

RNN

# ⇒ Named entity recognition

$x$  :- Mary ~~father~~ and Hermione Granger  
 $x^{(1)}$   $x^{(2)}$   $x^{(3)}$   $x^{(4)}$   $x^{(5)}$

$y$  :- 1 1 0 1 1 0 0 0 0  
 $y^{(1)}$   $y^{(2)}$   $y^{(3)}$

$T_x = 9$   
 $T_y = 9$

$\Rightarrow T_x^{(i)}$   
 $T_y^{(i)}$  }  $i^{th}$  training sample

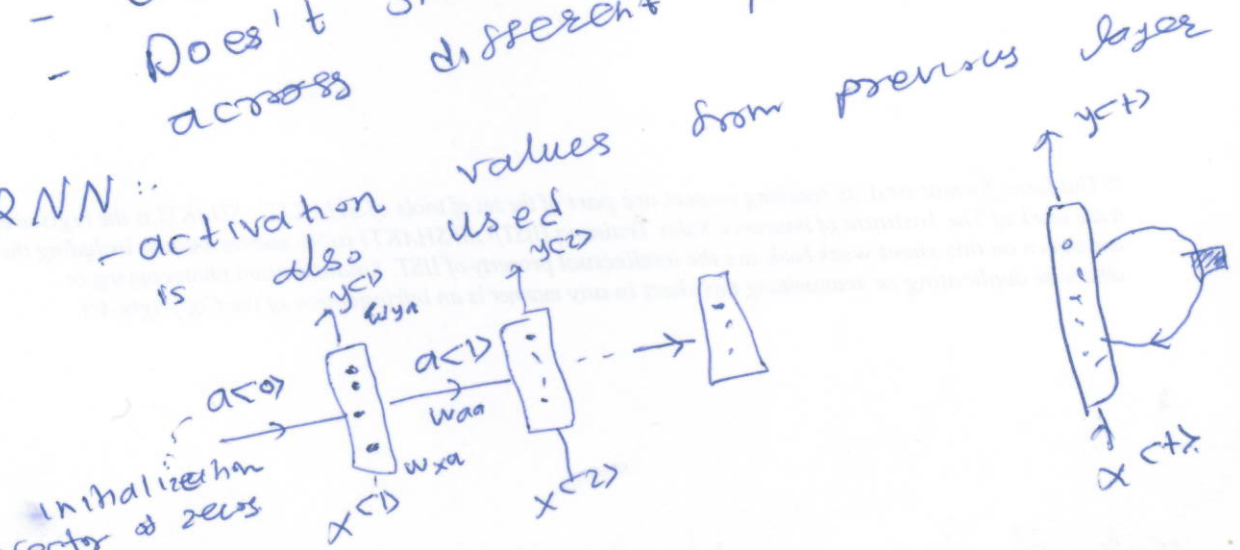
- ⇒ NLP
- dictionary of size 10k
  - Now a days 30k to 100k are used often as well
  - one-hot vector for each word
  - <unk> token for unknown word.

⇒ Why not standard network?

- Problems:
  - different length
  - Doesn't share features learned across different position of text

- RNN:

- activation is also used



- We also need to incorporate information from next words.
- We will address it later in bi-directional RNN

## \* Forward Propagation

$$\Rightarrow a^{(0)} = \vec{0}$$

$$\Rightarrow a^{(1)} = g(w_{aa} a^{(0)} + w_{ax} x^{(1)} + b_a)$$

$$\hat{y}^{(1)} = g(w_{ya} a^{(1)} + b_y)$$

→ sigmoid (binary)  
softmax (k-way)

tanh more popular in RNN

$$\Rightarrow a^{(t)} = g(w_{aa} a^{(t-1)} + w_{ax} x^{(t)} + b_a)$$

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

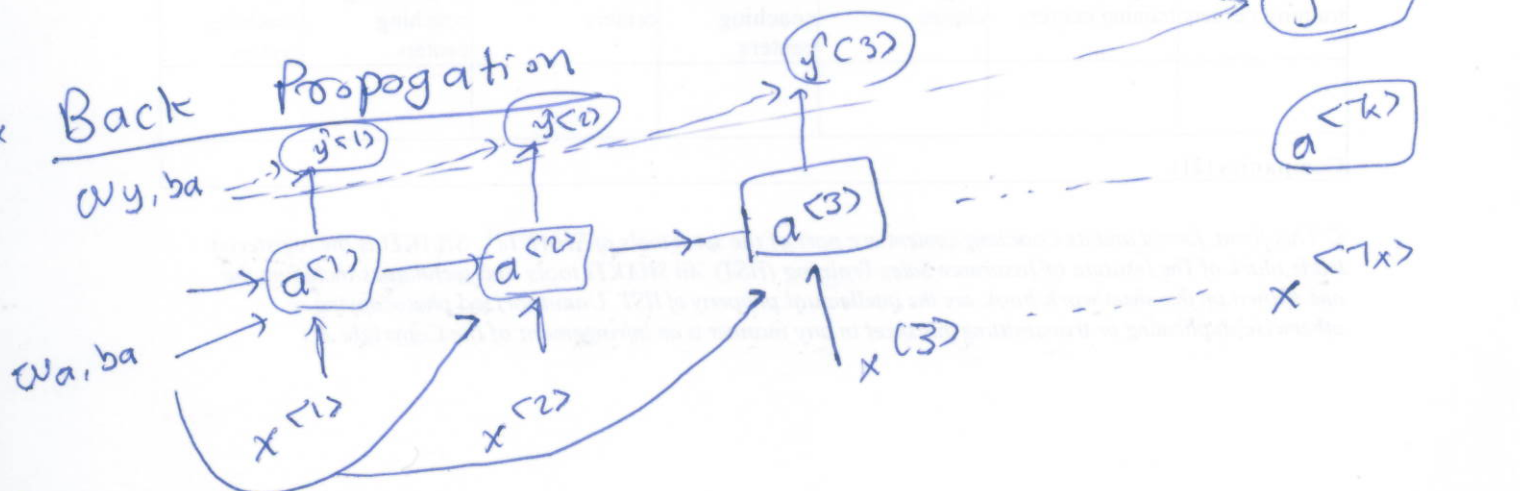
$$\Rightarrow a^{(t)} = g(w_a [a^{(t-1)}, x^{(t)}] + b_a)$$

