

# Week 1 Recurrent Neural Networks

Thursday, April 5, 2018 9:26 PM

# RNNs

Thursday, April 5, 2018 9:27 PM

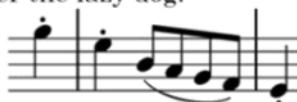
Why Sequence Models?

Speech recognition



→ "The quick brown fox jumped over the lazy dog."

Music generation



Sentiment classification

"There is nothing to like in this movie."



DNA sequence analysis

AGCCCCTGTGAGGAAC TAG

AGCCCCTGTGAGGAAC TAG

Machine translation

Voulez-vous chanter avec moi?

Do you want to sing with me?

Video activity recognition



Running

Name entity recognition

Yesterday, Harry Potter met Hermione Granger.

Yesterday, Harry Potter met Hermione Granger.

Screen clipping taken: 4/5/2018 9:34 PM

Screen clipping taken: 4/19/2018 1:23 PM

Notation

Screen clipping taken: 4/19/2018 1:22 PM

motivating example! Named entity recognition  
x<sup><1></sup> x<sup><2></sup> x<sup><3></sup> x<sup><4></sup> x<sup><5></sup>  
x: Harry Potter and Hermione Granger

y- 1 1 Ø 1 1 Ø Ø Ø Ø  
y<sup><1></sup> y<sup><2></sup> y<sup><3></sup> y y<sup><4></sup> y y<sup><5></sup> y

$x^{(i), t}$  -  $t^{\text{th}}$  value in example  $i$

$T_x^{(i)}$  - length of traing ex.  $i$

How do we represent words?

Vocab

a	1
aaan	2
!	
and	367
:	
harry	4075
:	
etter	6830
:	
Zulu	$10^4$

represent w/ one-hot Recall:  $\begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}_n^m \rightarrow [m]$

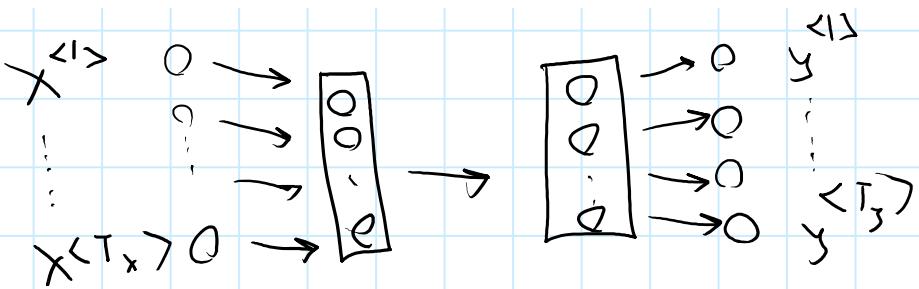
vector for  $x^{(t)}$

Words not in vocab denoted  $\langle \text{unk} \rangle$

Recurrent Neural Network Model

Why not use a standard NN?

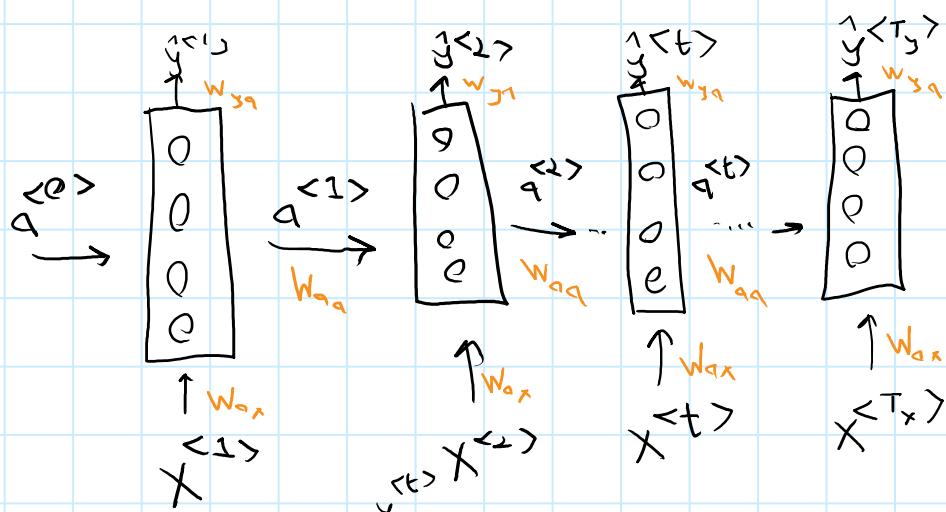




Problems:

- Inputs & outputs can be different lengths
- Doesn't share features learned across diff positions of text

What is an RNN?



thus relation



Can't get info on later words this way, issue case:

"He said, 'Teddy Roosevelt was a great President."

Hc said 'Teddy bears are on sale.'

Solve later with Bidirectional RNN (BRNN)

$$a^{<0>} = \vec{0} ; \quad a^{<1>} = g(W_{aa} a^{<0>} + W_{ax} X^{<1>} + b_a) \leftarrow \text{tanh/RLU}$$

