

Sequence Models

3/5/25

Week One 1/12 Why Sequence Models?

Sequence data/problems

- speech
- sentiment classification
- video activity recognition
- music
- DNA sequences
- name entity recognition
- text
- translation

Week One 2/12 Notation

c.x. Name entity recognition

x: Harry Potter and Hermione Granger went on a walk.
 $x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$ $x^{(5)}$ $x^{(6)}$ $x^{(7)}$ $x^{(8)}$ $x^{(9)}$
 y: $y^{(1)}$ $y^{(2)}$ $y^{(3)}$ $y^{(4)}$ $y^{(5)}$ $y^{(6)}$ $y^{(7)}$ $y^{(8)}$ $y^{(9)}$

$T_x = 9$, sequence length
 $T_y = 9$

$x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$ $x^{(5)}$ $x^{(6)}$ $x^{(7)}$ $x^{(8)}$ $x^{(9)}$
 $y^{(1)}$ $y^{(2)}$ $y^{(3)}$ $y^{(4)}$ $y^{(5)}$ $y^{(6)}$ $y^{(7)}$ $y^{(8)}$ $y^{(9)}$

Representing words in Natural Language Processing (NLP)

Vocabulary

aaron	2
and	367
harry	4075
potter	6830
zulu	10,000

<unk>

one-hot

used to represent words not in your vocabulary

x: Harry Potter and Hermione Granger went on a walk

$x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$ $x^{(5)}$ $x^{(6)}$ $x^{(7)}$ $x^{(8)}$ $x^{(9)}$
 \downarrow \downarrow \downarrow \downarrow \downarrow \downarrow \downarrow \downarrow \downarrow
 $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

Week One 3/12 Recurrent Neural Network Model

T = 20

Why not a standard network? 2 main problems.

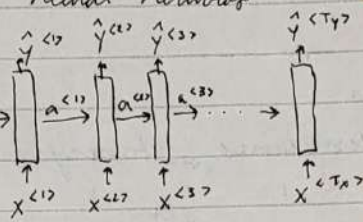
- 1) inputs, outputs can have different lengths
- 2) doesn't share features learned across different positions of text

Recurrent Neural Networks

where
 $T_x = T_y$

forward
prop

typically
vector of
zeros



W_{aa} W_{ax} W_{ya}

same
params used
each time step

limitation: uses info from
earlier in input
but not later

Solution: bi-directional

Forward propagation:

$a^{<0>} = \text{vector of zeros}$

$$a^{<1>} = g(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a)$$

$$y^{<1>} = g(W_{ya} a^{<1>} + b_y)$$

$\leftarrow g$: often tanh
occasionally ReLU

$\leftarrow g'$: sigmoid, softmax,
depends on what
task is

$$a^{<t>} = g(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a)$$

$$y^{<t>} = g(W_{ya} a^{<t>} + b_y)$$

Simplified notation

$$a^{<t>} = g(W_a [a^{<t-1>}, x^{<t>}] + b_a)$$

$$W_a = [W_{aa} \parallel W_{ax}]$$

\leftarrow stack matrices W_{aa} and W_{ax} horizontally
to create W_a matrix

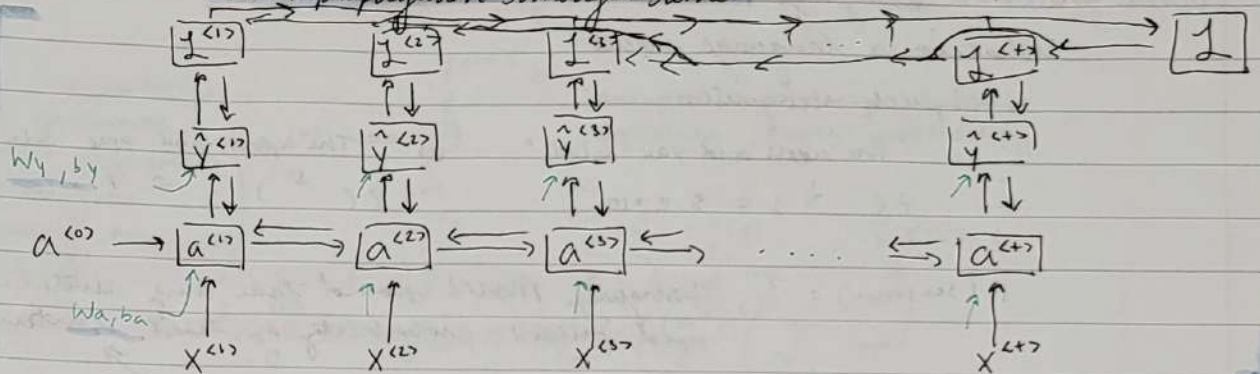
$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$$

