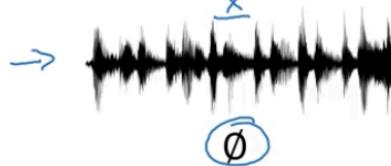


Sequence Models

Examples

Examples of sequence data

Speech recognition



Music generation



Sentiment classification

"There is nothing to like
in this movie."

DNA sequence analysis → AGCCCCCTGTGAGGAACTAG

Machine translation

Voulez-vous chanter avec
moi?

Video activity recognition



Name entity recognition

Yesterday, Harry Potter
met Hermione Granger.

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



AGCCCCTGTGAGGAACTAG

Do you want to sing with
me?

Running

Yesterday, Harry Potter
met Hermione Granger.
Andrew Ng

in some cases only x/y is a sequence

Sequence Model Example: Named Entity Recognition (NER)

Input (x): "Harry Potter and Hermione Granger invented a new spell."

Output (y): 1 1 0 1 1 0 0 0 0 (Indicates presence/absence of a named entity)

Notation:

- $x^{<t>}$: Word at position t in input.
- $y^{<t>}$: Output label (0 or 1) at position t .
- T_x : Input sequence length (here, 9).
- T_y : Output sequence length (here, 9).
- $x^{(i)}{<t>}$: Word at position t in the i -th input sequence.
- $y^{(i)}{<t>}$: Output label at position t of the i -th sequence.
- $T_x^{(i)}$: Length of the i -th input sequence.

Task: NER - Identify entities (names, dates, etc.) in text.

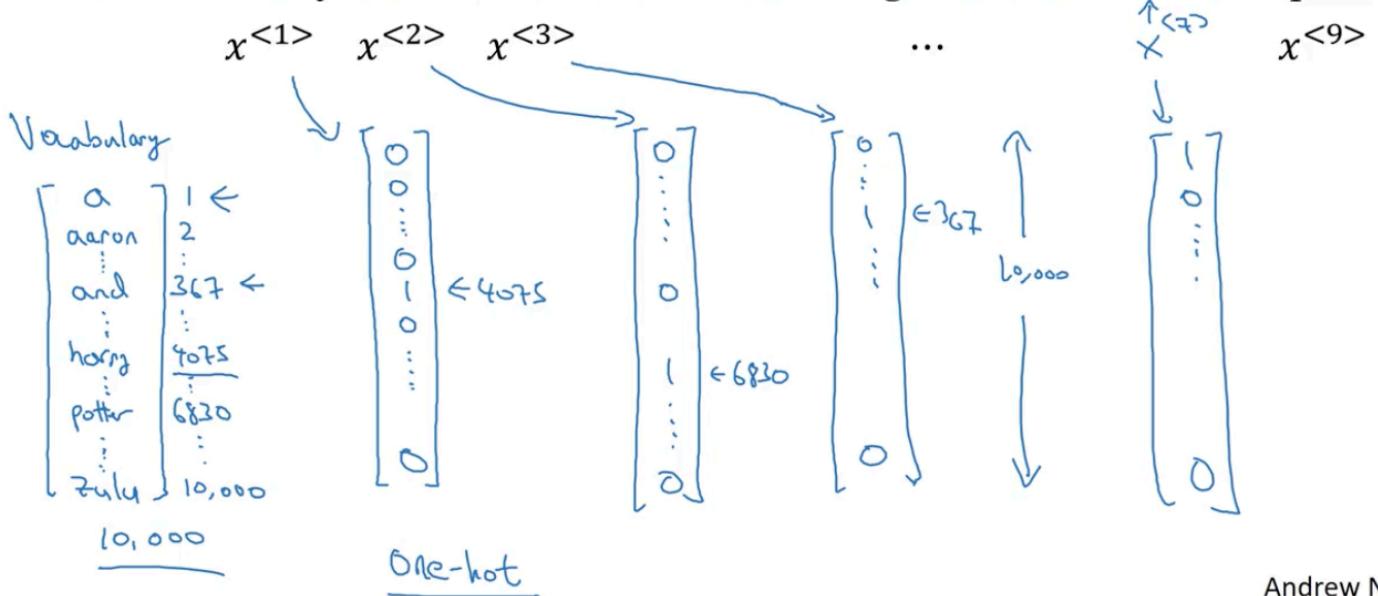
Note: Simple word-by-word classification is insufficient; context is crucial (e.g., "Harry Potter" vs. "Potter"). Better NER uses B- (beginning) and I- (inside) tags.