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

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

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

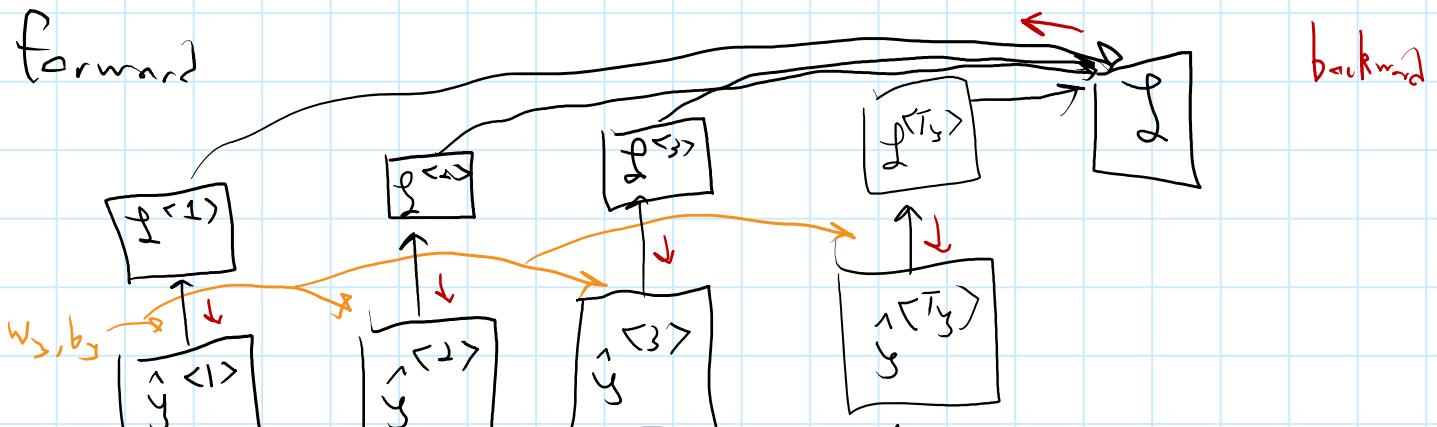
Shorthand notation

$$a^{<t>} = g(W_a [a^{<t-1>} \quad X^{<t>}] + b_a)$$

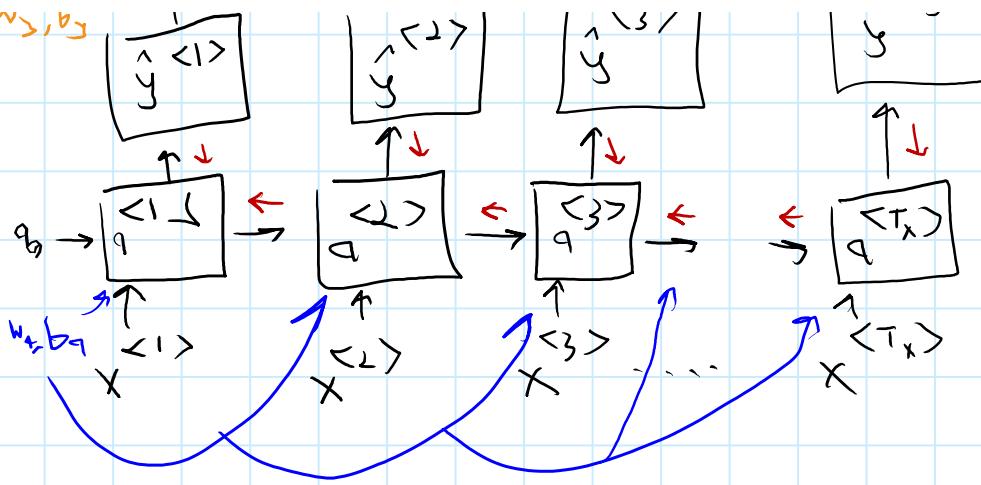
$$y^{<t>} = g(W_y a^{<t>} + b_y)$$

Back propagation through time

forward



backward



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

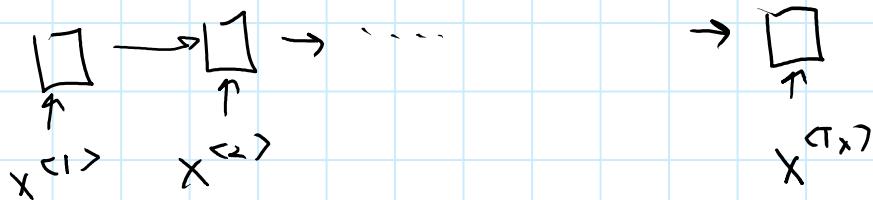
$$\mathcal{L}(y, \hat{y}) = \sum_{t=1}^{T_x} \mathcal{L}^{(t)}(y^{(t)}, \hat{y}^{(t)})$$

Different types of RNNs [Read "The unreasonable effectiveness of RNNs" - Andrej Karpathy]

What if  $x$  &  $y$  are different sizes? (See pic from L1)

Example:  $x = T_x, y = T_y$ , "many-to-many"

Sentiment classification  
 $x = \text{text}, y = 0/1 \dots 5$  "many-to-one"

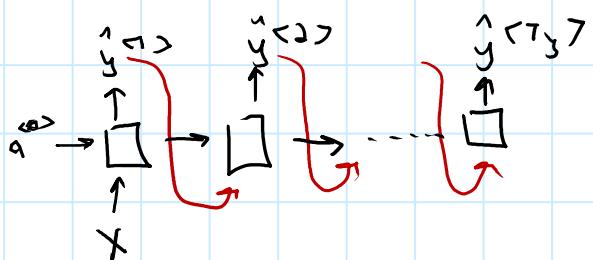


There is nothing to like in this movie

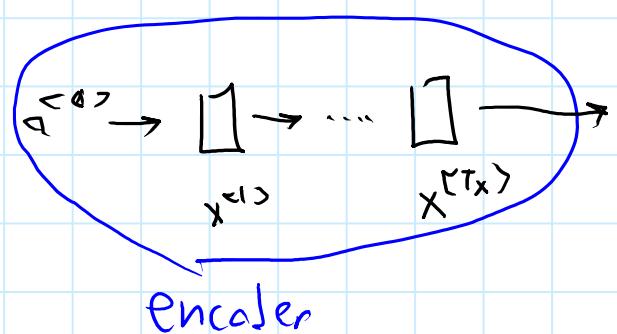
"One-to-one" like a standard NN from first 2 courses

$$x \rightarrow [] \rightarrow \hat{y}$$

Music Generation



Machine translation



"One-to-many"

"many-to-many"

decoder

Language Model & Sequence Generation

Speech recognition

Decide on homonyms

pear / pair

$P(\text{sentence}) = \text{output}$

$P(y^{<1>} | y^{<2>} | \dots | y^{<T_y>})$

Training set! Large corpus of English text

"Tokenize" sentence

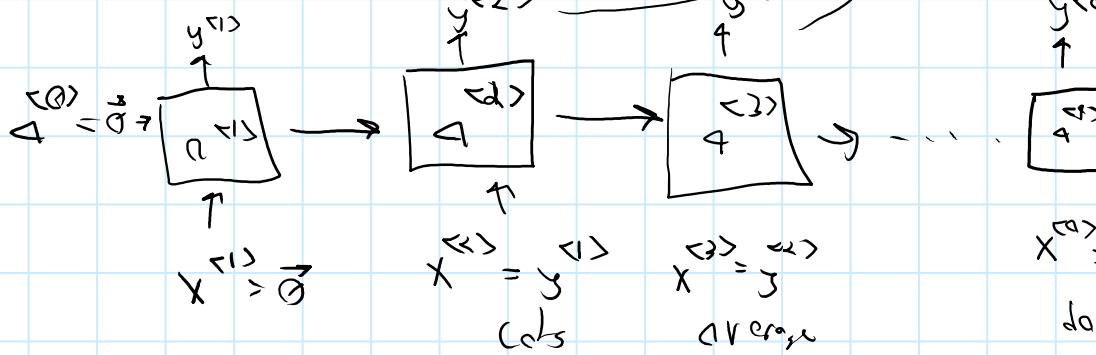
End of Sentence token

Cats average 15 hours of sleep a day. <EOS>  
 $y^{<1>} \quad y^{<2>} \quad \dots \quad y^{<8>} \quad y^{<9>} \quad y^{<10>}$

$p(\text{average} | \text{cats})$

RNN  $p(y^{<1>} = [\text{first}, \text{drat}, \dots])$

$p(\text{--- "cats average"})$



$$\mathcal{L}(y^{<t>}, \hat{y}^{<t>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>} ; \mathcal{L} = \sum_t \mathcal{L}(y^{<t>}, \hat{y}^{<t>})$$

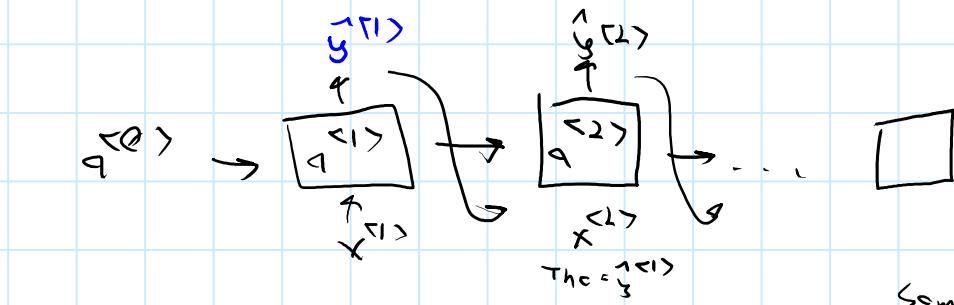
$$P(y^{<1>}, y^{<2>}, y^{<3>}) = P(y^{<1>}) P(y^{<2>} | y^{<1>}) P(y^{<3>} | y^{<1>}, y^{<2>})$$

Sampling More Sequences

Good way to check an trained networks

Sample a sequence from a trained RNN, gen sentence

Sample a sequence from a trained RNN, gen sentence



Sample until  $\langle \text{EOS} \rangle$  or word limit

$$y^{(1)} = \text{np.random.choice}([P(a_1), P(a_2), \dots])$$

Can also do this on character level

pro's: no  $\langle \text{UNK} \rangle$  (good for specialized vocabularies)

Con: much longer sequences, computationally expensive



## Vanishing Gradients With RNNs

Ex The cat, which already the ... , was full  
" cat" were ..

Basic RNN not good detecting long term dependencies,  
this is due to vanishing gradients

Values are mainly affected by other values near them

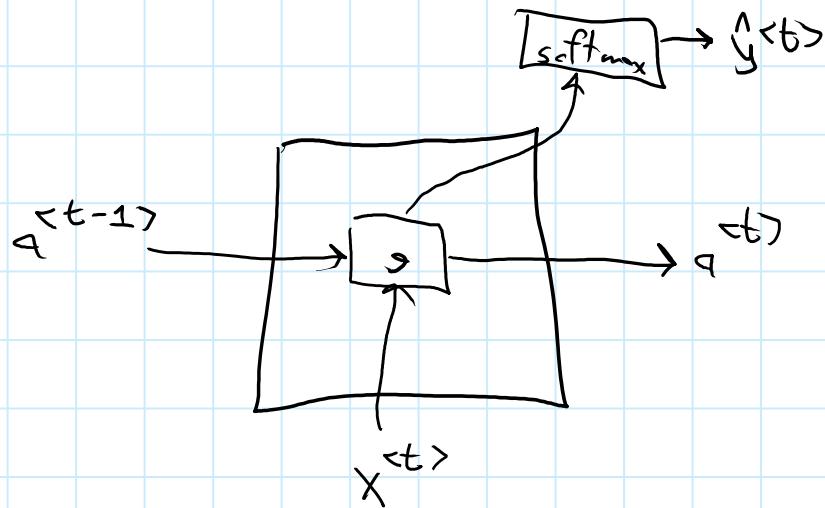
What if you got exploding gradients?

