

# SEQUENCE MODELS

PAGE NO.:  
DATE: / /

## # Why Sequence Models

## Examples of seq data

# Speech recognition

## music generation

## Sentiment classification

## DNA seq analysis

## \* Notation

## example

\* Harry Potter and Hermione Granger invented a potion

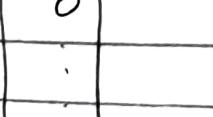
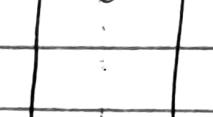
$y:$     1    0    1    0    1    0    0    0  
            $\langle 1 \rangle$      $\langle 2 \rangle$      $\langle 3 \rangle$           $\langle 4 \rangle$      $\langle 5 \rangle$      $\langle 6 \rangle$      $\langle 7 \rangle$      $\langle 8 \rangle$   
 indexing:  $x$      $x$      $x$      $x$     .    .    .     $x$      $x$      $x$      $x$   
            $\langle 1 \rangle$      $\langle 2 \rangle$      $\langle 3 \rangle$           $\langle 4 \rangle$      $\langle 5 \rangle$      $\langle 6 \rangle$      $\langle 7 \rangle$      $\langle 8 \rangle$   
           y    y    -    -    -    y    -    y  
           y    y    -    -    -    y    -    y  
 $T_f = 9$

$x^{(i)j}$  :  $j^{\text{th}}$  element of  $i^{\text{th}}$  training example

$T_x^{(i)}$ : length of  $i^{\text{th}}$  training example

Vocabulary: [a, aaron, ... harry ... potter ... zulu]  
all the possible word words. 10,000

are represented as one-hot form

$x =$ 	
$\begin{matrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} \rightarrow \begin{matrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix}$ Harry	$\begin{matrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix}$ Potter

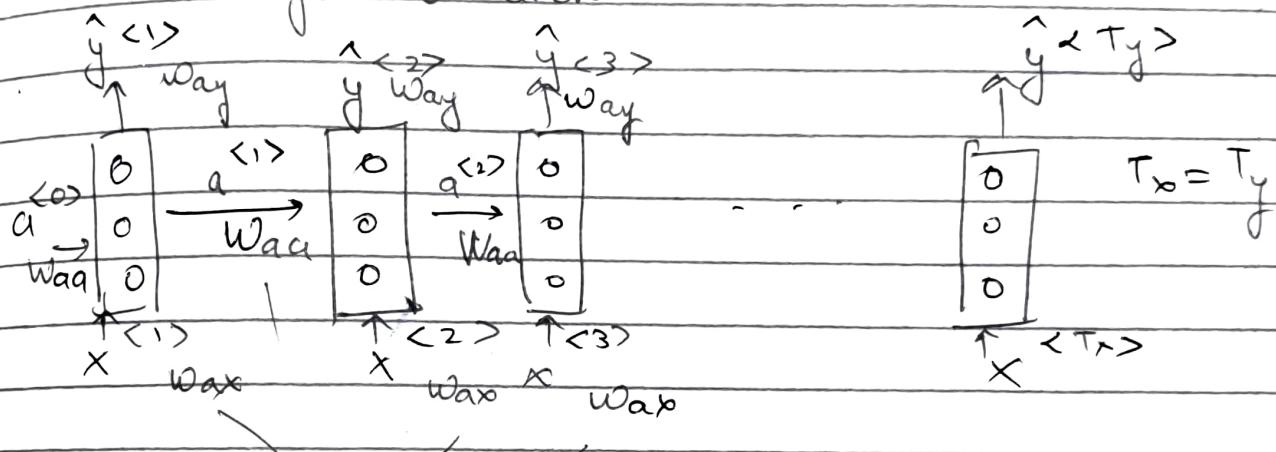
## # Recurrent Neural Network Model

Problems in a std network

- Input, outputs can be of diff lengths in diff examples
- doesn't share features learned across diff positions of text

RNN

general arch:



Same parameters

forward propagation

$$a^{<0>} = \vec{0}$$

$$a^{<1>} = g(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a) \leftarrow \tanh/\text{ReLU}$$

$$\hat{y}^{<1>} = g(W_{ya} a^{<1>} + b_y) \leftarrow \text{Sigmoid}$$

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

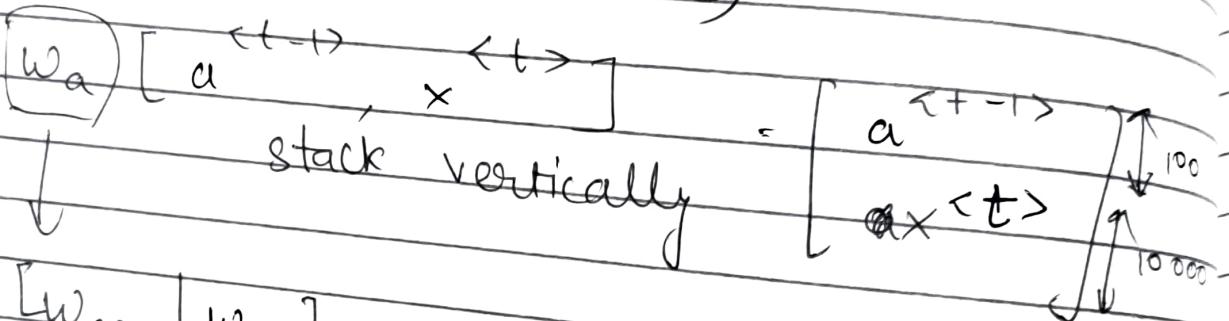
$$\hat{y}^{<t>} = g(W_{ya} a^{<t>} + b_y)$$

$$\mathbf{a}^{<t>} = g(\mathbf{w}_a [\mathbf{a}^{<t-1>}, \mathbf{x}^{<t>}] + \mathbf{b}_a)$$

$\mathbf{a}^{<t>} \quad \text{Shape}$

$$\mathbf{w}_{aa} = (1, 00, 100)$$

$$\mathbf{w}_{ax} = (1, 00, 10000)$$



$$[\mathbf{w}_{aa} \mid \mathbf{w}_{ax}]$$

stack horizontally

$$[\mathbf{w}_{aa} \mid \mathbf{w}_{ax}] \begin{bmatrix} \mathbf{a}^{<t-1>} \\ \mathbf{x}^{<t>} \end{bmatrix} = \mathbf{w}_{aa} \mathbf{a}^{<t-1>} + \mathbf{w}_{ax} \mathbf{x}^{<t>}$$

\* Back prop

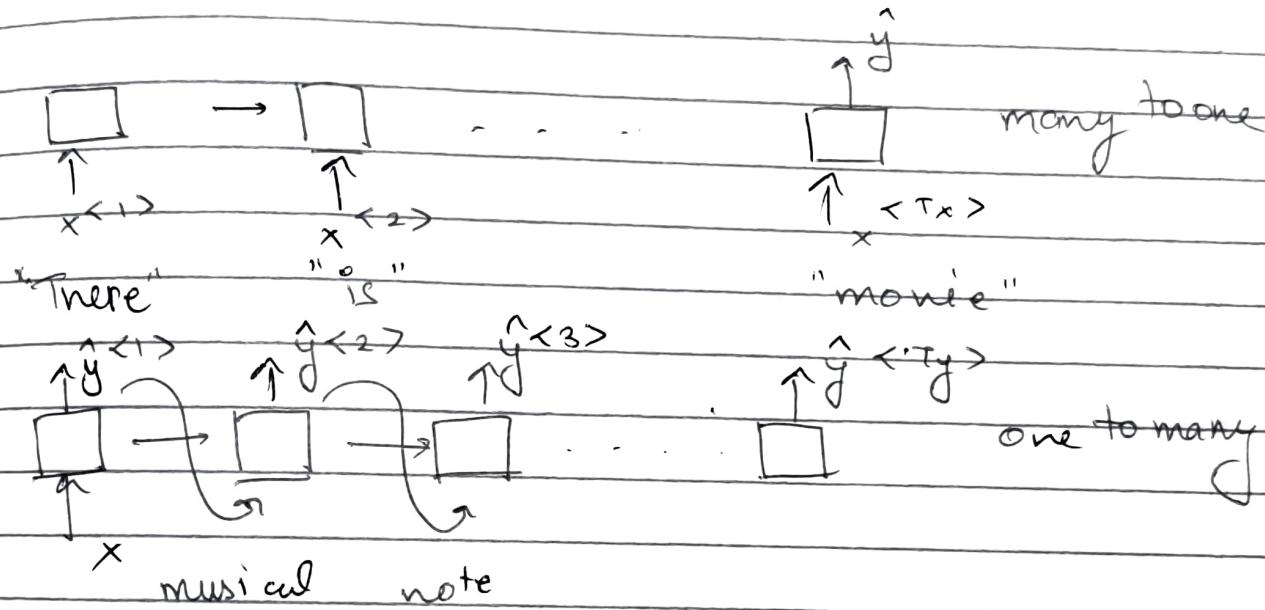
$$J^{<t>} \circ (\hat{\mathbf{y}}^{<t>} - \mathbf{y}^{<t>})$$