dictionary sizes: at least around (for commercial) 30k-50k, bigger uses/models vocabulary: millions.

Representing words

$$x^{<\leftrightarrow>} \quad x \rightarrow y$$

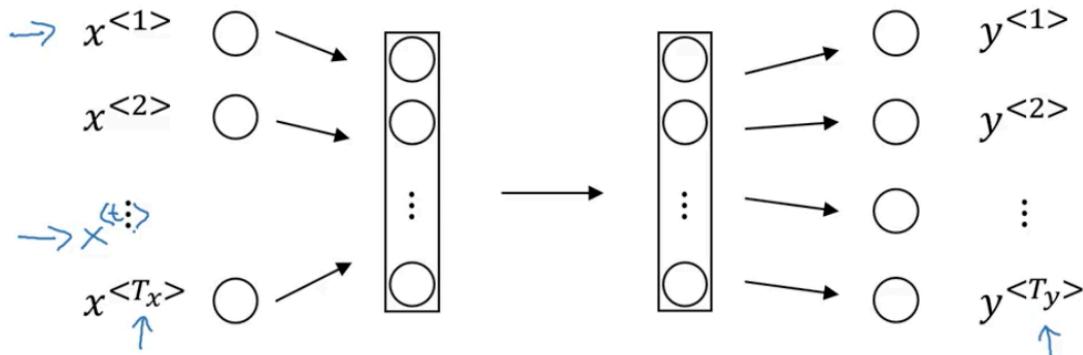
x : Harry Potter and Hermione Granger invented a new spell.



Andrew Ng

Recurrent Neural Network Model

Why not a standard network?



Problems:

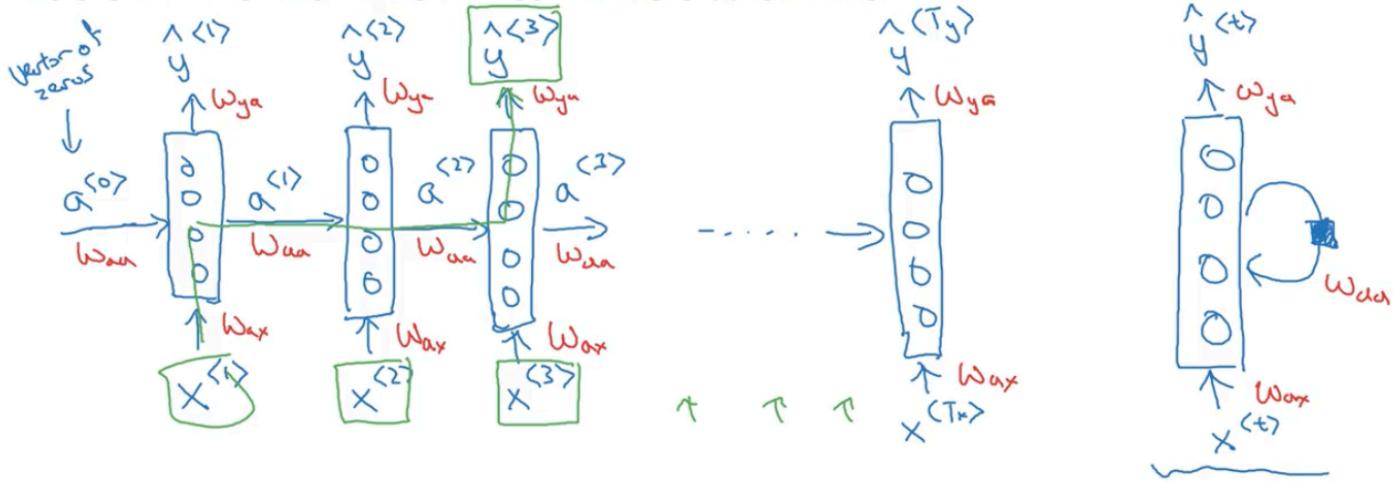
- Inputs, outputs can be different lengths in different examples.
- Doesn't share features learned across different positions of text.