Use gradient clipping (rescale gradient values over threshold)

## Gated Recurrent Unit (GRU)

RNN:

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



## GRU (Simplified)

$c = \text{memory cell}$   
 $c^{<t>} = a^{<t>}$

The  $\boxed{\text{cat}}$  ....  $\boxed{\text{was}}$  full

At every time step, consider overwriting  $c^{<t>}$  w/  $\tilde{c}^{<t>}$  thol ( $\tilde{c}^{<t>}$ )

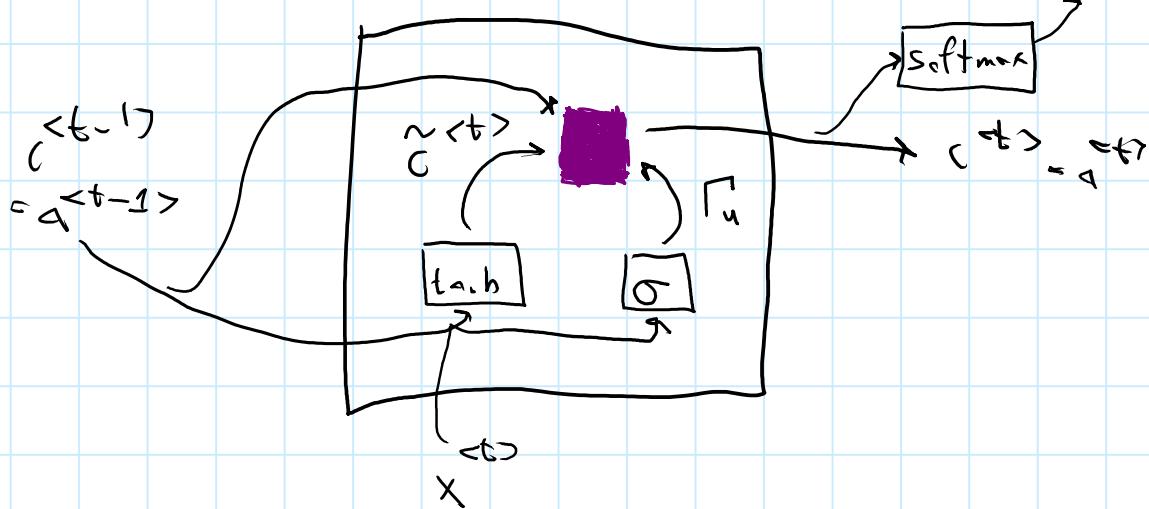
$$\tilde{c}^{<t>} = \tanh(W_c [c^{<t-1>}, x^{<t>}] + b_c)$$

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

in this example,  $c^{<t>} \leftarrow \begin{cases} 0 & \text{plural} \\ 1 & \text{singular} \end{cases}$

$$c^{<t>} = r_u * \hat{c}^{<t>} + (1 - r_u) * c^{<t-1>}$$

If this is 1, set  $\hat{c}^{<t>} \leftarrow \hat{c}^{<t>}$



## Full GRU

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

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

$\Rightarrow$   
 Long Short Term Memory (LSTM)

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

$$r = \sigma(a^{<t-1>} * c^{<t>} ..)$$

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

"forget"

$$\Gamma_f = \sigma(W_f [a^{<t-1>}, x^{<t>} + b_f]) \approx (I - \Gamma_u)$$

"output"

$$\Gamma_o = \sigma(W_o [a^{<t-1>}, x^{<t>} + b_o])$$

$$c^{<t>} = \Gamma_u * c^{<t>} + \Gamma_f * c^{<t-1>} ; \text{ can add old + new value}$$

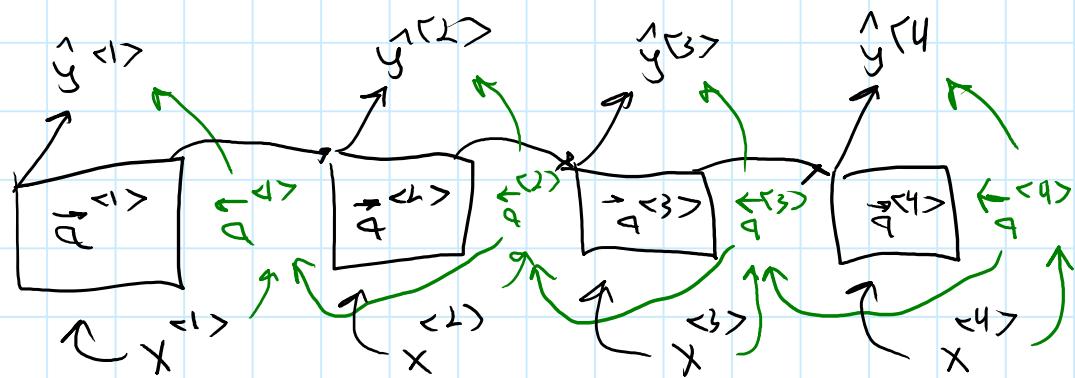
$$a^{<t>} = \Gamma_o * c^{<t>}$$

"recurrent connection"  $\sigma(W_- [a^{<t-1>}, c^{<t-1>}] + b_-)$

---

Bidirectional RNNs

Recall Teddy sentence, 3sec



Still forward propagation! Acyclic graph

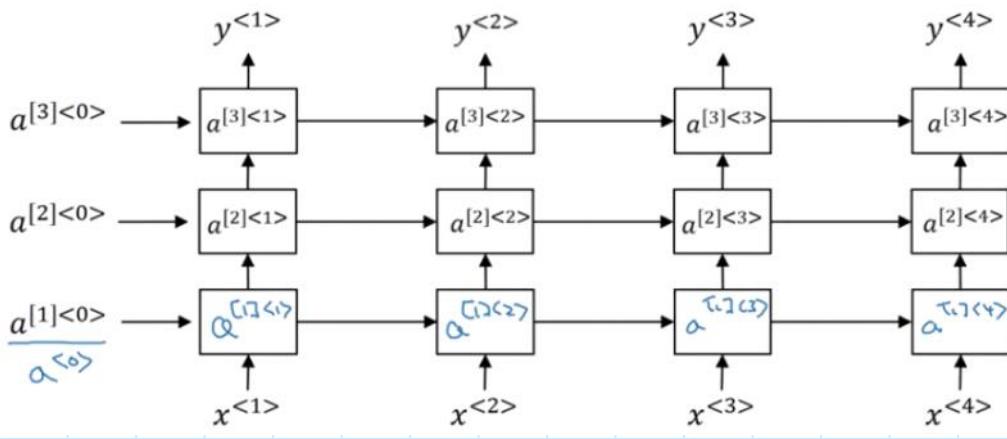
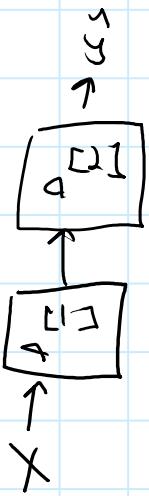
→ input-output → propagation, time step by step

$$\hat{y}^{(t)} = g\left(W_y \left[ \vec{a}^{(t)}, \vec{a}^{(t)} \right] + b_y\right)$$

Disadvantage: need whole sentence before computing  
(e.g., bad for real time speech processing)



## Deep RNNs



$$a^{[2]^{[3]}} = g\left(W_a^{[2]} \begin{bmatrix} a^{[2]^{[2]}} \\ a^{[1]^{[3]}} \end{bmatrix} + b_a^{[2]}\right)$$

# Week 2 NLP and Word Embeddings

Tuesday, April 24, 2018 10:53 AM

# Intro to Word Embeddings

Tuesday, April 24, 2018 10:53 AM

Helps understand analogies

Screen clipping taken: 4/24/2018 3:30 PM

## Word Representation

reviewly:

$$V = [a, aaron, \dots, zulu, \text{UNK}] \quad \text{Onehot vectors}$$

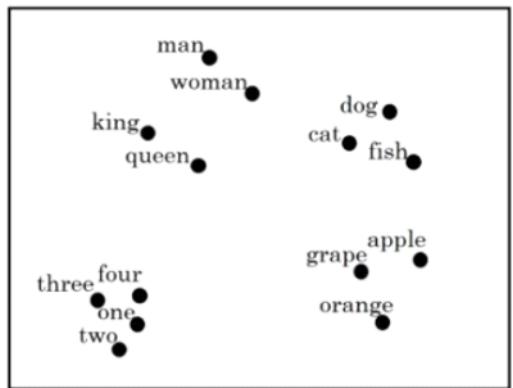
Downside! Each word self-contained, can't generalize across words

c.g. I want a glass of orange juice  
if  $\vec{r}^1 \vec{r}^2 \vec{r}^3 \vec{r}^4 \vec{r}^5$  apple  $\vec{\_}$ ? doesn't recognize

