

# LSTM (long Short term memory) Unit.

- first difference b/w GRU And LSTM.

$\Rightarrow$  GRU has 2 gates in it.

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

$$\Gamma_u = \sigma(W_u[\tilde{C}^{(t-1)}, x^{(t)}] + b_u) \quad \{ 2 \text{ gates in GRU}$$

$$\Gamma_f = \sigma(W_f[\tilde{C}^{(t-1)}, x^{(t)}] + b_f) \quad \{ \text{LSTM has 3 gates.}$$

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

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

$$\therefore a^{(t)} = C^{(t)} \text{ instead of } a^{(t)} = \Gamma_o * C^{(t)}$$

$\Rightarrow$  LSTM has 4 gates in it.

$$\text{② } \tilde{C}^{(t)} = \tanh(W_c[a^{(t-1)}, x^{(t)}] + b_c)$$

$$\Gamma_u = \sigma(W_u[a^{(t-1)}, x^{(t)}] + b_u) \rightarrow \text{update gate}$$

$$\Gamma_f = \sigma(W_f[a^{(t-1)}, x^{(t)}] + b_f) \rightarrow \text{forget gate}$$

$$\Gamma_o = \sigma(W_o[a^{(t-1)}, x^{(t)}] + b_o) \rightarrow \text{output gate}$$

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

$$a^{(t)} = \Gamma_o * \tanh C^{(t)}$$

- LSTM - the other type of RNN that can enable you to account for long-term dependencies. It's more powerful and general than GRU.
- In LSTM,  $c^{(t+1)} = a^{(t+1)}$
- In GRU we have update gate  $U(\Gamma_u)$ , a relevance gate  $r(\Gamma_r)$ , and a candidate cell variables  $\tilde{c}$ . While in LSTM we have an update gate  $U$  (Sometime its called input gate I), a forget gate  $F$ , and output gate  $O$ , and a Candidate cell variable  $\tilde{c}$ .
- Some Variants on LSTM includes:
  - LSTM with peephole connection.
  - The normal LSTM with  $c^{(t+1)}$  include with every gate.
- There isn't a universal Superior between LSTM and its variants. One of the advantages of GRU is that it's simpler and can be used to build much bigger networks but the LSTM is more powerful and general.

# Vanilla Long Short-Term Memory

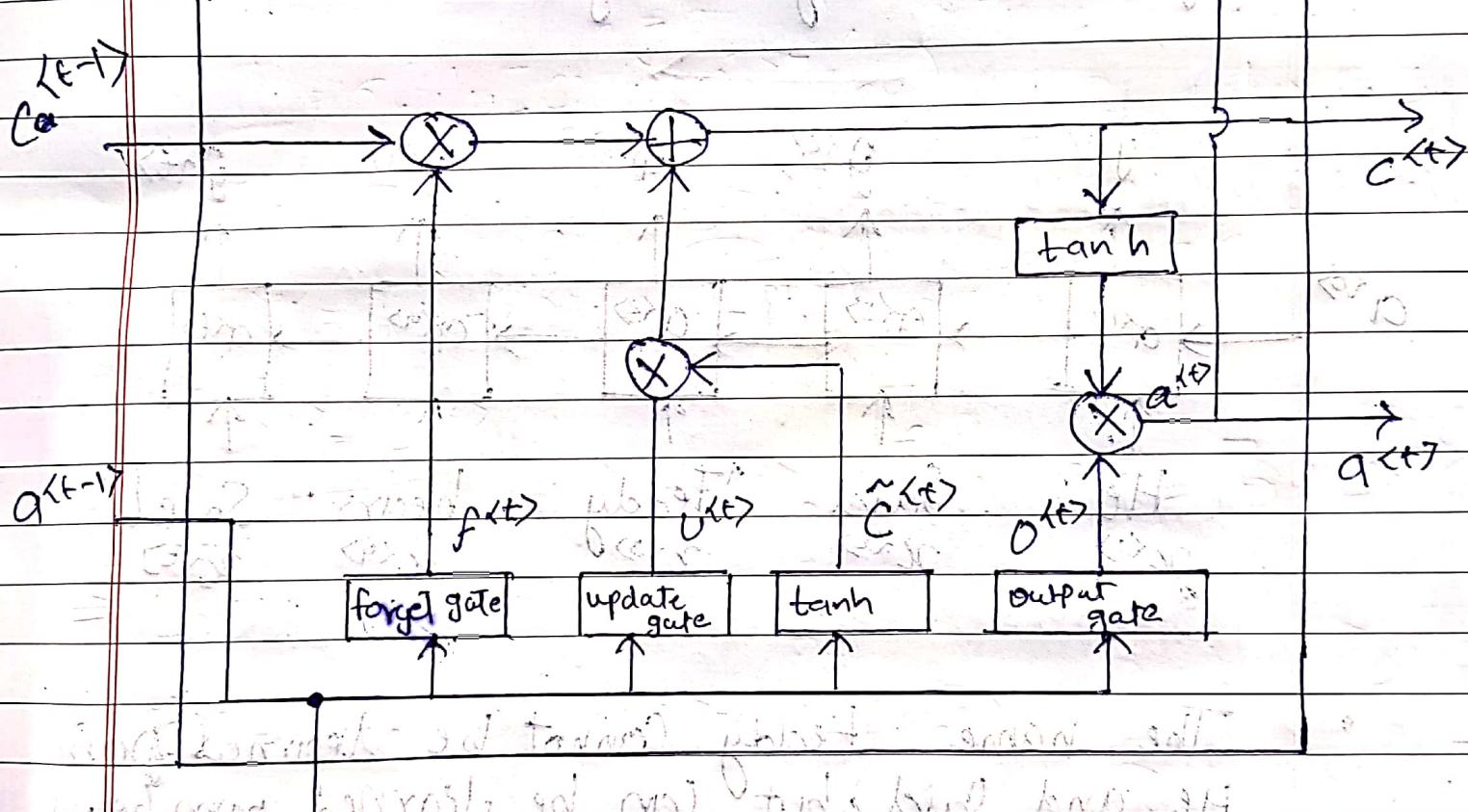
Int. at each step with own input

Previous hidden layer seen during update

Long short-term memory for each sequence  
just a certain time step  $\rightarrow$   $\text{softmax}$