$$J(\hat{\mathbf{y}}, \mathbf{y}) = \sum_t J^{<t>} (\hat{\mathbf{y}}^{<t>} - \mathbf{y}^{<t>})$$

## \* RNN Types

$T_x$  does not always equals to  $T_y$   
eg: sentiment classification:

- "There is nothing to like in the movie"  
converted to a star rating.
- generating music from initial notes.



## \* Language Model & Seq generation

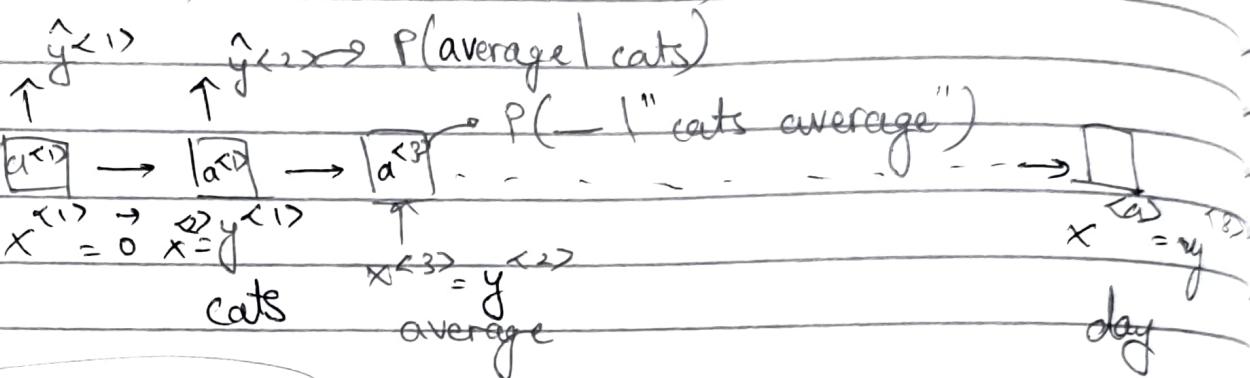
Training set: large corpus of english text

Tokenize

Cats average 15 hours of sleep a day. <EOS>  
 are hot tokens  
 y <sup><1></sup> y <sup><2></sup> . . . ~

The Egyptian Man is a breed of cat!  
 is not in the vocabulary  
 replace with token (UNK)

$P(a) P(aaron) \dots P(\text{cats}) \dots P(\text{zulu}) P(\text{UNK}) P(\text{EOS})$



$\rightarrow (\text{Cats average } 15)$  hours of sleep a day.  $\text{EOS}$

$$L(\hat{y}^{(t)}, y^{(t)}) = -\sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

\* Sampling Novel Sequences  
Word  
character

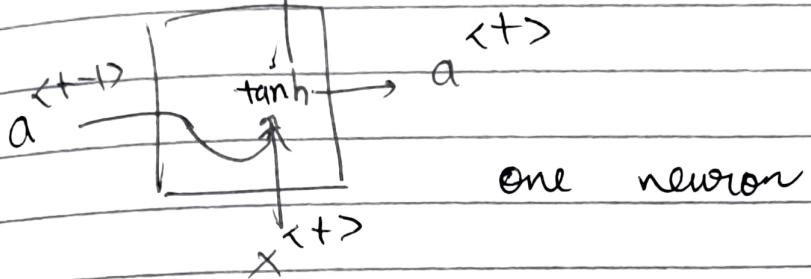
\* Vanishing Gradients in RNNs

In case of deep RNNs the backpropagation gradients obtained from the later layers may not significantly affect the earlier layers.