\leftarrow stack vectors on top of each other

$$y^{<t>} = g(W_y a^{<t>} + b_y)$$

Week One 4/12

Backpropagation through time

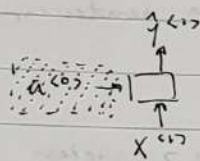


$$J^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1 - y^{(t)}) \log (1 - \hat{y}^{(t)})$$

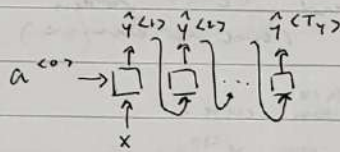
$$J(\hat{y}, y) = \sum_{t=1}^T J^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

Week One 5/12

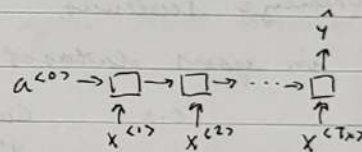
Different types of RNNs



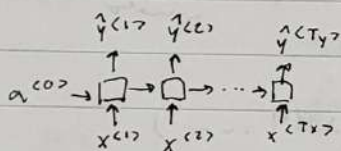
one to one



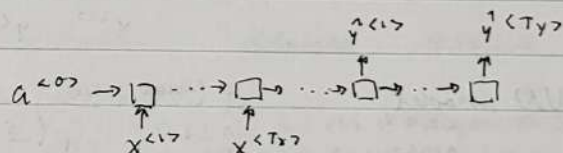
one to many
(sequence generation?)



many to one



many to many (where $T_x = T_y$)



many to many (when $T_x \neq T_y$)

Week One 6/12 Language model and sequence generation

T: 53

What is a language model?

Speech recognition

"The apple and pear salad."

vs "The apple and pear salad"

$$P(\downarrow) = 3.2 \times 10^{-15}$$

$$P(\downarrow) = 5.7 \times 10^{-10}$$

$P(\text{sentence}) = ?$, Language Model should take any sentence and output probability of that sentence

$P(y^{(1)}, y^{(2)}, \dots, y^{(T)})$ i.e. sequence of words

Language model w/ an RNN

Training set: large corpus of English text

Tokenize sentences, append $\langle \text{EOS} \rangle$ token, for sentences in your dataset (end of sentence)

e.g. Cats are cute.

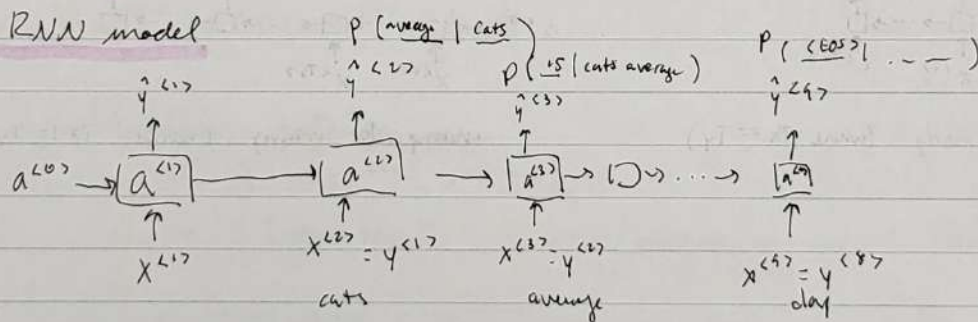
$y^{(1)} \quad y^{(2)} \quad y^{(3)}$

or
Cats are cute. $\langle \text{EOS} \rangle$

$y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad y^{(4)}$

$\langle \text{UNK} \rangle$ = token for unknown or uncommon words

RNN model



"Cats average 15 hours of sleep a day. $\langle \text{EOS} \rangle$ "