This arises b/c inner product between one-hot vectors is 0  
Euclidean distance always the same

Featureized		Representation					
index	(5391)	(9853)	(4914)	(7157)	(456)	(6257)	
word	Man	Woman	King	Queen	Apple	Orange	
Gender	-1	1	-0.95	0.97	0.00	0.01	
Royal	0.01	0.02	0.95	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	
;	↑	↑					
(simplification)	$e_{5391}$	$e_{9853}$					

~ ~



300D  
↓  
2D  
(PCA)

t-SNE

~~Using Word Embeddings~~

Use example

Sally Johnson is an orange farmer

Robert Lin is a durian cultivator

- 1) Train on 1Bn-100Bn words  
(or download pre-trained embedding online)
- 2) Use embedding for new task w/ smaller set
- 3) (Optional) Continue to fine tune word embeddings (if 2 hrs large set)

Relation to Face Encoding

$$f(x^{(i)}) \approx e_i$$

# Properties of Word Embeddings

## Analogies

index	(5311)	(9853)	(4914)	(7157)	(456)	(6257)
word	Man	Woman	King	Queen	Apple	Orange
Gender	-1	1	-0.95	0.97	0.00	0.01
Race	0.01	0.02	0.95	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
	$e_{\text{man}}$	$e_{\text{woman}}$	$e_{\text{King}}$	$e_{\text{Queen}}$		

Can I learn man is to woman as king is to queen?

$$e_{\text{man}} - e_{\text{woman}} = \begin{bmatrix} -2 \\ 0 \\ 0 \end{bmatrix}$$

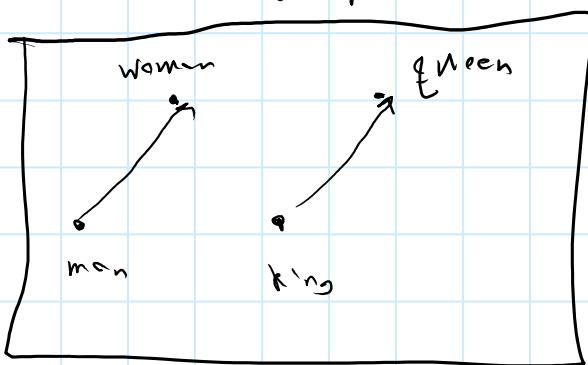
$$e_{\text{King}} - e_{\text{Queen}} = \begin{bmatrix} -2 \\ 0 \\ 0 \end{bmatrix}$$

[Algorithm!  $e_{\text{man}} - e_{\text{woman}} \approx e_{\text{King}} - e_{\text{Queen}}$ ]

Dr. N>: Find word in w:

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

in 300D space



(won't hold after t-SNE mapping)

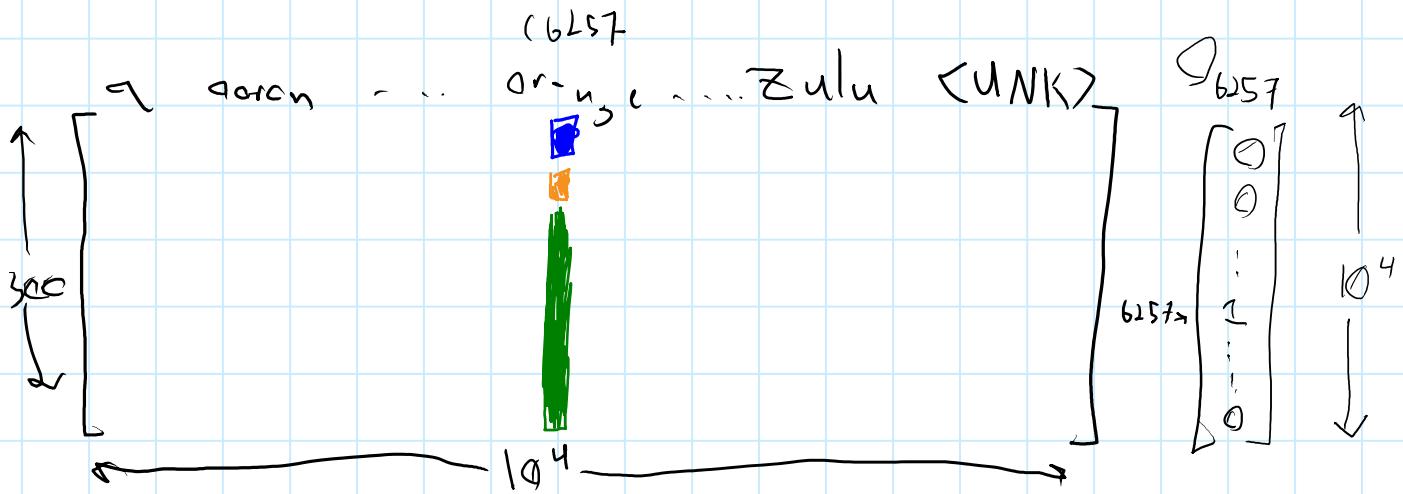
Cosine Similarity (inner product)

$$\text{Sim}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u}^T \mathbf{v}}{\|\mathbf{u}\|_2 \|\mathbf{v}\|_2} \equiv \cos \phi$$

[Could also use - Euclidean distance (dissimilarity)  $\|\mathbf{u} - \mathbf{v}\|^2$ ]



Embedding Matrix



$$E \cdot O_{6257} = \begin{bmatrix} \text{blue} \\ \text{orange} \\ \text{green} \end{bmatrix} \stackrel{300}{\uparrow} \stackrel{6257}{\downarrow} = e_{6257}$$

$$E \cdot O_j = e_j$$

In practice, this isn't very efficient, use specialized function

# Learning Word Embeddings

Wednesday, April 25, 2018 11:02 AM

Algorithms have gotten simpler over time

## Neural Language Model

I O<sub>4343</sub> → E → e<sub>4343</sub>

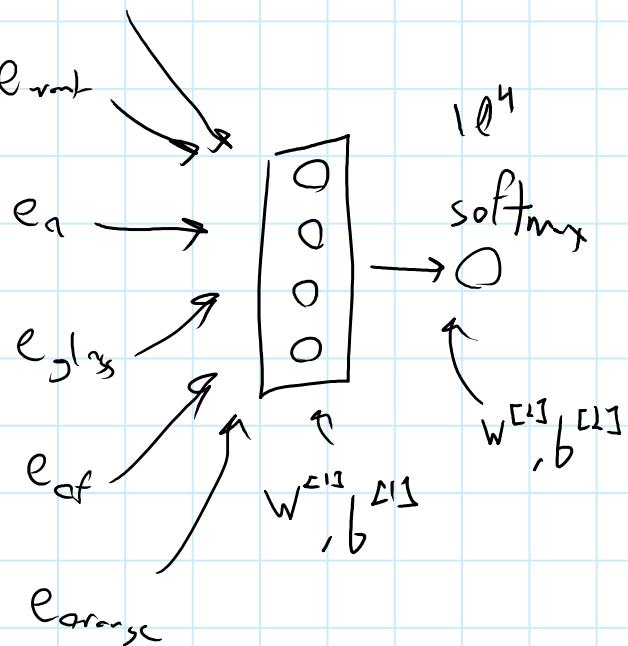
want Q<sub>want</sub> → E → e<sub>want</sub>

a

glass

of

orange



often use fixed history (4 words prior to prediction)

Other context/target pairs

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

Context target

Last 4 words

4 words on left + right

Last one word

orange ?