## \*GRU - Gated Recurrent Unit

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

$\xrightarrow{\text{softmax}}$   $\hat{y}^{(t)}$



The cat, which already ate..., was full  
 The cat, which already ate..., were full

Let's introduce  $c$  : memory cell  
 which will store cat/cats

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

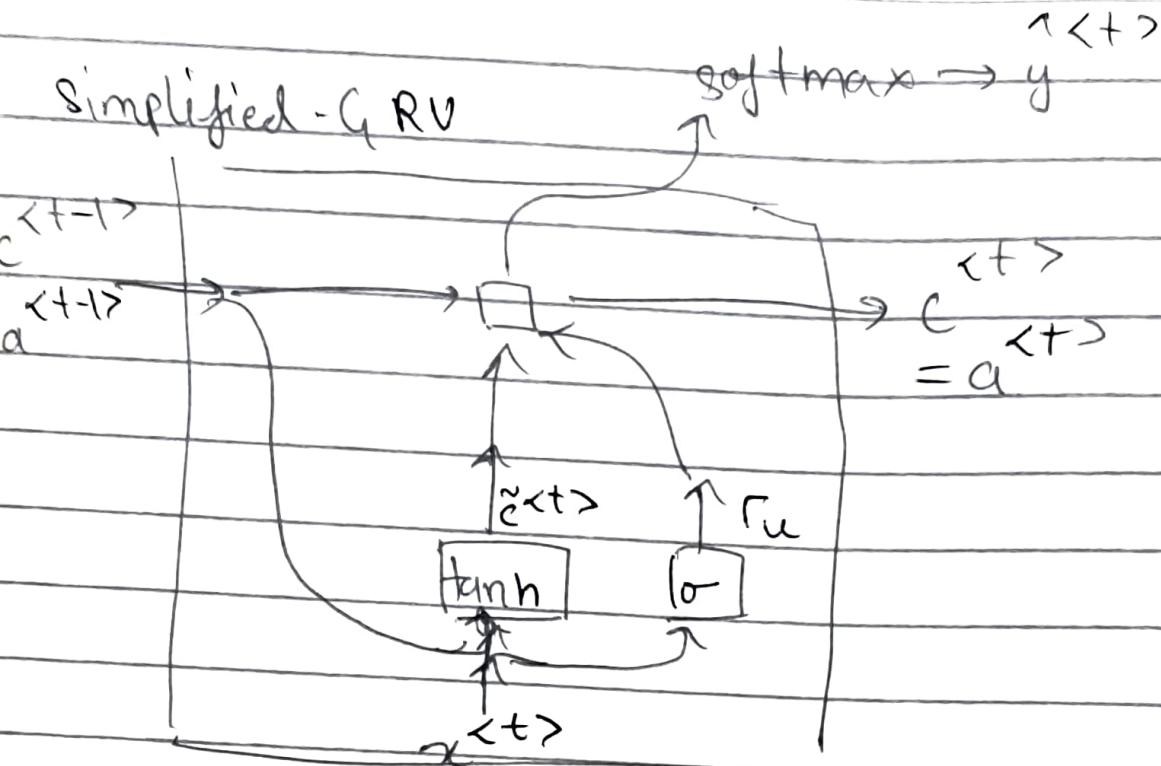
$$c^{(t)} = \tanh(W_c [c^{(t-1)}, x^{(t)}] + b_c)$$