Cost function

$$Z = (y^{(1)}, y^{(2)}, \dots, y^{(T)}) = - \sum_t y_i^{(t)} \log \hat{y}_i^{(t)}$$

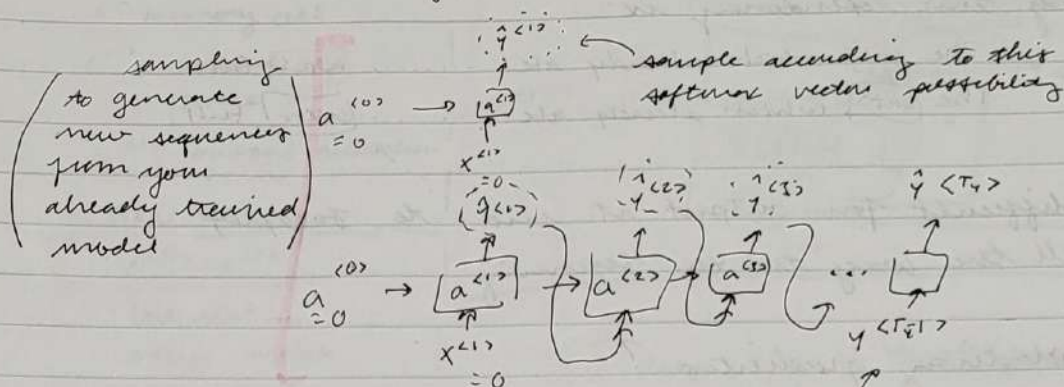
$$Z = \sum_t Z^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

softmax
loss

Week One 7/12 Sampling Novel Sequences

(in a RNN)

3/23/25



* model could output/generate unknown word token.

2 options

- 1) reject all output samples w/ unknown word token
- 2) leave it in if you don't mind

could be <EOS> token as until n amt of tokens

So far we have built a 'word level model' but you can also build 'character level model'

tokens are chars, think ASCII, instead of words

e.x vocabulary

a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x	y	z
A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
!	"	#	\$	%	&	'	()	*	+	=	<	>	[]	{	}	~	?	;	:	~	~	~	~

Pros + Cons

⊕ no unknown word tokens

⊖ much longer sequences

↓
char level models not as good at capturing how early parts of sentences affect latter parts

⊖ more computationally expensive

Week One 8/12 Vanishing Gradients of RNNs

T = 6.9

long term dependency is:

The cat which already ate was full

The cats which already ate were full

• difficult for output at end to backprop
all the way to the beginning

Exploding gradients?

↓
gradient clipping.

- rescale some gradient vectors
to be 'clipped' according to
some max value

Vanishing gradients?

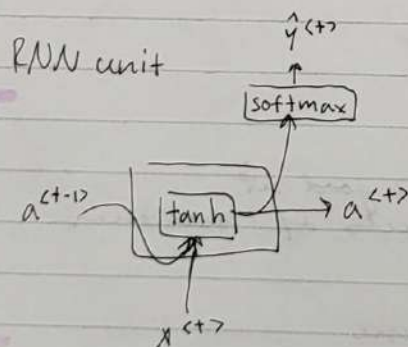
↓
harder to deal w/ :-

Week One 9/12 Gated Recurrent Unit (GRU)

T = 8.5

- a mod to RNN hidden layer

- helps capture long-range connections +
deal w/ vanishing gradient issue



tanh

$$a^{(t)} = g(W_a [a^{(t-1)}, x^{(t)}] + b_a)$$

GRU unit (simplified)