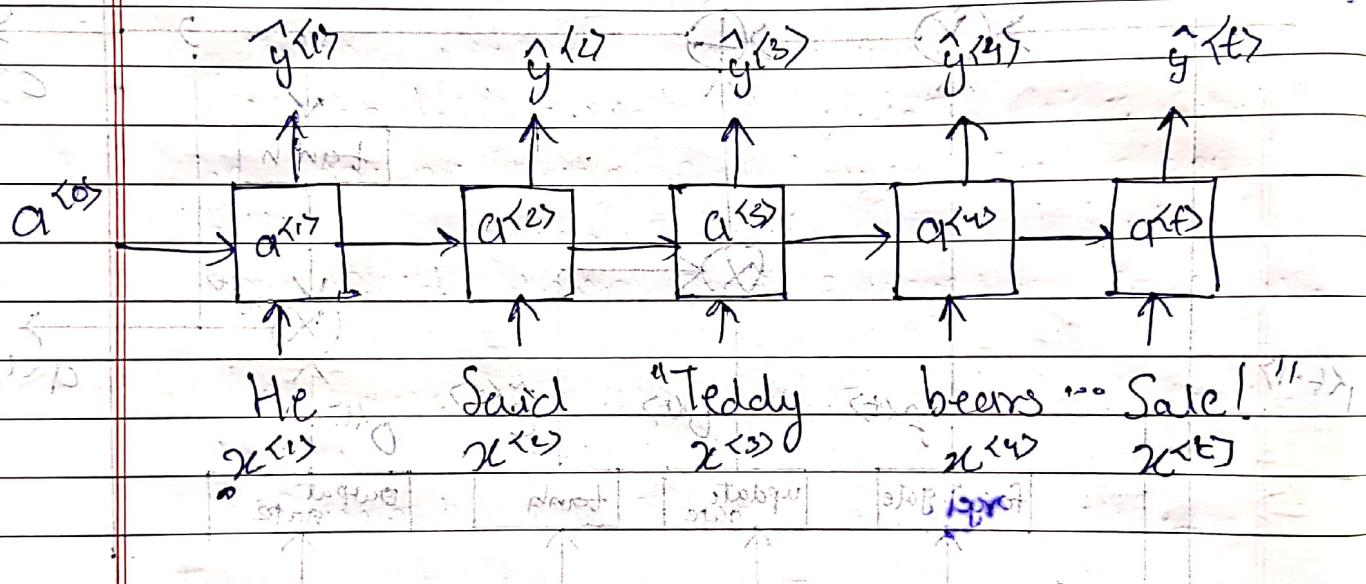
works like LSTM but without forget gate

most activation is hidden state



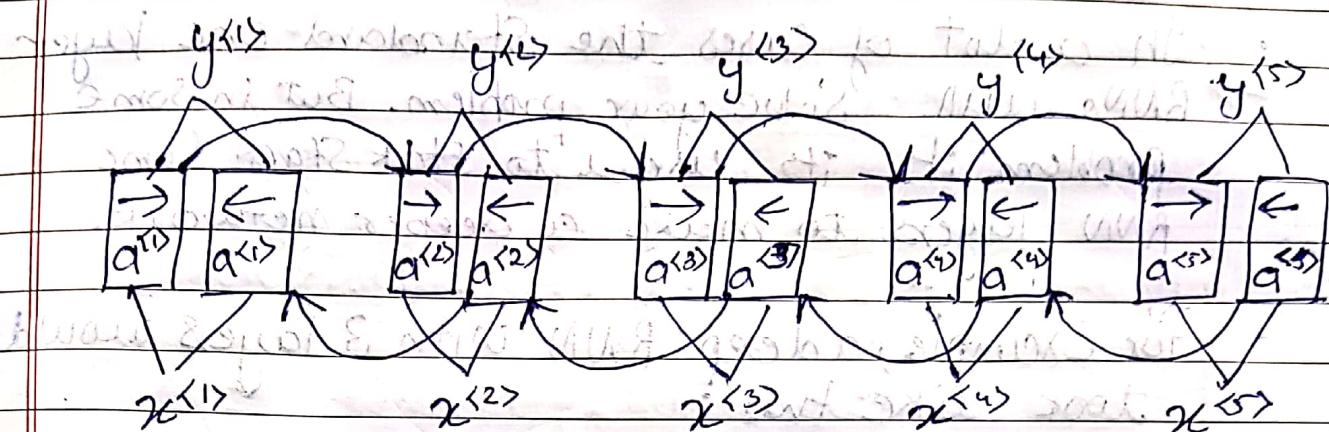
## \* Bidirectional RNN

- There are still some ideas to let you build much more powerful sequence models. One of them is bidirectional RNNs and another is Deep RNNs.
- As we saw before, here is an example of Name entity recognition task:



- The name teddy cannot be learned from He and Said, but can be learned from bears.
- BiRNN fixes this issue.