where  $w_a = [w_{aa} \mid w_{ax}]$

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

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

## \* Back Propagation



→ Cross entropy loss

$$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = - \hat{y}^{<t>} \log \hat{y}^{<t>} - (1 - \hat{y}^{<t>}) \log (1 - \hat{y}^{<t>})$$

$$L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

⇒ Backpropagation through time

\* Different types of RNN

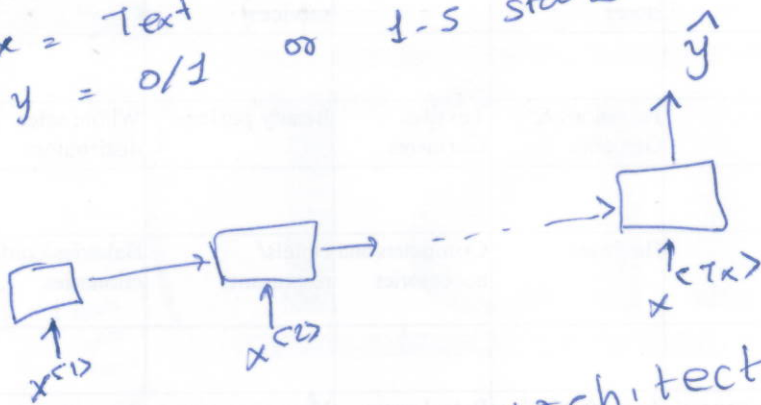
①  $T_x = T_y$

- we were talking about this case so far
- Many to many architecture.

② Sentiment classification

$x = \text{Text}$

$y = 0/1$  or 1-5 stars



- many to one architecture

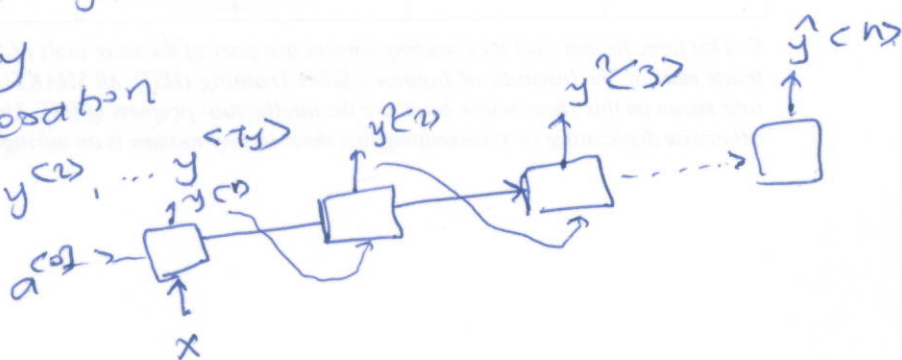
③ one to one

- Not much interesting
- standard general NN

④ one to many

- Music generation

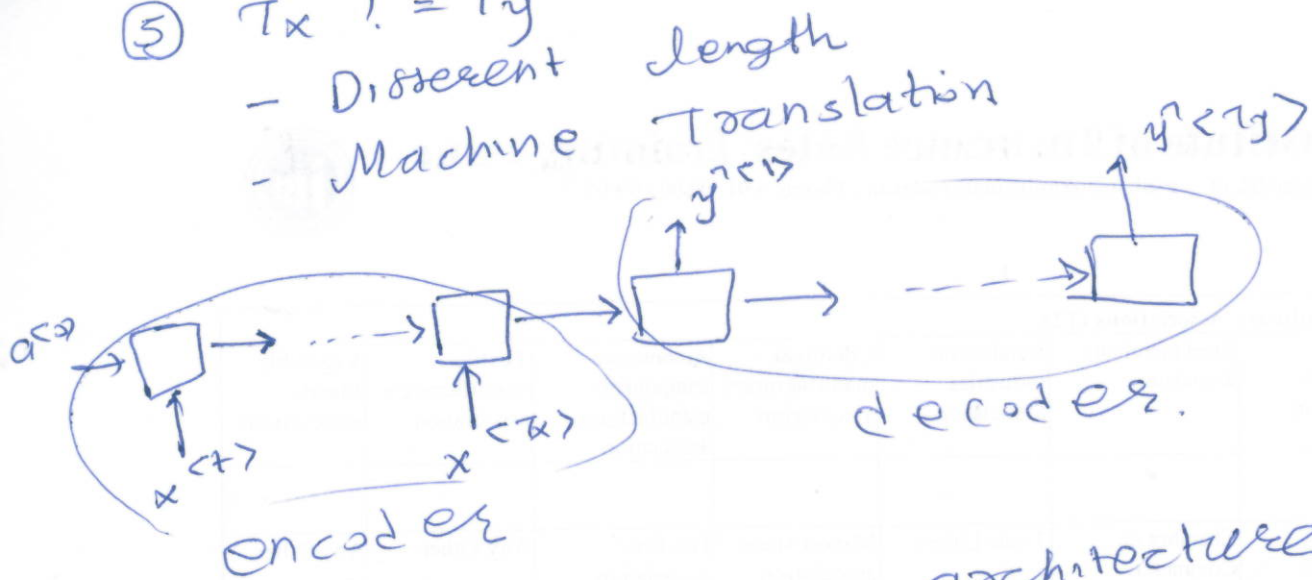
$x \rightarrow y^{<1>}, y^{<2>}, \dots, y^{<T_y>}$