c = memory cell

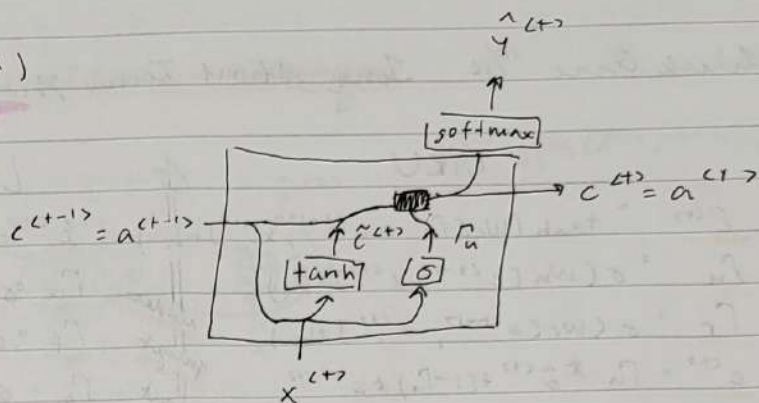
$c^{(t)}$ = value of c @ time t

$a^{(t)}$ = output activation value

in GRU

$c^{(t)} = a^{(t)}$

but not in other archts.



$$\tilde{c}^{(t)} = \tanh(W_c [c^{(t-1)}, x^{(t)}] + b_c)$$

$$\Gamma_u = \sigma(W_u [c^{(t-1)}, x^{(t)}] + b_u)$$

Γ_u = gate Γ = Gamma "update"

\tilde{c} = candidate to replace $c^{(t)}$

$$c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$$

$c^{(t)}$, $\tilde{c}^{(t)}$, & Γ_u can be vectors these are element wise multiplication in that case

Full GRU

$$\tilde{a}^{(t)} = \tanh(W_c [\Gamma_r * c^{(t-1)}, x^{(t)}] + b_c)$$

$$\Gamma_u = \sigma(W_u [c^{(t-1)}, x^{(t)}] + b_u)$$

$$\Gamma_r = \sigma(W_r [c^{(t-1)}, x^{(t)}] + b_r)$$

$$c^{(t)} = \Gamma_u * \tilde{a}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$$

alt. notation

Week One 10/12 Long Short Term Memory (LSTM) unit T=94

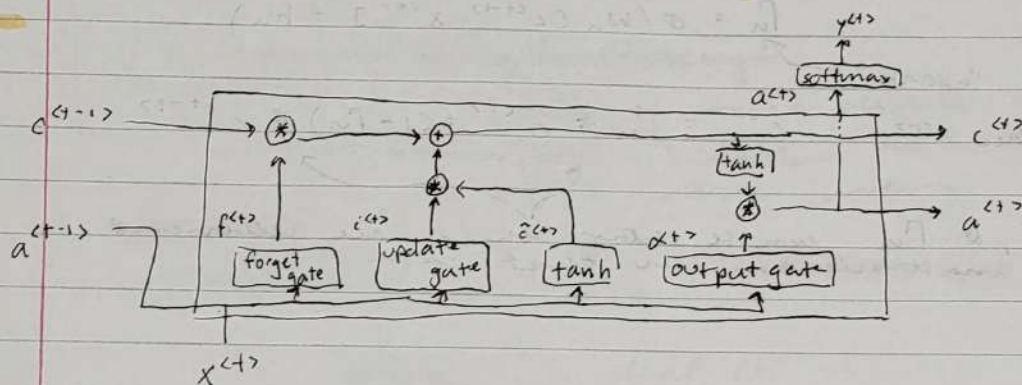
GRU

$$\begin{aligned}\tilde{c}^{(t)} &= \tanh(W_c [\Gamma_r * c^{(t-1)}, x^{(t)}] + b_c) \\ \Gamma_u &= \sigma(W_u [c^{(t-1)}, x^{(t)}] + b_u) \\ \Gamma_r &= \sigma(W_r [c^{(t-1)}, x^{(t)}] + b_r) \\ c^{(t)} &= \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)} \\ a^{(t)} &= c^{(t)}\end{aligned}$$

update
forget
output

LSTM

$$\begin{aligned}\tilde{c}^{(t)} &= \tanh(W_c [a^{(t-1)}, x^{(t)}] + b_c) \\ \Gamma_u &= \sigma(W_u [a^{(t-1)}, x^{(t)}] + b_u) \\ \Gamma_f &= \sigma(W_f [a^{(t-1)}, x^{(t)}] + b_f) \\ \Gamma_o &= \sigma(W_o [a^{(t-1)}, x^{(t)}] + b_o) \\ c^{(t)} &= \Gamma_u * \tilde{c}^{(t)} + \Gamma_f * c^{(t-1)} \\ a^{(t)} &= \Gamma_o * \tanh(c^{(t)})\end{aligned}$$