{ M-1 7 1-1 7 1 7 2

{ last one word  
 Nearby one word (skip gram)      orange ?  
 loss ?  
 Good for learning embeddings

Word  $\rightarrow$  Vec

Skip-gram

Create context-target pairs for supervised learning

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

↑      ↑      ↑

Context

orange  
orange  
orange

Target

juice  
glass  
my

Model

Vocab size = 10,000k

Context  $c$  ("orange")  $\rightarrow$  Target  $t$  ("juice")  
6257                          4834

$O_c \rightarrow E \rightarrow e_c \rightarrow \text{softmax} \rightarrow y$

$(f_t, \dots, f_{t-1}) - \frac{e^{e_t^T e_c}}{\text{norm}}$   $\Theta_t$  = parameter associated w/ output  $t$

$$\text{Soft max: } p(t|c) = \frac{e^{\Theta_t^T e_c}}{\sum_{i=1}^{10,000} e^{\Theta_i^T e_c}}$$

$\Theta_t$  = parameter associated w/ output  $t$

$$\text{Loss } L(\vec{y}, \vec{y}) = - \sum_{i=1}^{10,000} y_i \log(\vec{y}_i)$$

one-hot  $\vec{y}$  (10k dim vector output by softmax)

## Problems

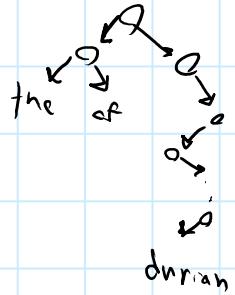
computational speed, vocab size due to summing

Solution:

heuristic soft max  
(split dataset)



in practice, not symmetric  
(common words near top)



How to sample context  $c$ ?

Don't want commonly occurring words dominating (the, of, a, and, etc.)  
Sampling not truly random

→

Negative Sampling

Make skip-gram more efficient

Given a pair, are they target-context pairs?

First row generated by neighbor sampling, next random dictionary words

<u>Context</u>	<u>Word</u>	<u>Target?</u>
orange	juice	1
orange	king	0
r	book	0
s	the	0
"	of	0

random dictionary words, labelled 0

5-10 for small data sets

2-5 for larger datasets

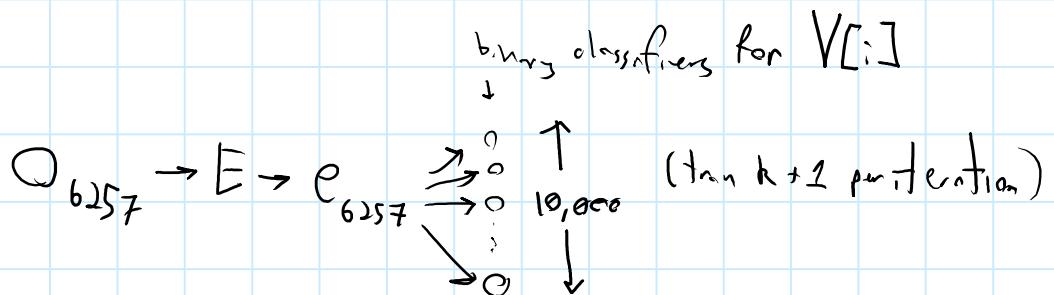
Model

```

    target      context
    ↓           ↓   ↗ word
    p(y=1 | c, t) = σ(θtTec)
  
```

NN

orange



Selecting Negative Examples

Best from paper!  $P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$

Fin between uniform & representative

# GloVe Word Vectors

Used less than word2vec + skip-gram

Global Vectors for Word representation

$$X_{ij} = \# \text{ times } i \text{ appears in context of } j$$

If context is symmetric,  $X_{ij} = X_{ji}$

$$\text{minimize}_{\Theta, b} \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\Theta_i^T e_j + b_i + b_j - l_0(X_{ij}))^2$$

$$\begin{aligned} \text{weight}_i &= \Theta_i \cdot \mathbf{f}(X_{ij}) = \Theta_i \\ \Theta l_0(\Theta) &= 0 \text{ in the context} \\ \text{this is of a} \end{aligned}$$

$\Theta$  &  $e_j$  are symmetric

$$e_w^{(f^{-1})} = \frac{e_w + \Theta_w}{2}$$

A note on the factorization view of a word

Cant guarantee individual components are interpretable

$$\underbrace{\Theta_i^T e_j}_{\text{[A is invertible]}}$$

$$(A \Theta_i)^T (A^{-T} e_j) = \Theta_i^T A^T \cancel{A^{-T}} e_j$$

axes may not be well aligned

# Applications Using Word Embeddings

Wednesday, April 25, 2018 2:49 PM

## Sentiment Classification

Sentence about restaurant  $\rightarrow$  star rating

Training sets  $<10k$  words  $\rightarrow 100k$  words

Ex:

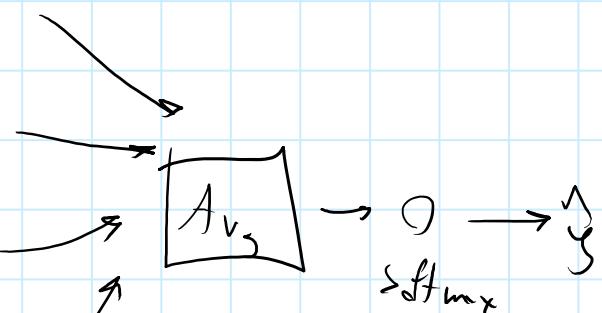
The  $\circ \rightarrow E \rightarrow e_{\text{The}}$

dessert  $\circ \rightarrow E \rightarrow e_{\text{dessert}}$

is  $\circ \rightarrow E \rightarrow e_{\text{is}}$

excellent  $\circ \xrightarrow{\text{q}} E \rightarrow e_{\text{excellent}}$

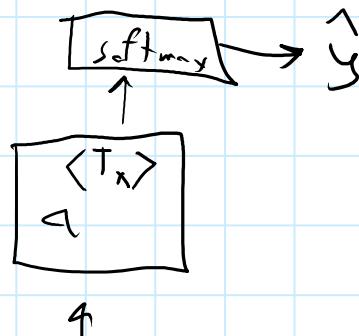
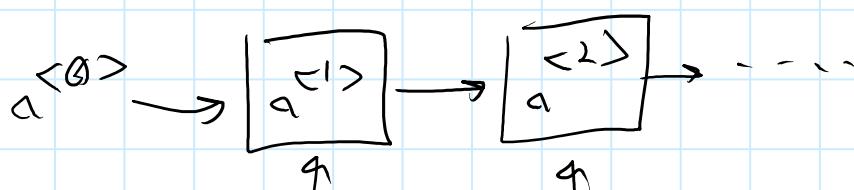
$100k$  words



Issue: ignores order

"Lacking in good service" marked high

Solution: RNN



$$a \rightarrow \boxed{a} \quad \boxed{a} \quad \boxed{\phantom{a}}$$

↓      ↓      ↓

e      e      e



## Debiasing Word Embeddings

(Discriminatory in the ethical, not statistical sense)

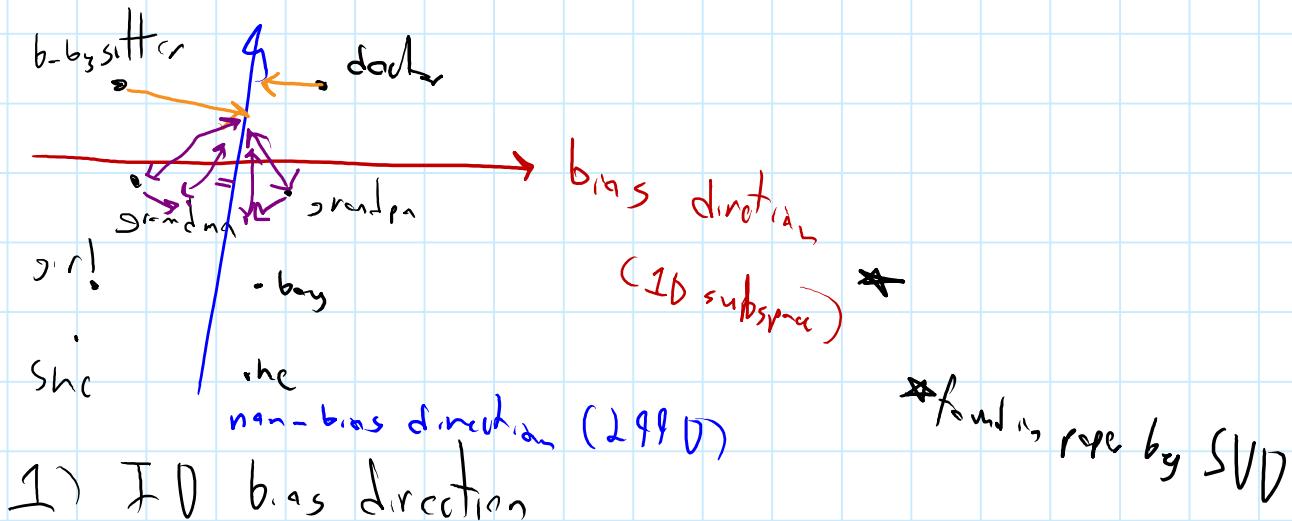
Man : Woman as King : Queen

[Bolukbasi 2016]

Man : Programmer as Woman : Homemaker

Father : Doctor as Mother : Nurse

## Addressing bias in word embeddings



$$\left\{ \begin{array}{l} e_{he} - e_{she} \\ e_{male} - e_{female} \end{array} \right.$$

or,

( ( ! )  
around  
average

2) Neutralize!