Andrew Ng

10k dimensional one-hot vector

- enormous amount of parameters

Recurrent Neural Networks



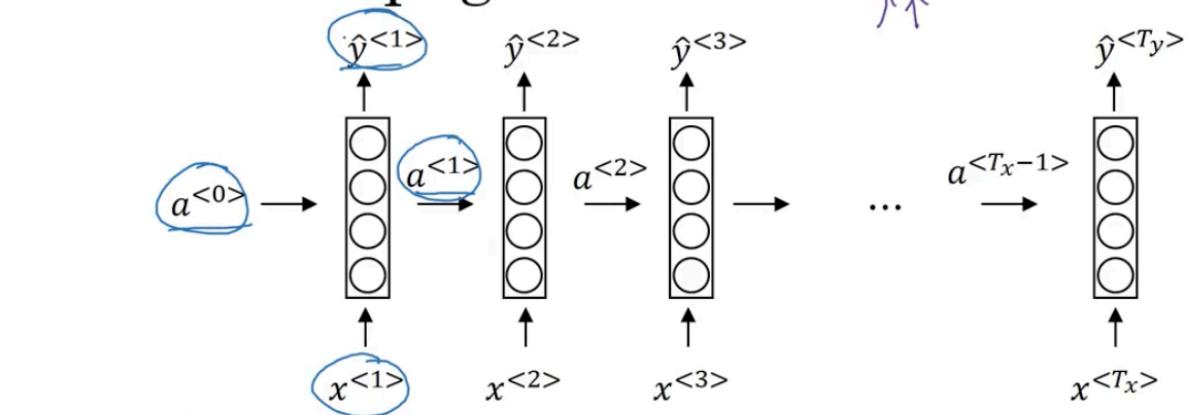
He said, "Teddy Roosevelt was a great President."

He said, "Teddy bears are on sale!"

Andrew Ng

uni directional recurrent neural network

Forward Propagation



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

$$a^{<t>} = g_1(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a) \leftarrow \text{tanh / ReLU}$$

$$\hat{y}^{<t>} = g_2(W_{ya} a^{<t>} + b_y) \leftarrow \text{sigmoid}$$

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

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

Andrew Ng

activation function, depends on what type of task - binary classification: sigmoid, though there are others for other tasks

Simplified RNN Notation

Simplified RNN notation

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

$\uparrow \quad \quad \quad \uparrow$
 $(100, 100) \quad (100, 10,000)$

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

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

\uparrow
 $[W_{aa}; W_{ax}]$

$= W_a$
 $(100, 10100)$

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

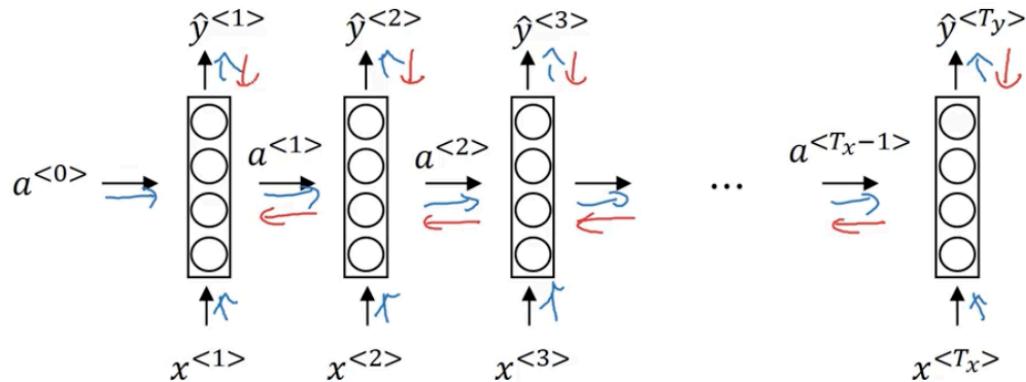
$\uparrow \quad \quad \quad \uparrow$
 $100 \quad 10000$