"peephole connection"

where in $\Gamma_u, \Gamma_r, \Gamma_o$, $W[a^{(t-1)}, x^{(t)}, c^{(t-1)}]$,
 $c^{(t-1)}$ used to affect gate value ↑

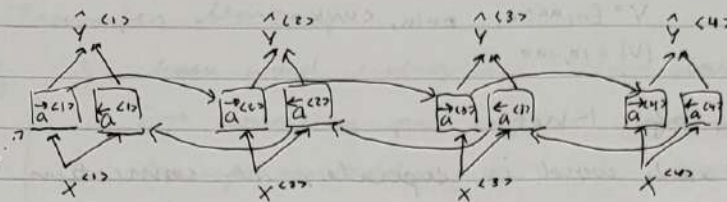
GRU v. LSTM

- neither is universally superior
- GRU is lighter, easier to scale
- LSTM more powerful, flexible

Week One 1/2 Bidirectional RNN (BRNN)

T=106

3/24/25



$$\hat{y}^{(t)} = g(w_y [a^{(t)}, \bar{a}^{(t)}] + b_y)$$

RNN blocks

G, U, V

(STM blocks) ← commonly used in BRNNs

note: need entire sequence before BRNN can use it, i.e. not suited for rolling speech transcription

Week One 1/2 Deep RNNs

$a^{[L] \langle t \rangle}$

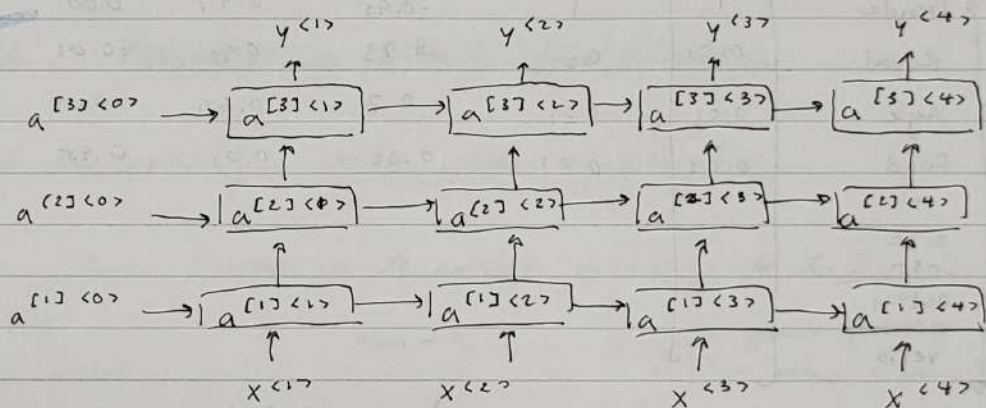
L-layer

t-time, i.e.

$a^{[2] \langle t \rangle}$

layer 2

time t



$$a^{[2] \langle t \rangle} = g(W_a^{[2]} (a^{[3] \langle t \rangle}, a^{[1] \langle t \rangle}) + b_a^{[2]})$$

- 3 layers is kinda already deep for RNN
- computationally expensive
- could connect y 's to standard neural network layers (not horizontally connected)

Intro to Word Embeddings

Week Two 1/4

Word Representation

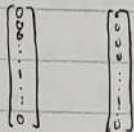
T = 12.1

1-hot representation

$V = [a, aaron, \dots, zulu, <UNK>]$

$|V| = 10,000$

Man (5391) Woman (9853)