- Here is BRNN architecture

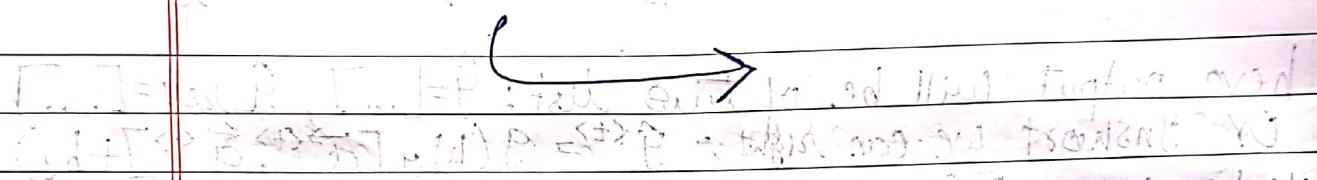


here output will be of two list:  $\hat{y} = [\dots]$ ,  $\hat{y}_{rev} = [\dots]$   
OR Inshort we can write:  $\hat{y}^{<t>} = g(W_y [\alpha^{<t>}, \alpha^{>t>}] + b_y)$

- Note, that BRNN is an acyclic graph
- Part of the forward propagation goes from left to right and part from right to left. It learns from both sides.
- To make prediction we use  $\hat{y}^{<t>}$  by using the two activations that come from left and right.
- The blocks here can be any RNN block including the basic RNN, LSTM or GRU.
- for a lot of NLP or text processing problems, a BRNN with LSTM appears to be commonly used.
- The disadvantage of BRNN that you need the entire sequence before you can process it. for Example, in Live Speech Recognition if you use BRNN you will need to wait for the person who speaks to Stop to take the entire Sequence and then make Prediction.

## \* Deep RNN

- In a lot of cases the standard one layer RNNs will solve your problem. But in some problems it's useful to stack them. Some RNN layers to make a deeper network.
- For example, a deep RNN with 3 layers would look like this :



- In feed-forward deep nets, there could be 100 or even 200 layers. In deep RNNs stacking 3 layers is already considered deep and expensive to train.
- In some cases you might see some feed-forward network layers connected after recurrent cell, just like in CNN.

### → Back propagation with Deep RNN

- In modern deep learning frameworks, you only have to implement the forward pass, and the framework takes care of the backward pass, so most deep learning engineers do not need to bother with the details of the backward pass.

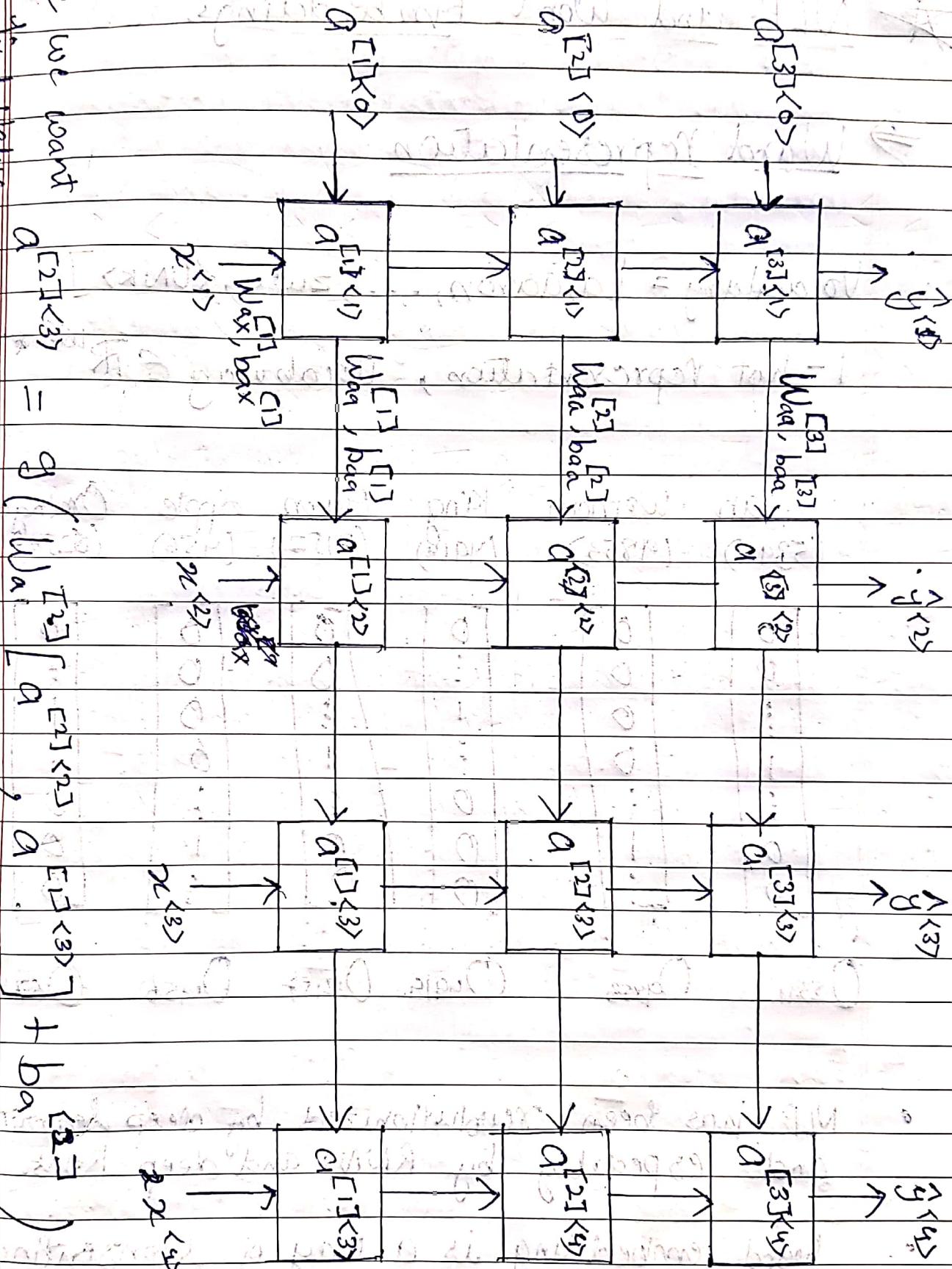
If however you expect to implement this yourself, here is a simple way to do it:

Start at the last layer and work your way back up to the first layer.