Project every non-definition word to eliminate bias

(  
Grandma, grandpa would be defined)

3) Equalize pairs

grandma → grandpa

boy ← girl

# Week 3 Sequence Models and Attention Mechanism

Thursday, April 26, 2018 1:22 PM

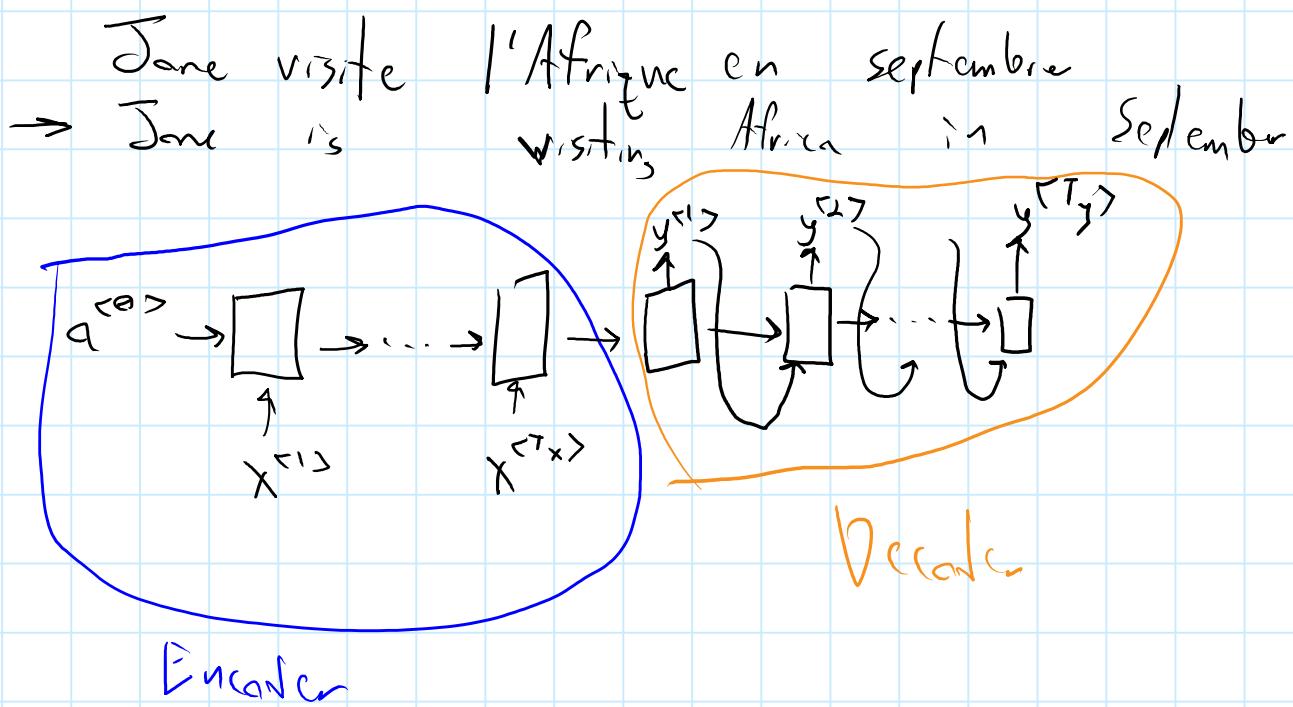
# Various Sequence to Sequence Architectures

Thursday, April 26, 2018

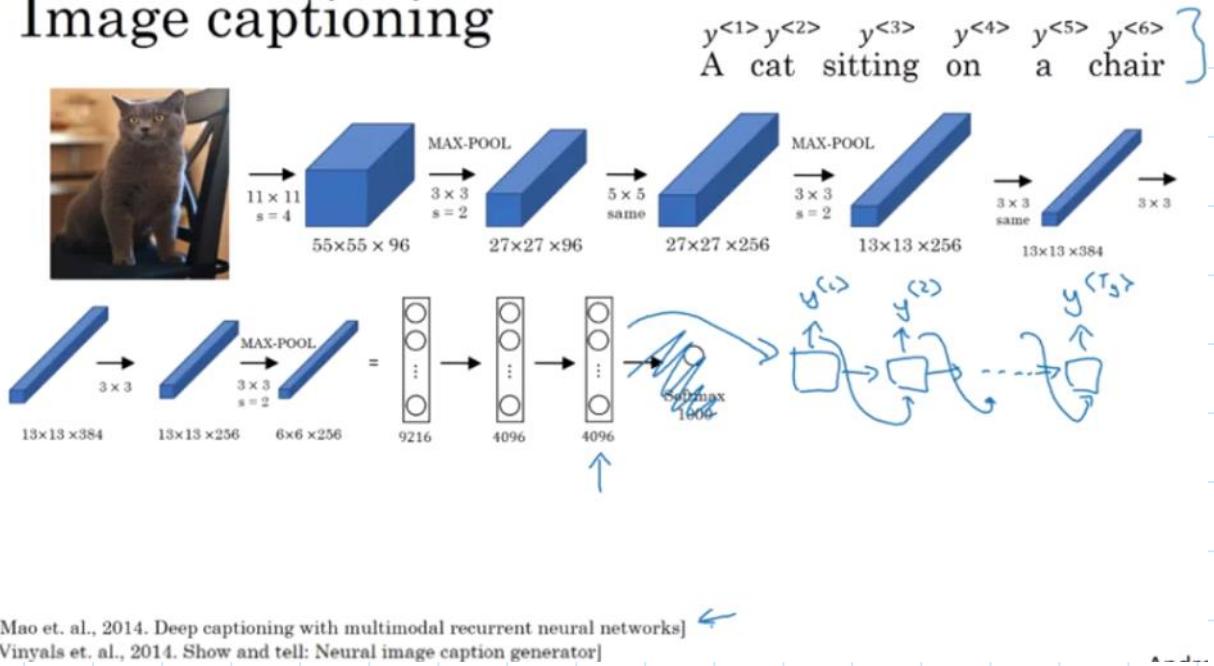
1:24 PM

## Basic Models

### Sequence-to-sequence model



# Image captioning



Screen clipping taken: 4/27/2018 8:10 AM

Picking the Most Likely Sentence

Decoder network resembles Language Model

Machine translation "conditional language model"

Where Language Model computes  $P(y^{<1>} y^{<2>} | \dots, y^{<T_y>})$

Machine Translation computes  $P(y^{<1>} \dots, y^{<T_y>} | x^{<1>} \dots, x^{<T_x>})$

$P(y^{<1>} \dots, y^{<T_y>} | x^{<1>} \dots, x^{<T_x>})$

French  
↑ English

Do not want random sampling!  $\arg \max_{y^{(1)}, \dots, y^{(n)}} P(y^{(1)}, \dots, y^{(n)} | X)$

Why not use greedy search? [most likely from position  $y^{(1)}$ ]

Jane is  $\begin{cases} \text{going} & \rightarrow \text{better } P(y^{(2)} | X) \\ \text{visiting} & \rightarrow \text{better sequence} \end{cases}$