③



- ⑤  $T_x \neq T_y$
- Different length
  - Machine Translation



- Many to many architecture

### \* Language Modeling

$\Rightarrow P(\text{sentence}) = ?$

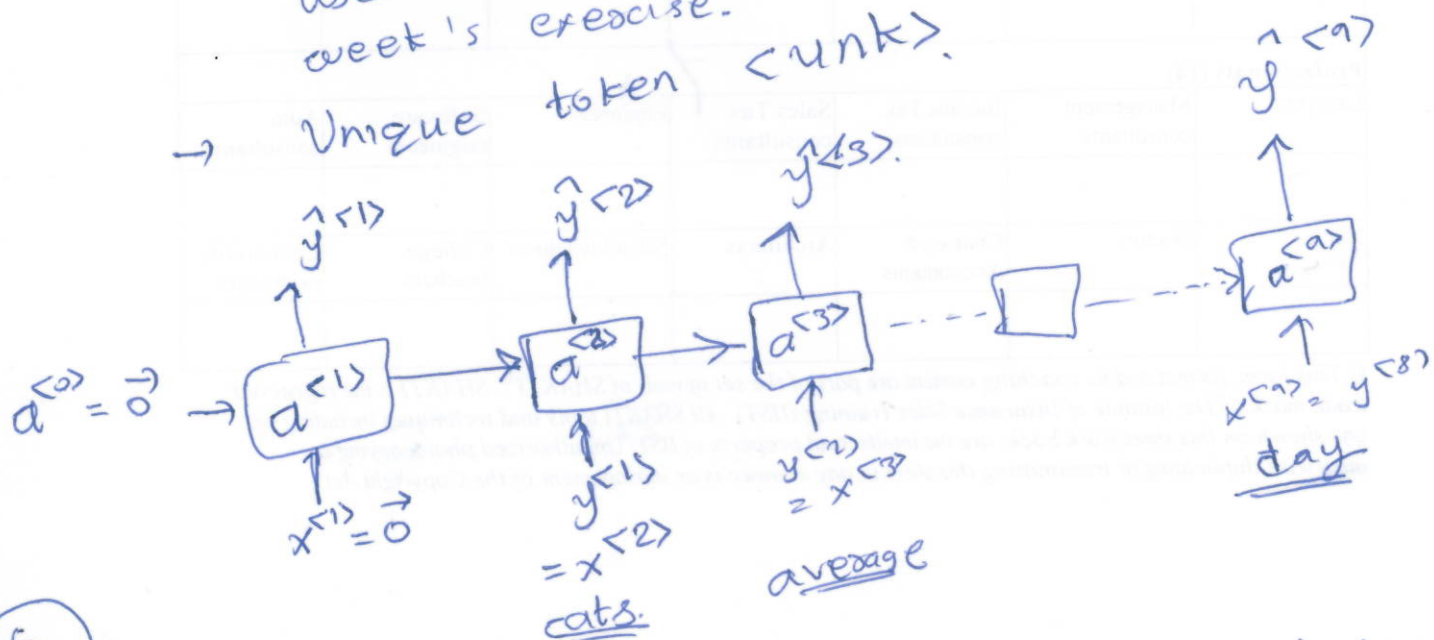
- Fundamental component of speech recognition and machine translation

→  $P(y^{(1)}, y^{(2)}, \dots, y^{(T)})$

- Training set: large corpus of English set

- End of sentence (Eos) taken can be used. we will not talk about it in this week's exercise.

- Unique token <unk>.



④

- cats average 15 hours of sleep a day

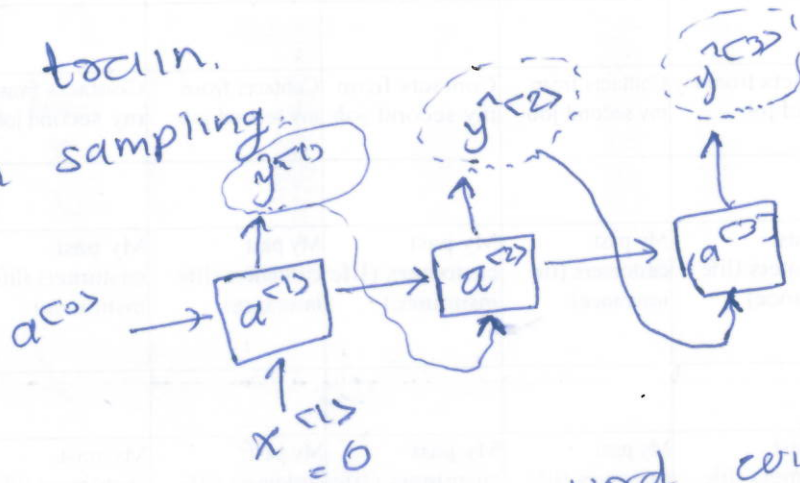
$$\Rightarrow L^{(t)}(\hat{y}^{(t)}, y^{(t)}) = - \sum_i y_i^{(t)} \log y_i^{(t)}$$

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

$$\begin{aligned} P(y^{(1)}, y^{(2)}, y^{(3)}) \\ = P(y^{(1)}) P(y^{(2)} | y^{(1)}) \\ P(y^{(3)} | y^{(1)}, y^{(2)}) \end{aligned}$$

### \* Sample Novel sequences

- First train.
- Then sampling.



- $\hat{y}^{(t)} \rightarrow$  choose the word with highest probability
- If  $\langle \text{unk} \rangle$  comes ignore it.
- When  $\langle \text{eos} \rangle$  comes sentence has been generated.