For each layer, calculate the gradients with respect to the previous layer's hidden states.

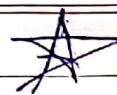
Then update the weights and biases of the current layer based on the gradients.

Suppose we want  
to find value  
of  $a$ :



Final answer is  $a[3][3]$

Divide and  
conquer



## NLP and Word Embeddings.

### ⇒ Word Representation

Vocabulary = [a, aaron, ..., zulu, <UNK>]

1-hot representation, Vocabulary  $\in \mathbb{R}^{10000}$

Man Woman King Queen Apple Orange  
 (5391) (9853) (4914) (7157) (456) (62527)

0	0	0	1	0	0	0
0	0	1	0	0	0	0
:	0	1	:	0	0	0
1	0	1	1	0	0	0
:	:	0	:	0	0	0
0	1	0	0	0	1	0
0	0	0	0	0	0	0

05391 09853 04914 07157 0456 062527

- NLP has been revolutionized by deep learning and especially by RNNs and deep RNNs.

- Word embedding is a way of representing words. It lets algorithm automatically find the analogies between words like "King" and "Queen".

- So far we have defined our language by a vocabulary. Then represent our words with a one-hot vector that represents the word in the vocabulary. As shown on previous page
- We will use the Annotation Index for any word that is represented with one-hot like (in the image).
- One of the weakness of this representation is that it treats a word as a thing itself and it doesn't allow an algorithm to generalize across words.

↳ For example: "I want a glass of apple orange.", a model should predict the next word as juice.

↳ Similar example: "I want a glass of apple", a model won't easily predict juice here if it wasn't trained on it. And if so the two examples aren't related although orange and apple are similar.

- Inner product between any one-hot encoding vector is zero. Also, the distances between them are same.

- So, instead of a one-hot encoding presentation, won't it be nice if we can learn a featurized representation associated with each of these words: man, woman, King, Queen, apple and orange?

man      woman      King      Queen      Apple      Orange  
 (5391)    (9853)    (4914)    (7157)    (456)    (6257)

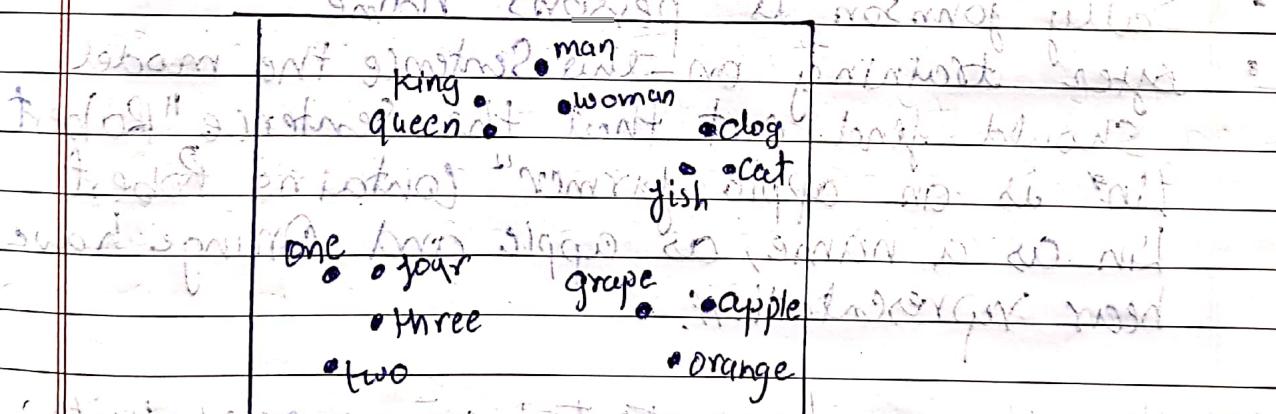
Genders	-1.01	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	0.03
Age	0.03	0.02	0.79	0.69	0.03
Food	0.04	0.01	0.02	0.01	0.95
Verb	0.5391	0.9853	0.4914	0.7157	0.456

- Each word will have a, for example, 300 feature with a type of float point number.
- Each word column will be a 300-dimensional vector which will be the representation

- We will use the notation  $C_{589}$  to describe 'man' word features vector.
- Now, if we return to the example we described again:
  - "I want a glass of orange"
  - "I want a glass of apple"
- Orange and apple now share a lot of similar features which makes it easier for an algorithm to generalize between them.

We call this representation 'Word Embedding'.

→ To visualize word embedding we use a t-SNE algorithm to reduce the feature to 2 dimension which makes it easy to visualize

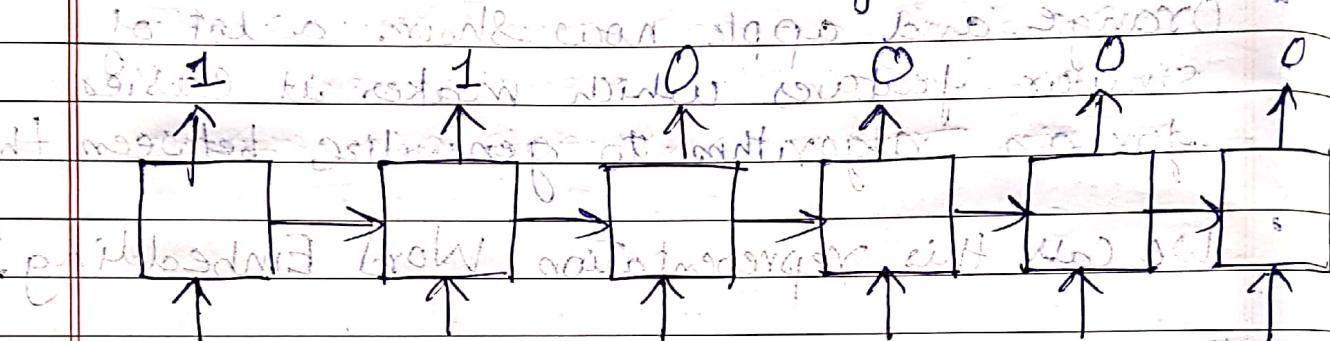


you will get sense that more related words are closer to each other.