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

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

element wise  
 multiplication

simplified - GRU



## \* Long Short Term Memory (LSTM)

$$c^{<t>} = \tanh(w_c [a^{<t-1>} \cdot x^{<t>}] + b_c)$$

$$r_u = \sigma(w_u [a^{<t-1>} \cdot x^{<t>}] + b_u) \quad \text{update gate}$$

$$f = \sigma(w_f [a^{<t-1>} \cdot x^{<t>}] + b_f) \quad \text{forget gate}$$

$$o = \sigma(w_o [a^{<t-1>} \cdot x^{<t>}] + b_o) \quad \text{output gate}$$

$$c^{<t>} = r_u \cdot c^{<t-1>} + f \cdot c^{<t-1>}$$

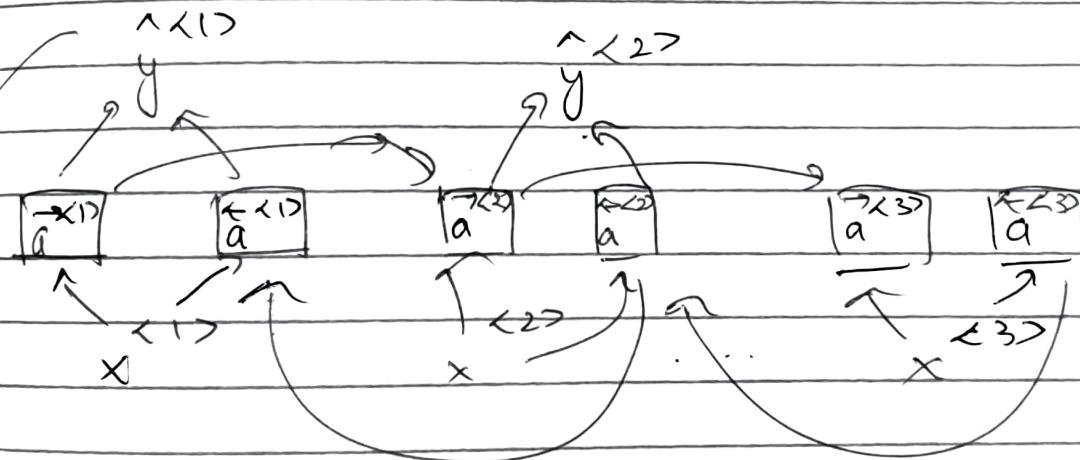
$$a^{<t>} = o \cdot \tanh(c^{<t>})$$

## \* Bidirectional RNN

He said " Teddy bears <sup>are</sup> ~~is~~ on sale!"  
 ) not person

He said " Teddy Roosevelt was a great president"  
 ) person

We need more info than the previous words (in both cases - same) to determine if Teddy is a person or not.  
 It will be obtained from the later part of the sentence.