Also computationally inefficient, better to use approximate

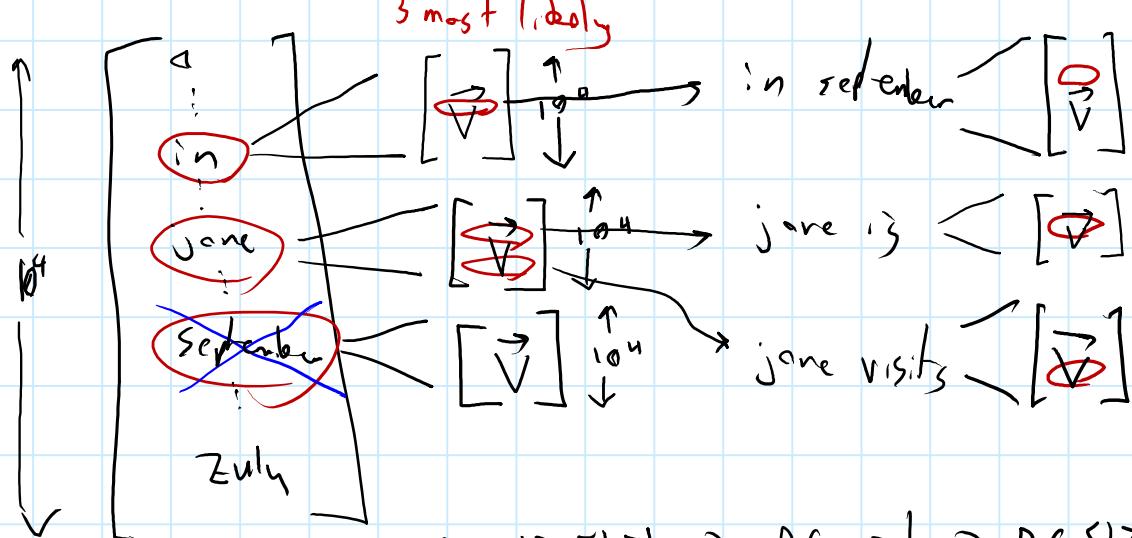
## Beam Search

Parameter  $B = \text{beam width}$

Picks  $B$  most likely possibilities for  $P(y^{(1)} | X)$   
keeps in memory

ex  $B=3$

Step 1 Step 2 Step 3



$$P(y^{(1)}, y^{(2)} | X) = P(y^{(1)} | X) P(y^{(2)} | X, y^{(1)})$$

$$V \leftarrow P(y^{<1>} | x) = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>})$$

Every step run model 3 copies of network

~~Refinement to Beam Search~~

Length Normalization

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

Soln:  $\arg \max_y \sum_{t=1}^{T_y} \log(P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>}))$

Small enough to run into underflow

[Same value will maximize both since  $\log$  monotonically increases]

Problem: Shorter sentences are preferred by nature of algorithm

Soln:  $\frac{1}{T_y} \alpha \sum_{t=1}^{T_y} \log(P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>}))$

How to choose  $B$ ?

large: better result, slower, more memory

small: worse, faster, less memory  
prod systems:  $10-100$ ; research:  $10^3 - 3 \times 10^3$

Graphs  $\approx$  logarithmic



## Error Analysis in Beam Search

Beam search doesn't **ALWAYS** find best result,  
how do we decide if it's an RNN issue or Beam Search?

Ex: Algorithm applies past tense erroneously ( $\hat{y}$ )

Plug in part of human translated ( $y^*$ )

Case 1:  $P(y^4|x) > P(\hat{y}|x)$

Beam search chose  $\hat{y}$ , tune BeamSearch

Case 2:  $P(y^4|x) \leq P(\hat{y}|x)$

Likely RNN at fault (maybe time length in BS)

Try multiple  $(\hat{y}, y^*)$  pairs & determine BS: RNN fail rate

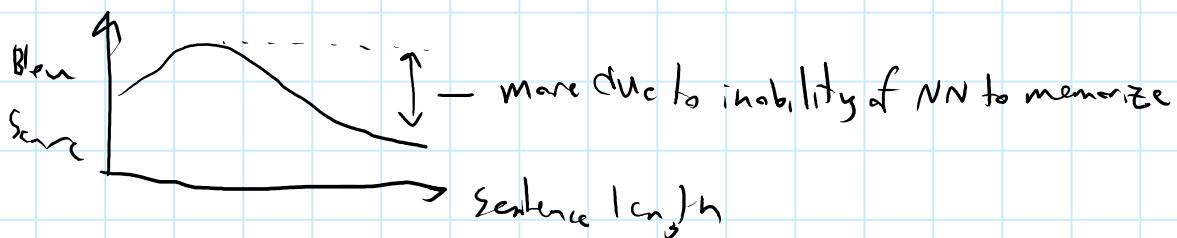


Bleu Score (optional) put later, as my interest in translation isn't THAT high

## Attention Model Intuition

Issue with long sequences!

People do not memorize an entire sentence & then decode it



Bidirectional RNN on sentence  $s^{<0>} \rightarrow \rightarrow \rightarrow$

another RNN  $s^{<0>} \rightarrow \rightarrow \rightarrow$

$s^{<i>}$  should weight  $x^{<i>}$  more heavily than  $x^{<i+1, \dots, c>}$

Use weights  $\alpha$ ,  $\alpha^{<i, t>}$

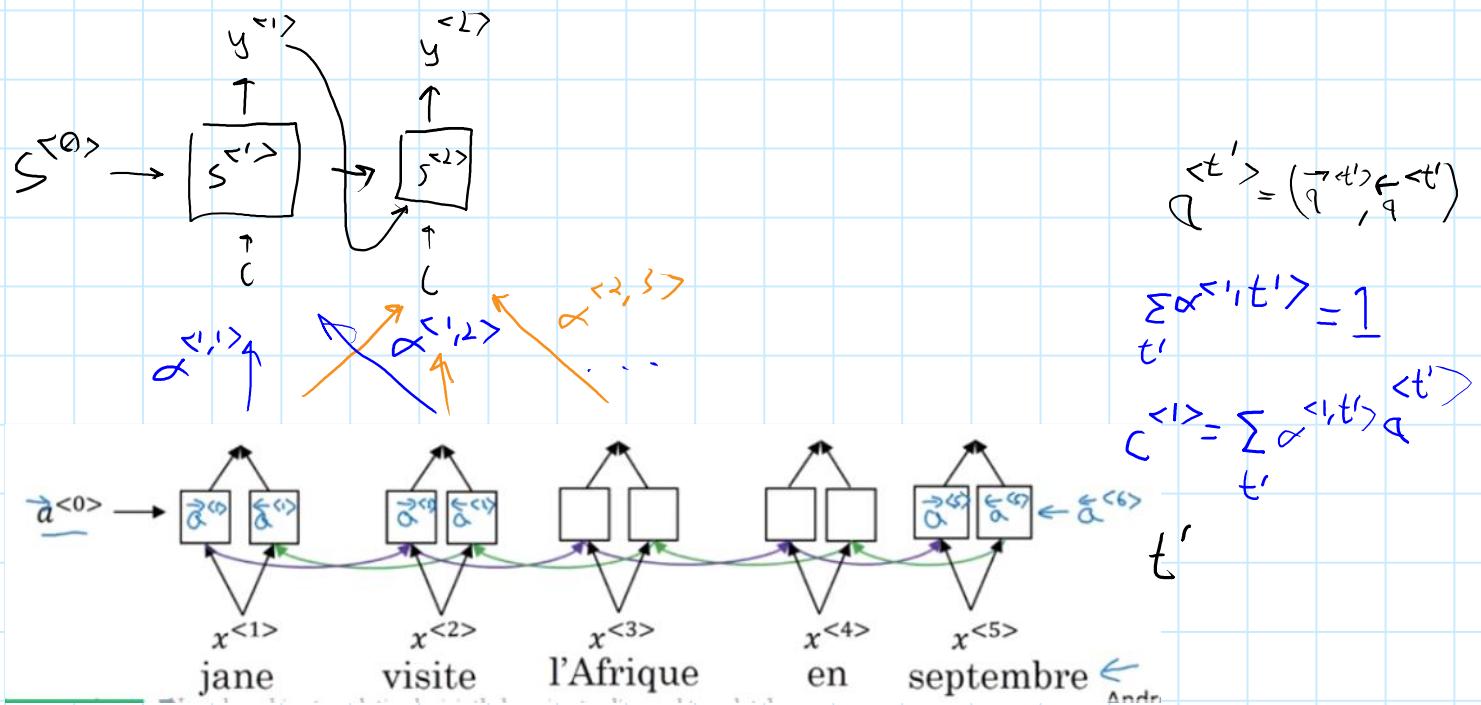
$$s^i \rightsquigarrow x$$

$\alpha^{<i, t>}$  will depend on  $s^{<i>}$ ,  $\overrightarrow{a}^{<t>}, \overleftarrow{a}^{<t>}$