- we can also build character level language model.

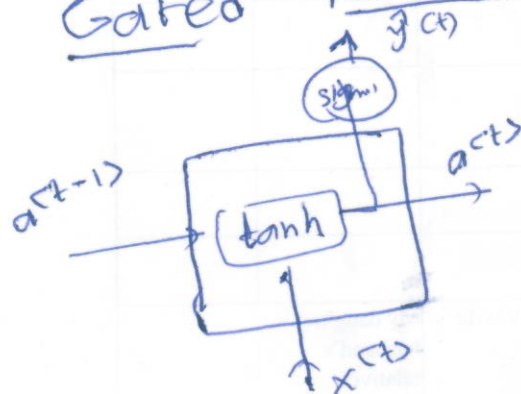
Pros: - don't need to worry about unknown tokens

cons: - Much longer sequences.  
- Not used generally.

## \* Vanishing gradient

- The cat, ..., was full
- The cats, ..., were full
- For DNN, back propagation can have little impact on earlier layers
- Idea is to have longer length dependency.

## => Gated Recurrent Unit (GRU)



- notion of cell  $c$  = memory cell
  - to remember whether it was cat or cats

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

- It will be different LSTM

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

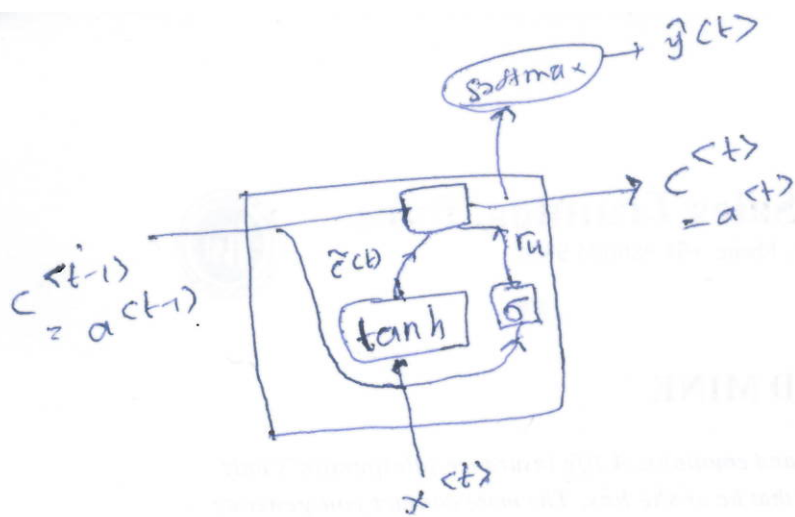
- Gate  $\Gamma_u = \sigma(W_u [c^{(t-1)}, x^{(t)}] + b_u)$

- Intuition: think it as 0 or 1

- Gate tells when to update  $c^{(t)}$

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



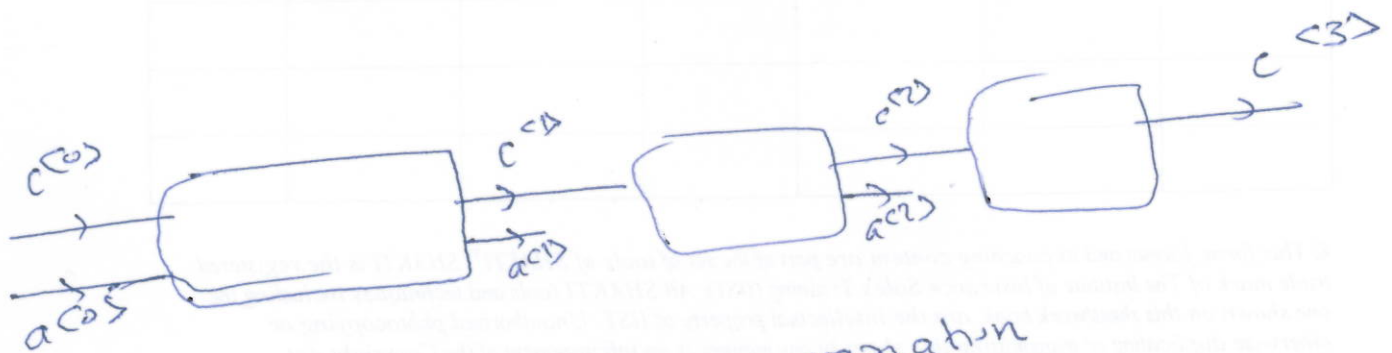


$\Rightarrow r_u$  is generally close to zero and hence mostly  $c^{(t)} = c^{(t-1)}$  it helps with vanishing gradient

$\Rightarrow$  if you have too activation  $c^{(t)}$ ,  $r_u$  will also be too dimensional

## \* LSTM

- we don't have  $a^{(t)} = c^{(t)}$  here
- Two gates  $r_u$  and  $r_f$  ( $r_f$  serves as  $1 - r_u$ )
- one more output-gate  $r_o$
- $c^{(t)} = r_u * c^{(t-1)} + r_f * c^{(t-1)}$ 
  - $u$  - update
  - $f$  - forget
  - $o$  - output
- $a^{(t)} = r_o * c^{(t)}$