→ The Word Embedding came from that we need to embed a unique vector inside a n-dim Space.

## \* Using Word Embeddings.

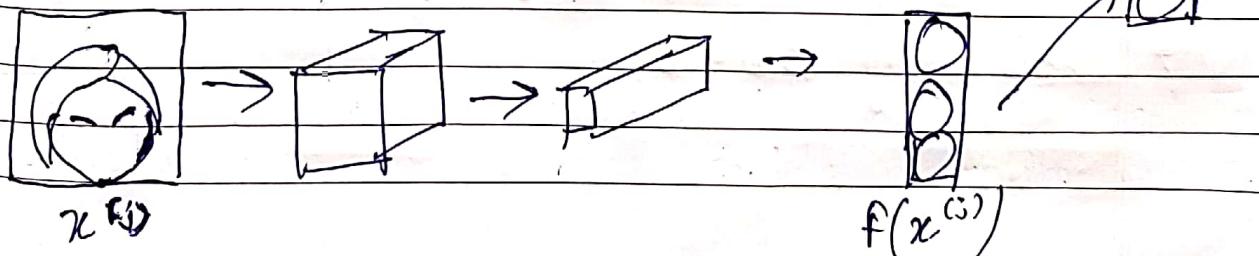
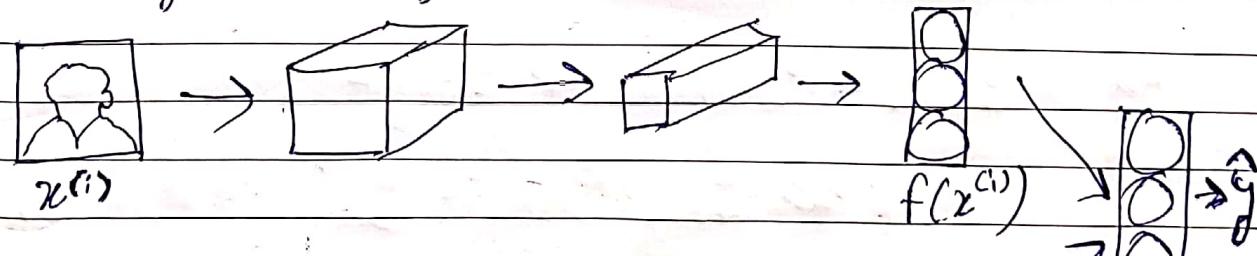
- Let's see how we can take the feature representation we have extracted from each word and apply it in the named entity recognition problem.
- Given this example (from named entity recognition)



Sally Johnson is an apple orange farmer at Robert Lin's farm

- Sally Johnson is person's name
- After training on this sentence the model should find out that the sentence "Robert Lin" is an apple farmer" contains Robert Lin as a name, as apple and Orange have near representations.
- Now if you have tested your model with this sentence "Mahmoud Baldry is a durian cultivator" the network should learn the name even if it hasn't seen the word durian before (during training). That's the power of word representations.

- The algorithms that are used to learn word embedding can examine billions of words of unlabeled text - for example, 100 billion words and learn the representations from them.
- Transfer learning and word embeddings:
  - i) Learn word embeddings from large text corpus. (1-100 Billion of words)
  - ii) Or, download pre-trained embedding online.
  - iii) Transfer embedding to new task with the smaller set (say, 100k words)
- Word embedding tend to make the biggest difference when you're trying to carry out has a relatively smaller training set.
- Also, one of the advantages of using word embedding is that it reduce the size of the input!
  - ↳ 10,000 one-hot compared to 300 feature vector.
- Word embedding have an interesting relation to the face recognition task:



- In this problem, we encode each piece into a vector and then check how similar are these vectors.
- Words encoding and embeddings have similar meaning here.

→ In the word embedding task, we learning ways a representation for each word in our vocabulary (unlike in image encoding where we have to map each new image into some n-dimensional vector). ~~we will discuss the algorithms in next section.~~

Forward pass, action at first, without loss of generality

Let's start with forward propagation

→ Forward propagation consists of two steps

1. Linear transformation step:  $y = w_1x + b_1$

2. Non-linear activation function step:  $y = \sigma(w_1x + b_1)$

action of softmax at horizontal layer shown at

at bottom of handwritten material mentioned above

→ Action of softmax at horizontal layer shown at

at bottom of handwritten material mentioned above

→ Action of softmax at horizontal layer shown at

at bottom of handwritten material mentioned above

→ Action of softmax at horizontal layer shown at

at bottom of handwritten material mentioned above

→ Action of softmax at horizontal layer shown at

at bottom of handwritten material mentioned above

→ Action of softmax at horizontal layer shown at

at bottom of handwritten material mentioned above

→ Action of softmax at horizontal layer shown at

at bottom of handwritten material mentioned above

## \* Properties of Word Embedding.

- One of most fascinating properties of word embedding is that can also help with analogy reasoning. While analogy reasoning may not be by itself the most important NLP application, but it might help convey a sense of what these word embedding can do.
- Analogies Example:  
Given this word embedding table.