$$[W_{aa}; W_{ax}] \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = W_{aa}a^{<t-1>} + W_{ax}x^{<t>}$$

10100

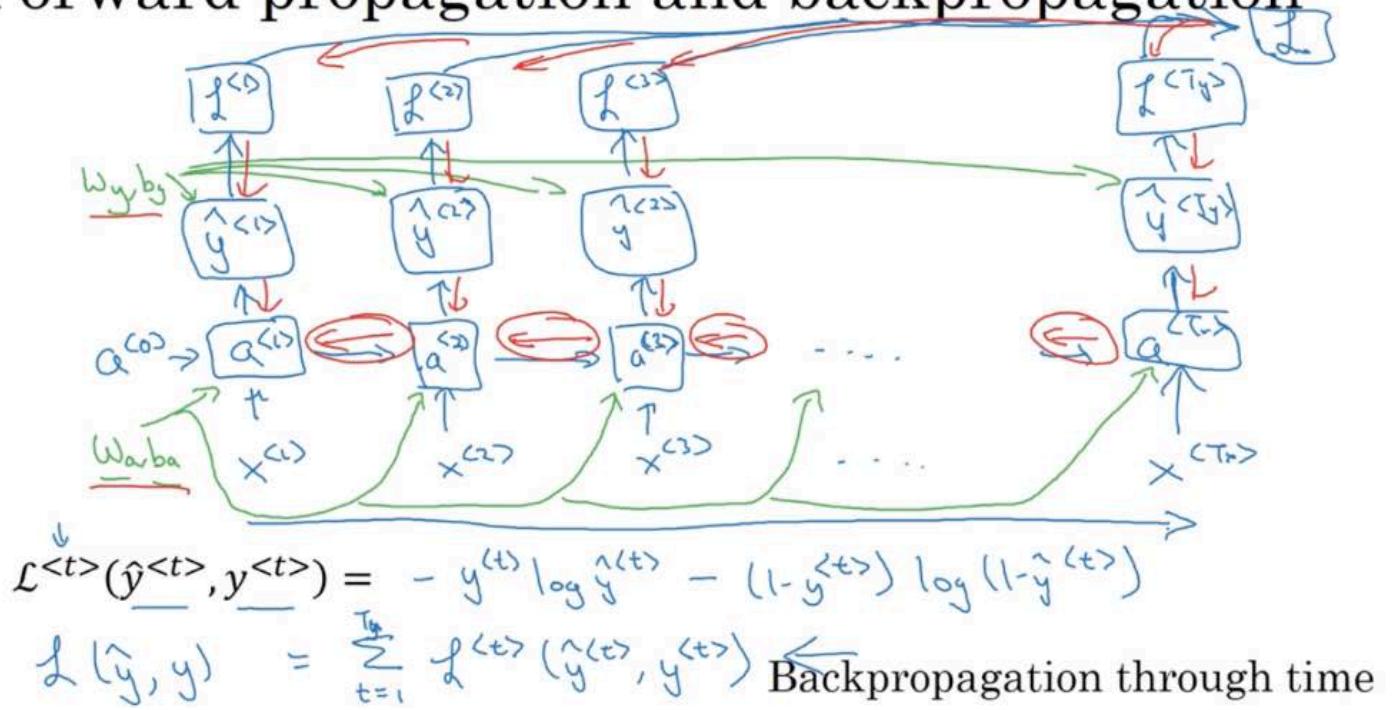
Backpropagation Through Time

Forward propagation and backpropagation



basically the opposite direction of forward propagation

Forward propagation and backpropagation



Forward Propagation:

- Input sequence: $x^{<1>} , x^{<2>} , \dots , x^{<T_x>}$
- Hidden states: $a^{<t>}$ (calculated sequentially, incorporating previous hidden state $a^{<t-1>}$ and current input $x^{<t>}$. W_{aa}, W_{ax}, b_a are involved, image is a simplified version).
- Outputs: $\hat{y}^{<t>}$ (calculated from hidden state $a^{<t>}$. W_{ya}, b_y are involved, image is a simplified version).
- Loss at each time step: $L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log(\hat{y}^{<t>}) - (1-y^{<t>}) \log(1-\hat{y}^{<t>})$ (Cross-entropy loss)

Backpropagation Through Time (BPTT):

- Total Loss: $L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{<t>}(\hat{y}^{<t>}, y^{<t>})$
- Gradients are calculated and propagated backward through time, from $t = T_y$ down to $t = 1$.
- The key characteristic is the backpropagation *through* the recurrent connections (the red arrows in the image, representing the flow of gradients through the hidden state dependencies across time).
- This unrolling represents the model computation dependencies.

Note that the provided image represents an *unrolled* Recurrent Neural Network. $W_{ax}, b_a, W_{aa}, W_{ya}$ and b_y parameters are omitted.

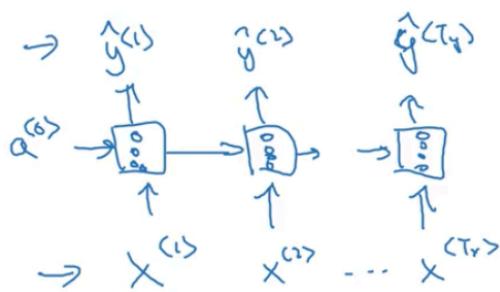
Different Types of RNNs

output could be n - 0 - or a fixed number like 5 (star rating from a sentence - sentiment)
input output numbers do not match (english to french)