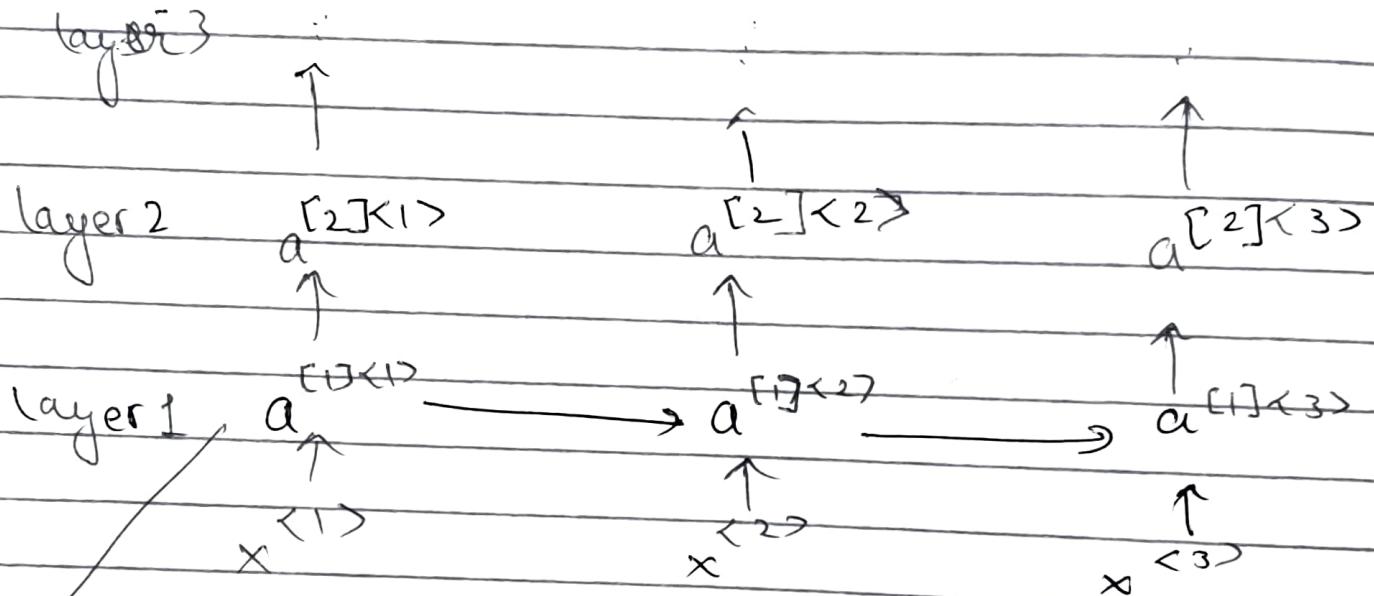


$$y^{(t)} = g(w_y [a^{(x_t)}, a] + b_y)$$

It is often used with LSTM for NLP

You need to process entire data ~~seq~~ before making any prediction

## \* Deep RNNs



can be  
RNN,  
GRU,  
LSTM

$$g^{[2]<3>} = g\left(\omega_a^{[2]} [a^{[2]<2>}, a^{[1]<3>}] + b_a^{[2]}\right)$$

## \* Word Embeddings

problem with one hot representation

$$\begin{array}{l} \text{'apple'} = \\ 456 \end{array} \quad \left[ \begin{array}{c} 0 \\ 0 \\ -1 \\ 0 \\ \vdots \\ 0 \end{array} \right] \quad \begin{array}{l} \text{'orange'} : \\ 6257 \end{array} \quad \left[ \begin{array}{c} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ -1 \\ 0 \end{array} \right]$$

I want a glass of apple juice

I want a glass of orange juice.

there is a lot large distance (less similarity) in apple & orange according to one hot representation.

## Featureized representations

	Man 5391	Woman 9853	King 49140	Queen 7157	Apple 456	Orange 6250	
Gender	1	1	-0.95	0.97	0.0	0.01	?
Royal	0.01	0.02	0.93	0.95	-0.01	0.00	
Age	0.03	0.02	0.7	0.69	0.03	-0.02	
Food	0.04	0.01	0.02	0.01	0.95	0.97	
Size							
Cost							

300 feature word embedding

embedding  $\approx$  encoding

PAGE NO.:  
DATE: / /

These 300 dimensions can be reduced to 2D using t-SNE

### \* Using Word Embeddings

Joe David is an apple farmer

Don Pablo is an orange farmer

↳ it is easy to learn that  
Don Pablo is a person

Michael Scott is a kiwi cultivator

↳ never seen

before words

What to do now?

Use transfer learning

1. Learn word embeddings from a larger text corpus (1-100B words)

or download pre trained embedding

2. Training Transfer embedding to new task with smaller training set. (100k words)

# Properties of word Embeddings:

## Analogies

Man : Woman :: King : ?

from the previous table

$$\mathbf{e}_{\text{man}} = \begin{bmatrix} -1 \\ 0.01 \\ 0.03 \\ 0.04 \end{bmatrix} \quad \mathbf{e}_{\text{woman}} = \begin{bmatrix} 1 \\ 0.02 \\ 0.02 \\ 0.01 \end{bmatrix}$$

$$\mathbf{e}_{\text{man}} - \mathbf{e}_{\text{woman}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

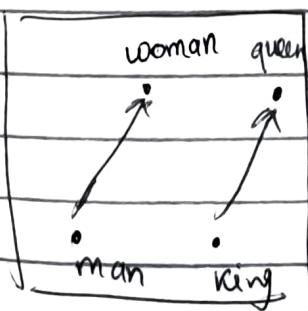
$$\mathbf{e}_{\text{King}} - \mathbf{e}_{\text{Queen}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

find  $\mathbf{e}_w$  such that

$$\mathbf{e}_{\text{man}} - \mathbf{e}_{\text{woman}} \approx \mathbf{e}_{\text{King}} - \mathbf{e}_w$$