	Man	Woman	King	Queen	Apple	Orange
(5391)	(9853)	(4914)	(7157)	(456)	(62579)	
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	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97
	eman	ewoman	king	ueen	pple	orange

- Can we conclude this relation:
  - ↳ Man → Woman
  - ↳ King → ??

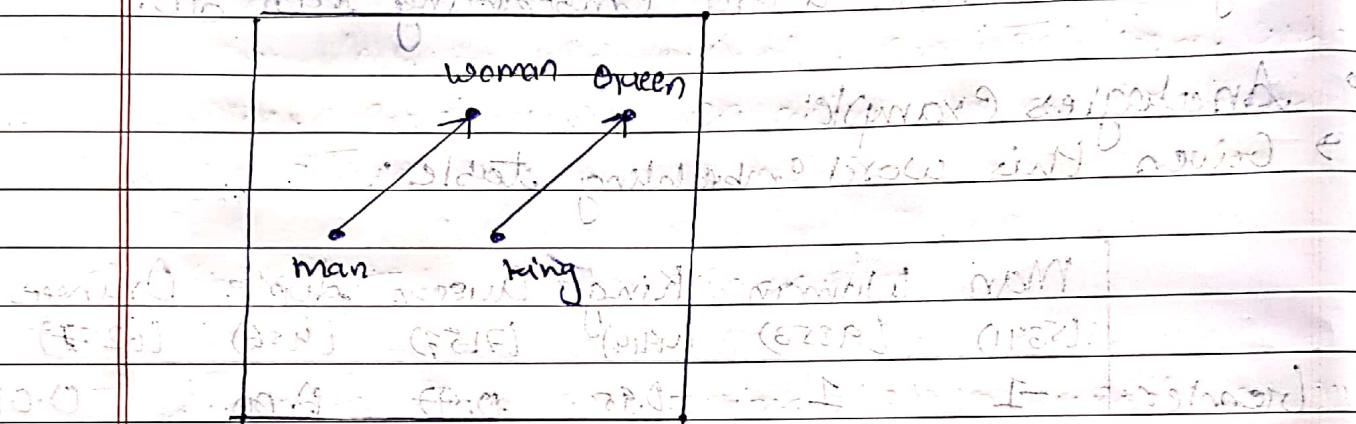
- Lets Subtract Eman from Ewoman.

This will equal to Vector:

$$\text{Cman} - \text{Cwoman} = [-2 \ 0 \ 0]$$

$$\text{Sking} - \text{Cman} = [2 \ 0 \ 0]$$

So the difference is about the gender in both



→ This vector represent the gender

→ The drawing visualization of the 4D vector in 2D form that has been extracted by a t-SNE algorithm. It's a drawing just for visualization. Don't rely on the t-SNE algorithm for finding parallel parallels.

So we can reformulate the problem to find:

$$\rightarrow \text{Cman} - \text{Cwoman} \approx \text{Cking} + P??$$

$$P?? \approx \text{Cking} - \text{Cman} + \text{Cwoman}$$

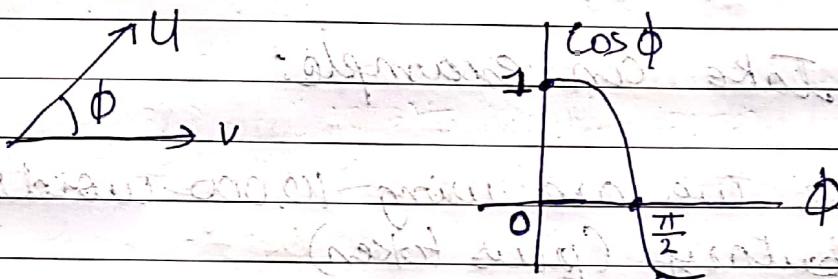
$$\rightarrow \text{let's say } P_W = P??$$

∴ it can also be represented mathematically by:

$$\rightarrow \text{argmax } \text{Sim}(P_w, P_{\text{king}} - P_{\text{man}} + P_{\text{woman}})$$

$$\Rightarrow \text{Sim}(U, V) = \frac{U^T V}{\|U\|_2 \|V\|_2} = \frac{U^T V}{\sqrt{\sum_{i=1}^k U_i^2} \times \sqrt{\sum_{i=1}^k V_i^2}}$$

→ this is called Cosine Similarity.



• Cosine Similarity - the most commonly used Similarity function.

→ In formula, The top part represent the inner product of U and V vector. It will be large if the vectors are very similar.

• We can also use Euclidean distance as a Similarity function (but it rather measures a dissimilarity, so we should take it with negative sign.)

• We can use this equation to calculate the similarities between word embedding and on the analogy problem where  $U = P_w$  and  $V = P_{\text{king}} - P_{\text{man}} + P_{\text{woman}}$



## Embedding Matrix

⇒ Learning word embeddings

• Only start learning some algorithms that can

- When you implement an algorithm to learn a word embedding, what you end up learning is a embedding matrix.
- Let's take an example:

→ suppose we are using 10,000 words as our vocabulary. (plus token)

- The algorithm should create a matrix  $E$  of the shape  $(300, 10000)$  in case we are extracting 300 features.

miniaturization orange generalized (univ)