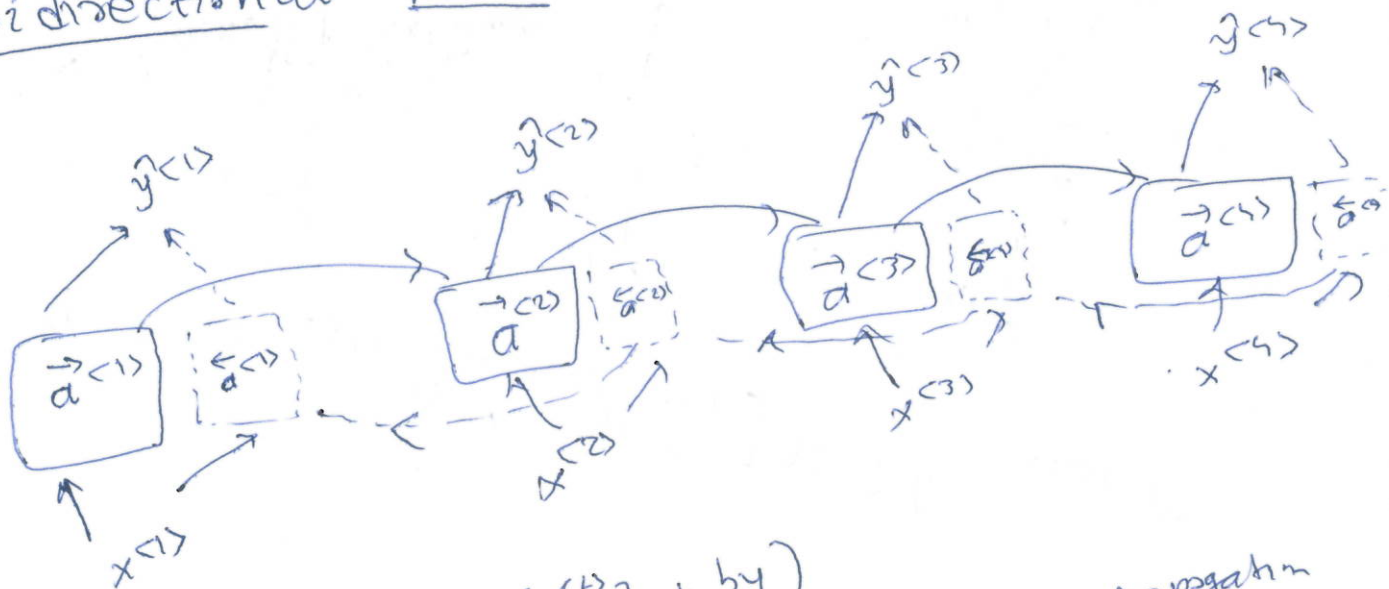


- There are some variations



- GRM is more faster
- LSTM is more powerful
  - people use it as default first thing to try.

## \* Bidirectional RNN (BRNN)



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

- This is still in backward direction

⇒ He said, Teddy Resault

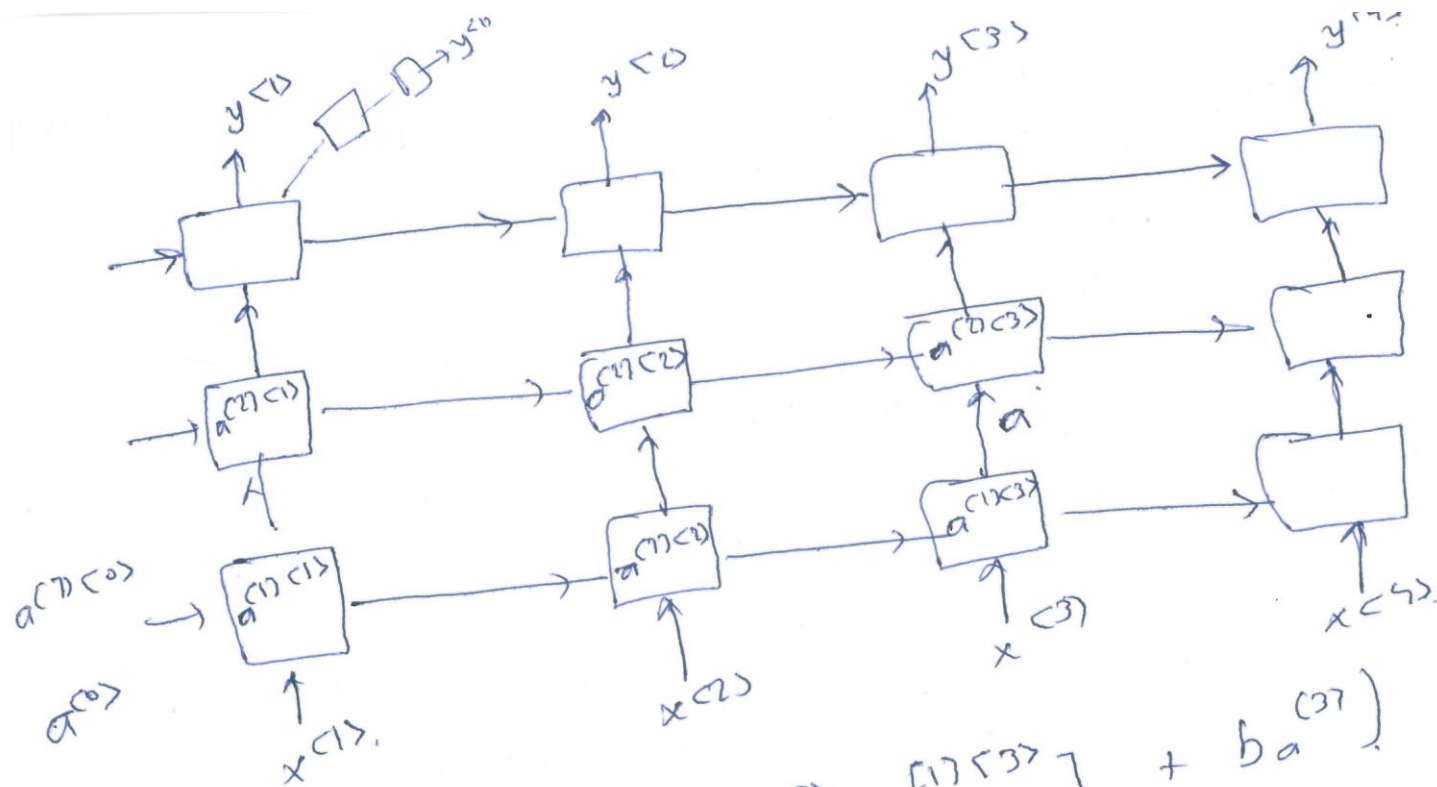
- blocks can be RNN, GRU or LSTM

⇒ Disadvantage:

- you need entire sequence of data

## \* Deep RNN





$$a^{(2)}(t) = g(Wa^{(1)}(t) + ba^{(0)}(t))$$

# NLP & Word Embedding

## \* Word Representation

- 1 hot vector
  - each token is independent
  - inner product of all words is 0
  - orange juice / apple juice
- featurized representation
  - features for gender, royal, age, food
  - 300 features
- t-SNE algo
  - 300 dimension to 2 dimension and visualize

## \* Using word embeddings

- ⇒ Named entity recognition
- learn embedding from billion words
  - transfer this knowledge to task where named entity where training data is sample
    - Transfer learning
    - we also use less dimensional dense vector
  - Generally used when you have relatively small dataset for task B.
  - Relation to face encoding
    - getting vector for a face
    - vector of activation of one layer of classifier

## \* Properties of word embeddings

- Man → woman as king → ?
- $\arg \max_w \text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$
  - King word  $e_w$  is in 300 D, not in 2D (after t-SNE).
  - Where works in 300 D, not in 2D (after t-SNE).
    - t-SNE is non-linear



- ⇒ Similarity function
- cosine similarity → used often
  - Euclidean

\* Embedding Matrix

- getting vector from matrix (Multiply with one-hot)

$$\text{word} = E \cdot \text{O}_{\text{word}}$$

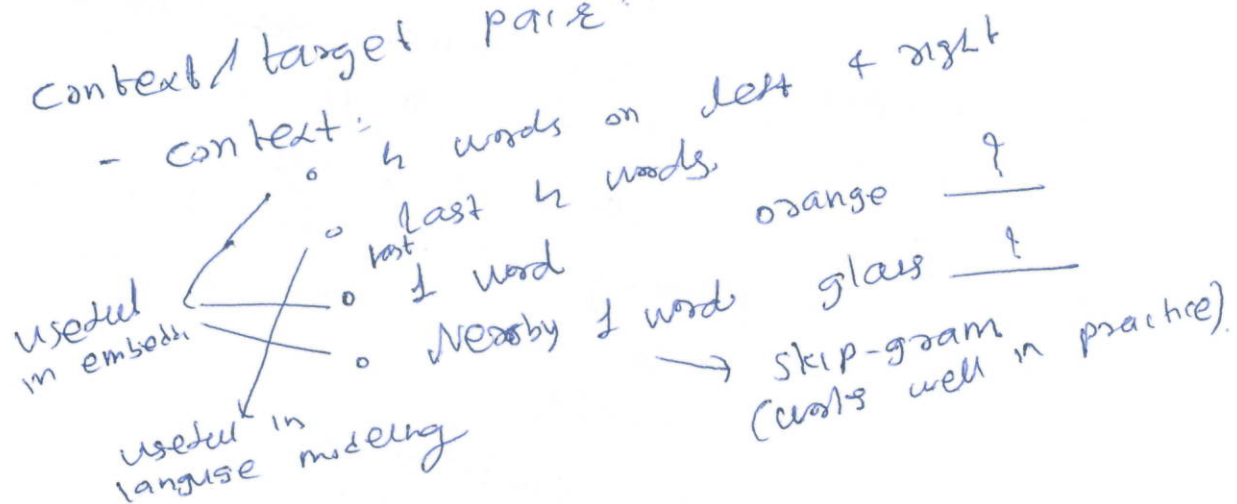
$$(300 \times 10,000) \cdot (10,000 \times 1)$$

- while learning you initialize  $E$  randomly and learn values via gradient descent

\* Learning word embeddings

- starting with historic models before going for word2vec or glove
- classifier to predict next word
- we will do repeated training
  - parameters to learn are  $E, w^{(n)}, b^{(n)}$

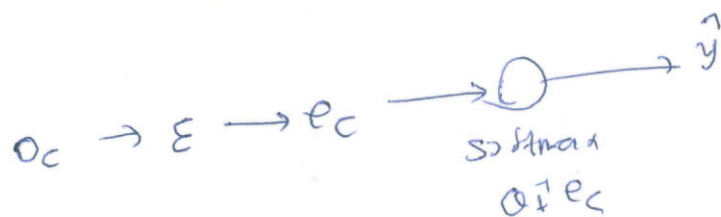
⇒ context / target pairs



## \* Word 2 Vec

⇒ skip-grams -

- context / target pair for supervised learning
- randomly pick up target in a window of  $\pm 5$  words



$$p(t/c) = \frac{e^{o^T e_c}}{\sum_{j=1}^{10000} e^{o^T e_c}}$$

$$L(\hat{y}, y) = - \sum_{i=1}^{10000} y_i \log \hat{y}_i$$

$\hat{y}_i \rightarrow$  one hot vector  
 $y_i \rightarrow$  probabilities of 10000 target words

⇒ Hierarchical softmax

- solves the problem of summing 10000 exponential
- In practice it does not use symmetric tree
  - more common words on top, less used words in deep.

⇒ How to sample context c?

- the, of, a, and
- orange, apple, durian
- some heuristic to sample less frequent words equally.

