-covm

$O \rightarrow$ Sum of all embed
to word oh

$1 \rightarrow \dots \dots a_2$

$\begin{matrix} L & \rightarrow \\ P & \rightarrow \\ 1 & , \end{matrix}$

a_1
 a_2
 c



① Sum + embodi. \rightarrow term d

② Fast $\rightarrow d$ with all abstract

covid $\rightarrow d$ / los \rightarrow with all abstract

Vaccin' $\rightarrow d$ / los \rightarrow with all d

Figure it out

Fast $\rightarrow a_1 \xrightarrow{0.3} a_2 \xrightarrow{0.2} a_3 \xrightarrow{0.4}$

covid $- a_1 \xrightarrow{} a_2 \xrightarrow{} a_3$

$$\begin{array}{ccc}
 \xrightarrow{0.2} & \xrightarrow{0.1} & \xrightarrow{0.1} \\
 \overline{0.4} & \overline{0.15} & \overline{0.25}
 \end{array}$$

$a_1 > a_3 > a_2$