300

exit

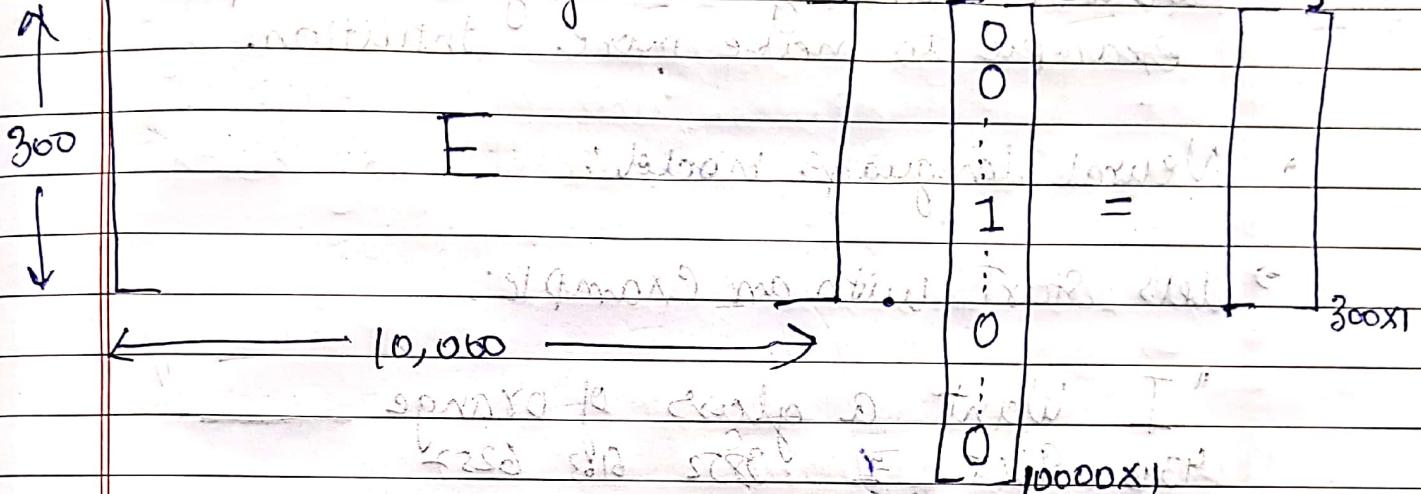
10,000

- If  $O_{6257}$  is the one-hot encoding of the word Orange of shape  $(10000, 1)$ , then,

$$\text{np.dot}(E, O_{6257}) = C_{6257} \text{ which shape is } (300, 1)$$

- Generally  $\text{np.dot}(E, O_i) = C_i$

a vector ... orange ... zulu <unk> ...  $O_i$  ...  $C_i$



- In the next Sections, we will see that we first initialize  $E$  randomly and then try to learn all the parameters of this matrix.
- In practical it's not efficient to use a dot multiplication when you are trying to extract the embedding of a specific word instead we instead we will use slicing to slice a specific column. In Keras there is an embedding layer that extracts this column with no multiplication!

# \*enlive

## Learning Word Embedding: Word2Vec & GLOVE

Date : \_\_\_\_\_  
Page: \_\_\_\_\_

⇒ learning word embeddings

• Let's start learning some algorithms that can learn word embeddings.

• At the start, word embeddings algorithms were complex but then they got simpler and simpler.

• We will start by learning the complex example to make more intuition.

• Neural language model:

⇒ Let's start with an example:

"I Want a glass of orange"

4343 9665 I 3852 6183 6252

⇒ we want to build a language model so that we can predict the next word.

⇒ So we use this neural network to learn the language model.

⇒ we get  $C_j$  by np.dot( $E, O_j$ )

↳ NN layer has parameter  $W_1$  and  $b_1$  while

softmax layer has parameters  $W_2$  and  $b_2$

↳ Input dimension is  $(300 * 6, 1)$  if the window size is 6 (six previous words)

One-hot

E.O.

Encoding

I

$$O_{4343} \rightarrow E.O_{4343} \rightarrow P_{4343}$$

Want

$$O_{9665} \rightarrow E.O_{9665} \rightarrow P_{9665}$$

a

$$O_1 \rightarrow E.O_1 \rightarrow P_1$$

glass

$$O_{3852} \rightarrow E.O_{3852} \rightarrow P_{3852}$$

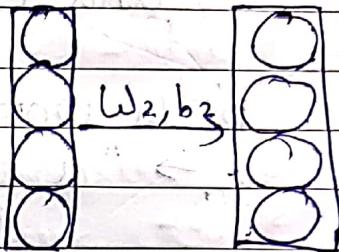
of

$$O_{6163} \rightarrow E.O_{6163} \rightarrow P_{6163}$$

orange.

$$O_{6257} \rightarrow E.O_{6257} \rightarrow P_{6257}$$

Softmax



(10000x1)

(300x6, 1)

→ here  $\hat{y}$  will be  $(10000 \times 1)$  vector which represent any word from Vocabulary. As  $y = \text{'Juice'}$ , it will backpropagate and learn all the weights according to true output.

- We are optimizing E matrix and layer parameters.
- We need to maximize the likelihood to predict the next word given the context (previous words)
- This model was build in 2003 and works pretty decent for learning word E.
- In this ~~first~~ Example we took a window of 6 words that fall behind that we want to predict. There are other choices when we are trying to learn word embeddings.
- ↳ Suppose we have an example:

"I Want a glass of orange juice to go along with my Cereal"

- To learn 'juice', choice of context etc:

a) Last 4 words.

↳ we use a window of last 4 words  
( $h$  is a hyperparameter), "a glass of orange" and try to predict the next word from it.

b) 4 words on the left and on the right.

↳ "a glass of orange" and "to go along with" a glass of orange to go along with.

c) Last 1 word

↳ orange

d) Nearby 1 word. Around word black and

↳ "glass" word is near juice.

↳ This is other idea of Skipgrams model.

↳ The ~~old~~ idea is much simpler and

works remarkably well.

↳ we will talk about this in the next

Section.

• Researchers found that if you really want to build a language model, it's natural to use the last few words as a context. But if your main goal is really to learn a word embedding then you can use all of these other contexts and they will result in very meaningful word embeddings as well.