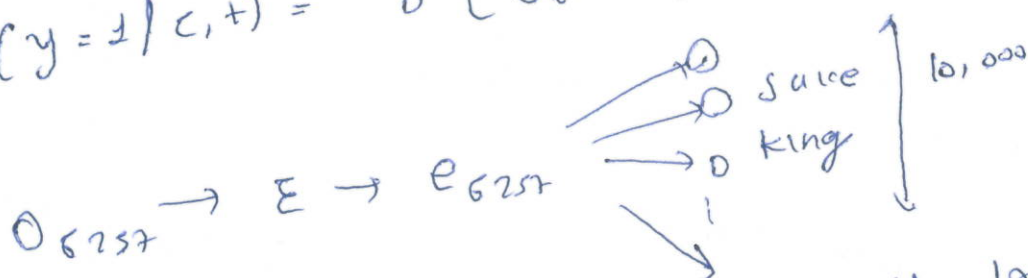
⇒ Problem of skip-gram is to calculate denominator of softmax.  
 - solution is negative sampling

## \* Negative sampling

context	word	target ?	
orange	juice	1	← Pick-up soon → to window
orange	king	0	← Pick up randomly
orange	book	0	
orange	the	0	
orange	of	0	
	y		

$k = 5$  to  $20$   
 smaller dataset  
 $k = 2-5$   
 for larger dataset

$$p(y=1|c, t) = \sigma(\mathbf{Q}_t^T \mathbf{e}_c)$$



- instead of 10,000 we will learn 5 of them
- 10,000 binary classification is simpler than 10000 way softmax

→ Sampling negative sample

- Upon seq.
  - high chances of a, an, the.
- $$p(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10000} f(w_j)^{3/4}}$$

## \* Glove

- Even more simpler, ~~it~~  
- slowly picking up but not as popular as word2vec
- $x_{ij}$  = # times  $i$  appears in context of  $j$
- $x_{ij} = x_{ji}$  or  $x_{ij} \neq x_{ji}$  depending upon your context

$$\Rightarrow f(x_{ij}) = \begin{cases} 0 & \text{if } x_{ij} = 0 \\ \text{empirical weight} & \text{otherwise} \end{cases}$$

weight frequent / infrequent word

$$\text{minimize } \sum_{i=1}^{100000} \sum_{j=1}^{100000} f(x_{ij}) (\mathbf{O}_i^T \mathbf{e}_j + b_i + b_j' - \log x_{ij})^2$$

$$e_w^{\text{final}} = \frac{e_w + \mathbf{O}_w}{2}$$

because here  $e$  and  $\mathbf{O}$  plays symmetric role

## \* Sentiment classification

- Output is 5 star review
- You can quantify comments on social media as well
- Generally training data is less

$\Rightarrow$  Simple model

- o get 300 D vector and average it
- o and pass it to softmax classifier (1-5)
- o Problem is it ignores word-order



⇒ RNN for sentiment classification

- feed embedded vectors in RNN
- Many to one RNN

## \* Debiasing word embeddings

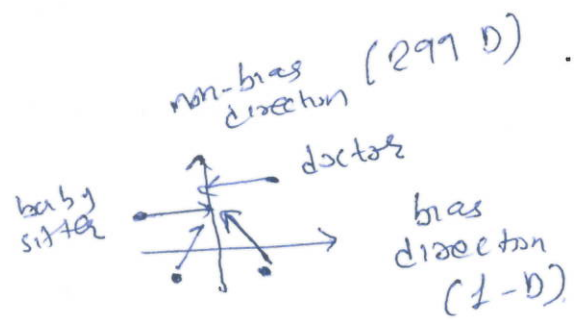
1. Identify bias direction

ehe - eshe

emale - esenal

⋮

average



2. Neutralize  
- for words like doctor, baby-sister,  
project them on non-bias direction

3. Equalize pair  
- grand-mother & grand father  
- move them to equal distance from  
non-bias direction

⇒ One researcher trained classifier to and out  
which words to neutralize or equalize

## Sequence models & Attention mechanism

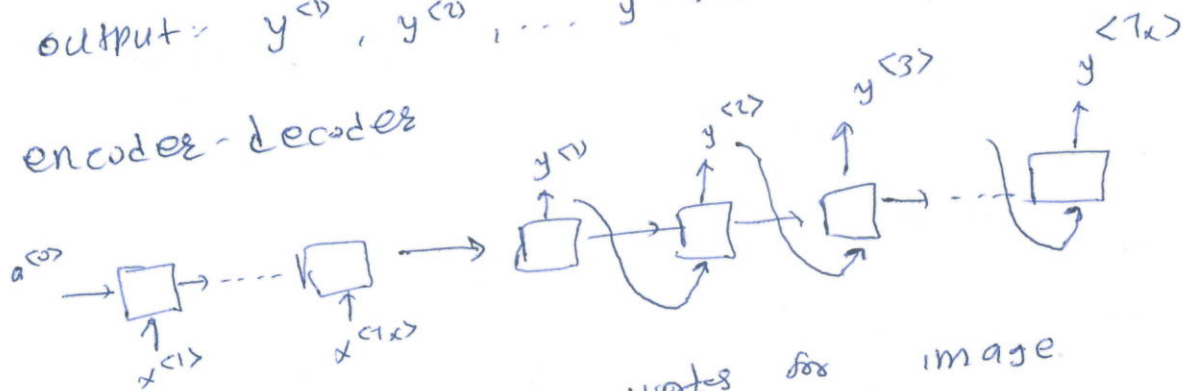
\*

Week-3

## Sequence to sequence Model

⇒ Input:  $x^{(1)}, x^{(2)}, \dots, x^{(T_x)}$   
 output:  $y^{(1)}, y^{(2)}, \dots, y^{(T_y)}$

⇒ encoder-decoders



⇒ Similar architecture works for image captioning

• CNN followed by RNN

→ You don't want ~~most likely~~ <sup>or</sup> best sentence, but

## \* Picking most likely sentence

→ In sentence generators we would pick up words based on their probability.