O_{5391}

O_{9853}

Cons of 1-hot

each word is separate, no connection between them

Featurized representation

features

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	-0.05	0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
size						
cost						
action						
noun						
verb						

e_{5391}

e_{9853}

e_{4914}

t-SNE - visualizing word embeddings

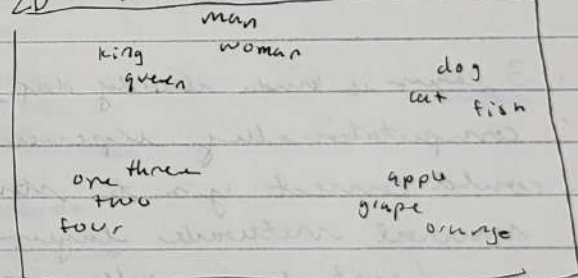
3D



300D orange vector

300D \rightarrow 2D

2D t-SNE



Week Two ^{2/4} Using Word Embeddings

T=130

Transfer Learning

- (A) 1. learn word embeddings from large text corpora (1-100B words)
(or download pre-trained embedding)
- (B) 2. transfer embedding to new task of smaller training set (e.g. 100k)
3. optional: fine-tune embeddings of new data
[only in practice when training set is large]

most helpful when tens of data for (A) and less for (B)

Week Two ^{3/4} Properties of Word Embeddings

T=141

word embeddings help with analogies

e.g. "man is to woman as king is to queen"

model can calculate \rightarrow
then look for a word
that creates same result
for $e_{\text{king}} - \boxed{?}$

$$e_{\text{man}} - e_{\text{woman}} \approx \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \left\{ \begin{array}{l} \leftarrow \text{gender feature} \\ \leftarrow \text{other features} \end{array} \right.$$

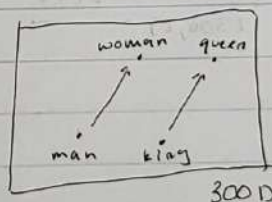
$$e_{\text{king}} - e_{\text{queen}} \approx \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \quad \leftarrow \text{same}$$

$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{?}$$

find word w that maximizes

similarity of $e_{\text{king}} - e_w$ & $e_{\text{man}} - e_{\text{woman}}$

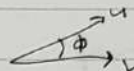
$$\arg \max \text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$$



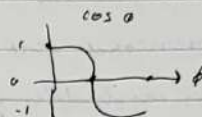
visual / parallelism
relationships likely
won't hold thru
t-SNE

Cosine similarity

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$



Euclidean distance /
square distance
 $\|u - v\|^2$

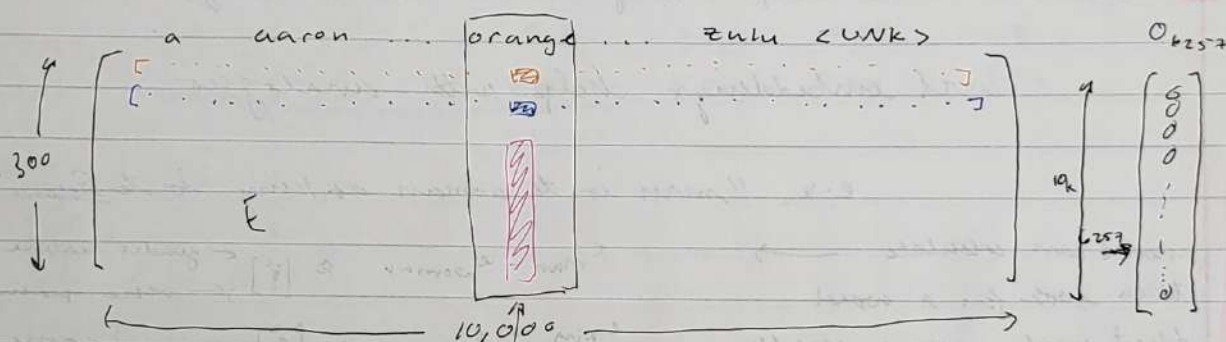


Week Two 4/4 Embedding Matrix

T=145

~ 2.4 h

When you implement an algo to learn word embedding,
what you learn is an embedding matrix



$$E \cdot \mathbf{e}_{6257} = \begin{bmatrix} \text{orange} \end{bmatrix} = \mathbf{e}_{6257}$$

(300, 10k) (10k, 1) (300, 1)