- To summarize, the language modeling problem poses a machine learning problem where you input the context (like the last four words) and predict some target words. And posing that problem allows you to learn good word embeddings.

## \* Word2Vec

- Before presenting Word2Vec, let's talk about SKIP-GRAMS:

For example, we have the sentence: "I want a glass of orange juice to go along with cereal".

↳ We will choose Context and target

↳ The target is chosen randomly based on a window with a specific size.

Context      target      how far

Orange	juice	+1
Orange	glass	-2
Orange	my	+6

We have converted the problem into a supervised problem.

This is not an easy learning problem. Learning within 40/10 words (10 is an example) is hard. We want to learn this to get our word embeddings model.

- Word2Vec model.

→ Vocabulary size = 10,000 words

→ Let's say that the context word are 'c' and the target word is 't'

→ We want to learn 'c' to 't'

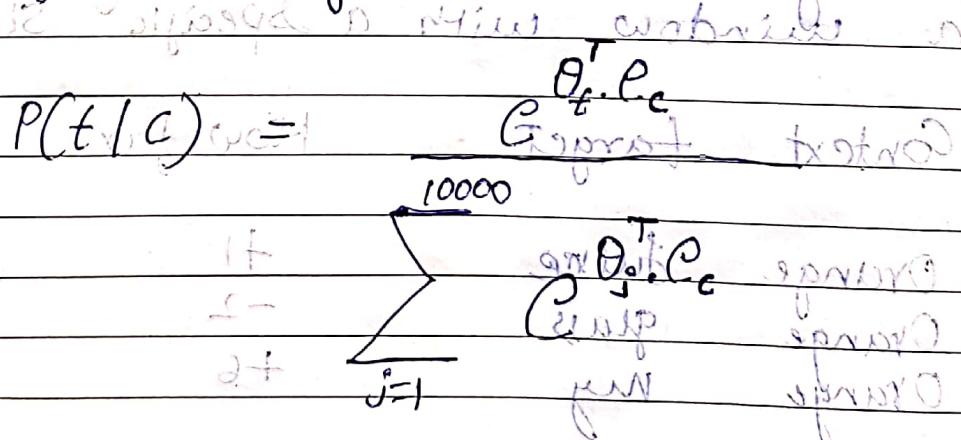
→ We get  $\theta_c$  by  $E_{\theta_c}$

→ We then use softmax layer to get  $P(t|c)$  which is

→ Also we will use the cross-entropy loss function.

→ This model is called Skip-gram model

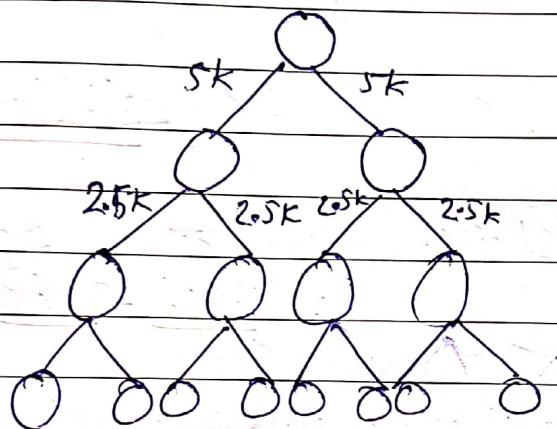
• The last model has a problem with the



→ Here we are summing 10000 numbers which corresponds to the number of words in our vocabulary.

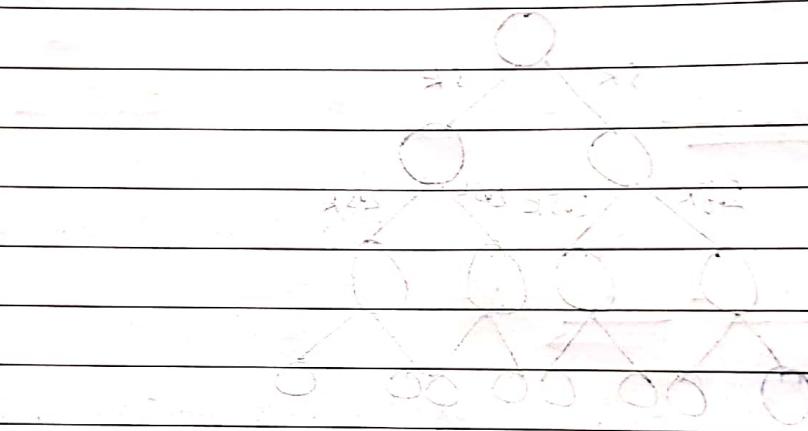
→ if this number is larger say 1 million, the computation will become very slow.

- One of the solutions for the last problem is to use "Hierarchical softmax classifier" which works as a tree classifier.



- In practice, the hierarchical softmax classifier doesn't use a balanced tree like the drawn one. Common words are at the top and less common are at the bottom.
- How to sample the context C?
- One way is to choose the context by random from your corpus.
- if you have done it that way, there will be frequent words like "the, of, a, and, to, .." that can dominate other words like 'Orange, apple, durian, ..'.
- In practice, we don't take the context uniformly random instead there are some heuristics to balance the common words and the non-common words.

- Word2Vec paper includes 2 ideas of learning word embedding. One is Skip-gram model and another is CBOW (Continuous bag-of-words)



difficult question because it's often not a  
good basis for normalizing the hidden  
unit to get desired movement and much  
difficulty with the most important part of the

skip-gram unit. Several other approaches  
to this problem have been proposed.  
One approach is to use a softmax function  
with a small batch of words. Another approach  
is to use a neural network to predict the  
next word. This can be done by training a  
recurrent neural network (RNN) to predict