→ which is prediction  $P(y^{(1)}, y^{(2)}, \dots, y^{(T_y)})$

— Machine translation is "conditional language model", which models  $P(y^{(1)}, y^{(2)}, \dots, y^{(T_y)} | x^{(1)}, x^{(2)}, \dots, x^{(T_x)})$

$$\arg \max_{y^{(1)}, \dots, y^{(T_y)}} P(y^{(1)}, \dots, y^{(T_y)} | x)$$

⇒ greedy search:

— Most likely first word

⇒ Exact search has time complexity

— so we go for approximate search

## Beam Search

- Parameter beam-width,  $B=3$
- choose 3 most likely choices for first word.
- <sup>keep</sup> top 3 choices after 2<sup>nd</sup> step as well
- It is not exponential
- After every step we have 3 choices only
- At every ~~step~~ step you have 3 copies of network
  - you need just 3 copies, not 30000 copies
- setting  $B=1$ , makes it greedy search

## Refinements to Beam Search

- ⇒ Length normalization :-
- Multiplying multiple probabilities make it a tiny no.
  - $\arg \max_y \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$
  - $\arg \max_y \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$
  - $\frac{1}{T_y} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$
- $\alpha = 1$  → full normalization  
 $\alpha = 0$  → No normalization  
 $\alpha = 0.7$  → can be tuned

⇒ How to choose beam width  $B$  ?

- very large :-
- better result
  - slower

→ small  $B$  :-

- worse result
- faster
- lower memory requirement

→ On products we see beam width around 10

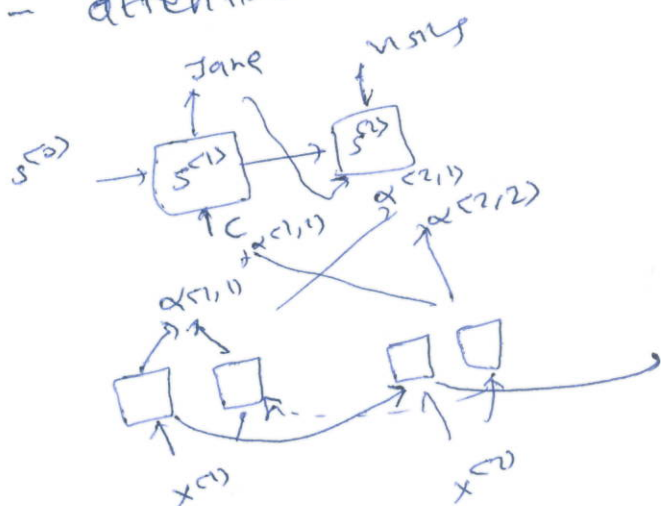


## Error Analysis in Beam Search

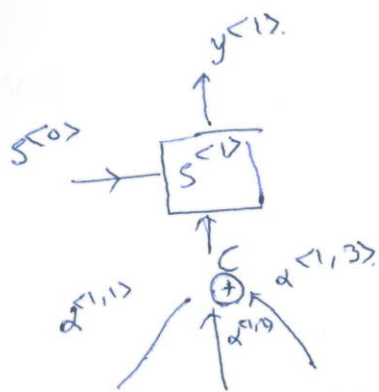
- where to spend time on?
  - beam search or RNN
- Compute  $P(y^*/x)$  and  $P(\hat{y}|x)$  using RNN
- $P(y^*/x) > P(\hat{y}|x)$ 
  - Beam search chose  $\hat{y}$
  - But  $y^*$  attains higher  $P(y/x)$
  - Beam search is at fault
- $P(y^*/x) \leq P(\hat{y}|x)$ 
  - $y^*$  is better than  $\hat{y}$
  - But according to RNN it is worse
  - RNN is at fault
- Rare error analysis for each sample in Dev set.
  - And see who is at fault more times

## \* Attention Model

- Human translator would work part by part
- encoder-decoder works well for smaller sentences
- attention model computes attention weight



$$a^{<t>} = (\vec{a}^{<t>}, \overleftarrow{a}^{<t>})$$

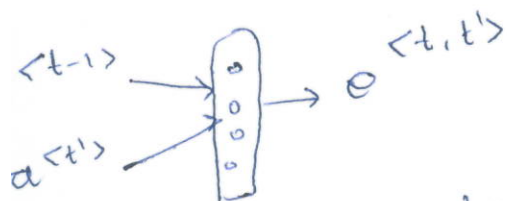


$$\sum_{t'} \alpha^{<1,t'>} = 1$$

$$e^{<1>} = \sum_{t'} \alpha^{<1,t'>} a^{<t'>}$$

$$a^{<t>} = (\vec{a}^{<t>}, \overleftarrow{a}^{<t>})$$

$$\alpha^{<t,t'>} = \frac{\exp(e^{<t,t'>})}{\sum_{t'=1}^T \exp(e^{<t,t'>})}$$



Runs in quadratic cost:  
 $O(n \times m)$

- Also applied in image captioning
- Date in various format as input, standard format as output
- visualization of attention weights

# Speech Recognition

MFCC in pre-processing

Attention model is one way

⇒ Another is CTC

- CTC = Connectionist temporal classification.

- # input timestamp  $\xrightarrow{\text{much}}$  # output timestamp

→ the quick brown fox  
→ btt - h - eee - --- l - --- zzz - -

## Trigger Word Detection

- Alexa, Siri, Google Home

- 0/1 label for y

- It creates un-balanced training set

- few 1's when it triggered