$\downarrow$

$$\mathbf{e}_w = \mathbf{e}_{\text{Queen}}$$



find word  $w$ :

$$\arg \max_w \text{sim}(\mathbf{e}_w, \mathbf{e}_{\text{King}} - \mathbf{e}_{\text{Man} + \frac{\mathbf{e}_{\text{Woman}}}{\|\mathbf{e}_{\text{Woman}}\|}})$$

300D

Cosine similarity

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

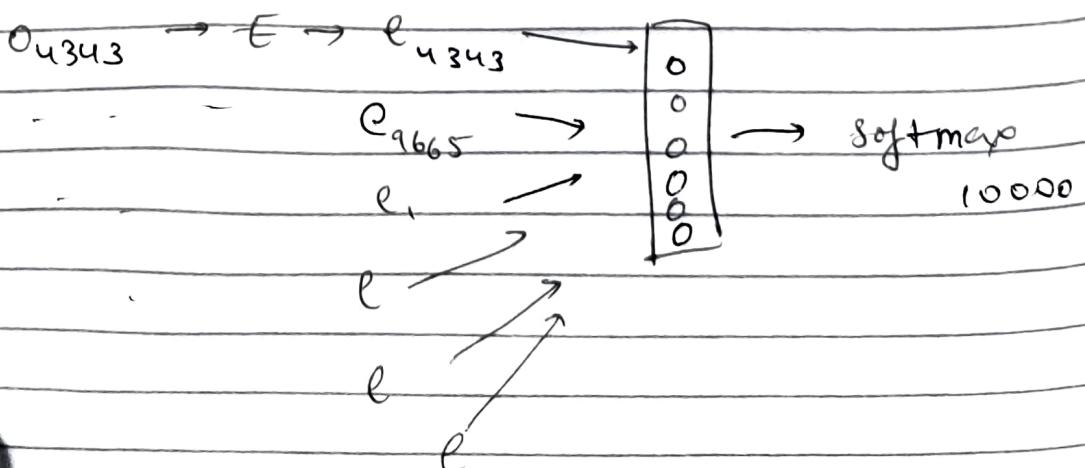
\* Embedding Matrix

$$E_{\text{mat}} = \begin{bmatrix} 6257 & 0_{6257} \\ \text{aaron} & \begin{bmatrix} \text{Orange} \\ \text{zulu} \\ \langle \text{UNK} \rangle \end{bmatrix} \\ 300 & 6257 \end{bmatrix} \xrightarrow{\text{10000}} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}_{10000}$$

$$\overset{e_{6257}}{(300, 10k)} \xrightarrow{\epsilon} \overset{e_{6257}}{(10k, 1)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} = e_{6257}$$

\* Neural Language Model

I want a glass of orange juice  
 4343 9665 1 3852 6163 6257



## \*Word2Vec

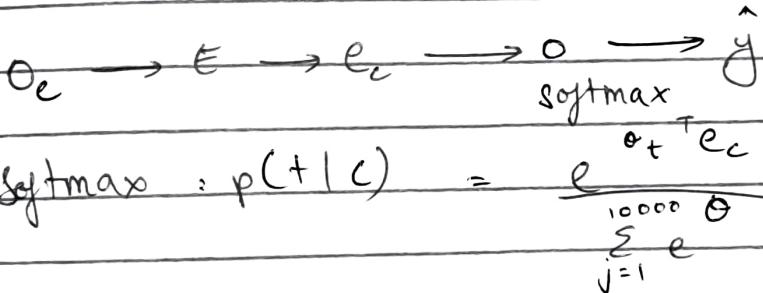
I want a glass of orange juice to go along with my cereal

<u>context</u>	<u>target</u>
orange	glass
orange	juice
orange	my

context c ("orange") → target t ("juice")

6257

4834



$\theta_t$ : parameter associated with t

$$L(\hat{y}, y) = -\sum_{i=1}^{10K} y_i \log \hat{y}_i$$

## \*Negative Sampling

<u>context</u>	<u>target</u>	<u>target</u>
orange	juice	1
orange	king	0
orange	of	0