## Attention Model

# Attention Model

(Typically GRU/LSTM)

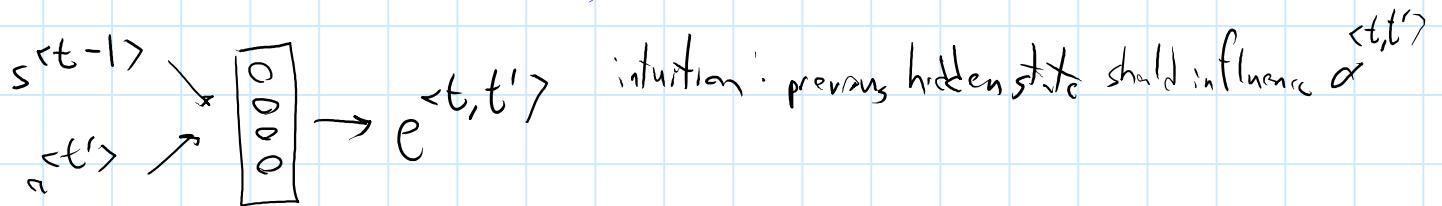


$$\alpha^{<t,t'} = (\vec{q}^{(t)}, \vec{c}^{(t')})$$

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

$$\vec{c}^{<1>} = \sum_{t'} \alpha^{<1,t'} \vec{c}^{(t')}$$

$\alpha^{<t,t'>}$  is amt of attention y<t> should pay to a<t'>



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

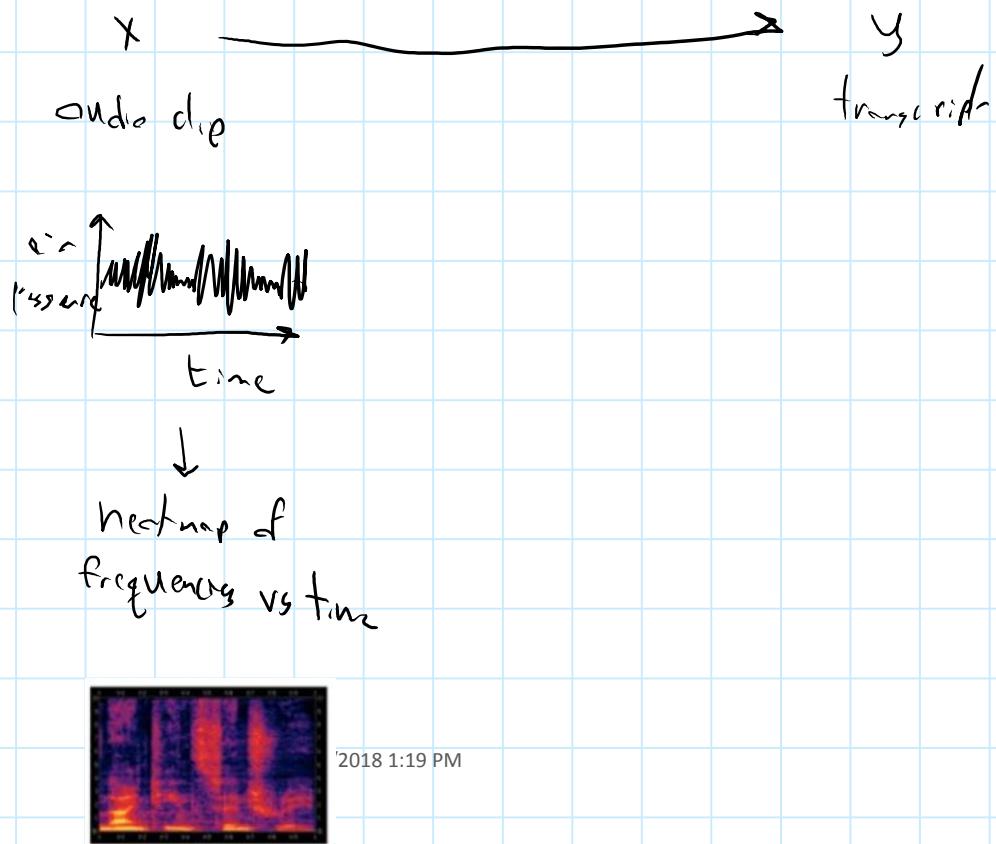
ensures α components  $\sum$  to 1

Downside: Quadratic cost, but sentences are short so typically negligible

[This can also be used for image captioning]

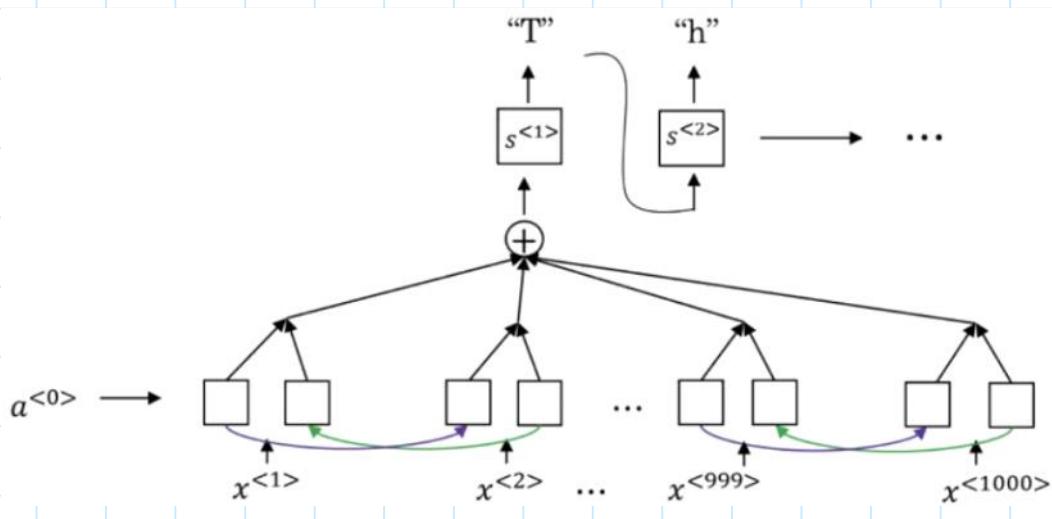
# Speech Recognition

Friday, April 27, 2018 8:23 AM



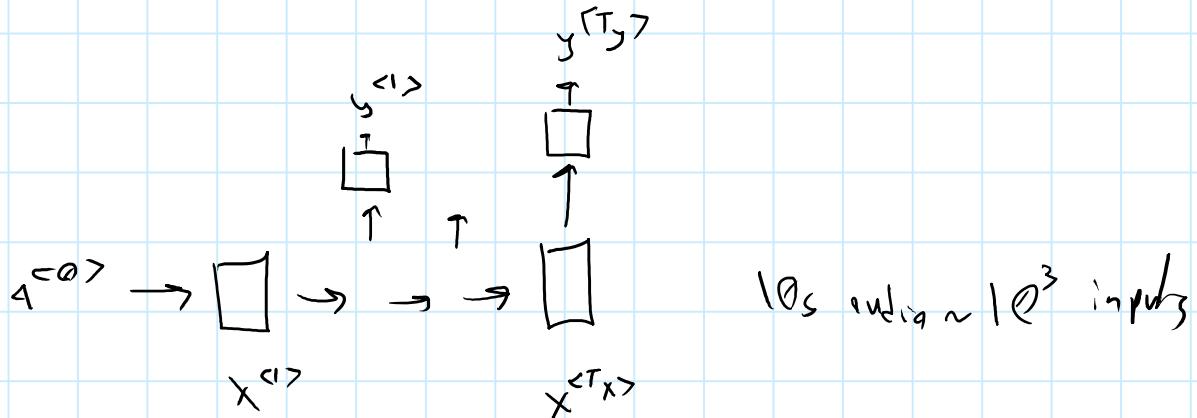
Screen clipping taken: 4/27/2018 1:20 PM

Attention model possible



(CTC) cost for speech recognition  
(connectionist temporal classification)

Bi-RNN typically used



"the quick brown fox" 19 characters from  $10^3$  inputs

g t t t t t t - h - e e e e \_ \_ \_ L \_ \_ \_ q \_ \_  
(collapse repeated characters not separated by "blank" (-))

## Trigger Word Detection

Algorithm not yet set in stone  
↓ from ratio of 0's to 1's