- we can modify the basic RNN architecture to address all these problems

Examples of RNN architectures

$$T_x = T_y$$

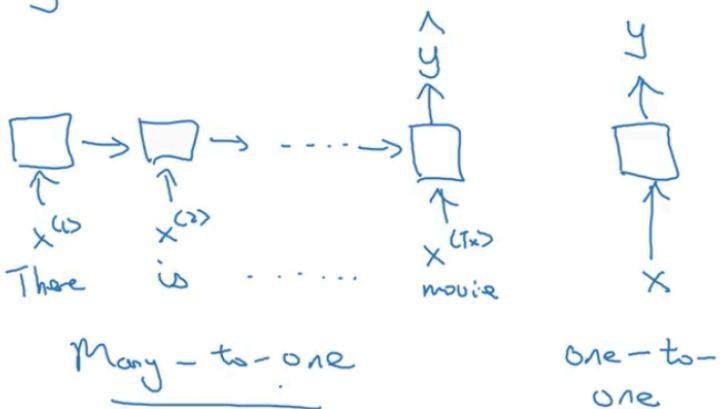


Many-to-many

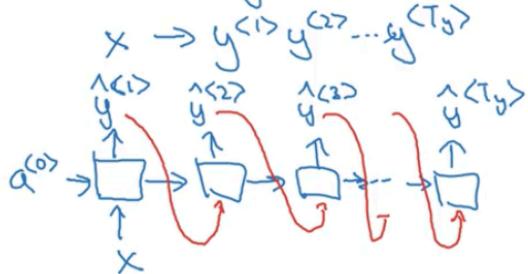
Sentiment classification

$$x = \text{text}$$

$$y = 0/1 \quad 1 \dots 5$$



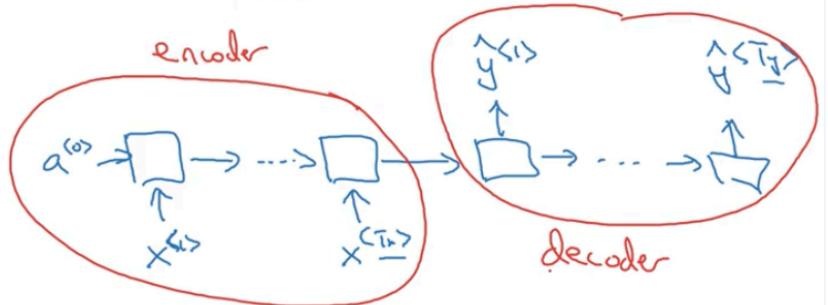
Music generation



One-to-many

$$x = \emptyset$$

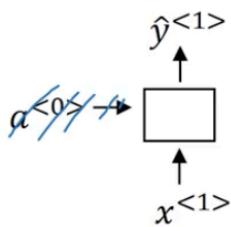
Machine translation



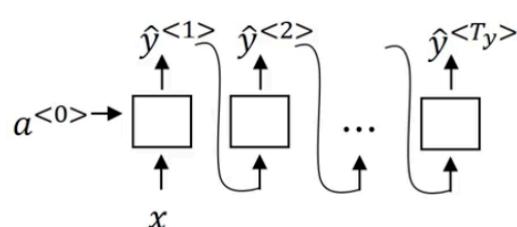
Many-to-many

the famous encoder/decoder architecture

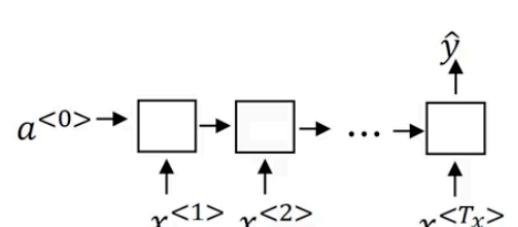
Summary of RNN types



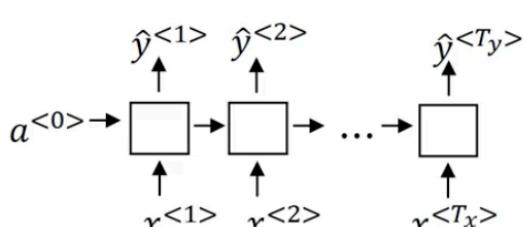
One to one



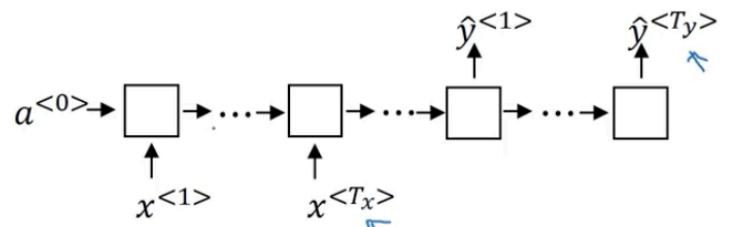
One to many



Many to one



Many to many



Many to many

Language Model and Sequence Generation

Language Modeling

Definition: A language model estimates the probability of a sequence of words, $P(y^{<1>} , y^{<2>} , \dots , y^{<T_y>})$.

Application Example: Speech Recognition

- Given an audio input, a speech recognition system might generate multiple possible transcriptions.
- A language model helps choose the most likely (most probable) sequence.

Example:

Two possible transcriptions of the same audio:

- "The apple and pair salad."
- "The apple and pear salad."

Language Model Probabilities:

- $P(\text{"The apple and pair salad"}) = 3.2 \times 10^{-13}$
- $P(\text{"The apple and pear salad"}) = 5.7 \times 10^{-10}$

Conclusion: The language model assigns a higher probability to "The apple and pear salad," making it the more likely correct transcription. The model learned that "pear" as fruit is a likely word in this context, rather than "pair".

Tokenization

Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

Cats average 15 hours of sleep a day. $\downarrow <\text{EOS}>$

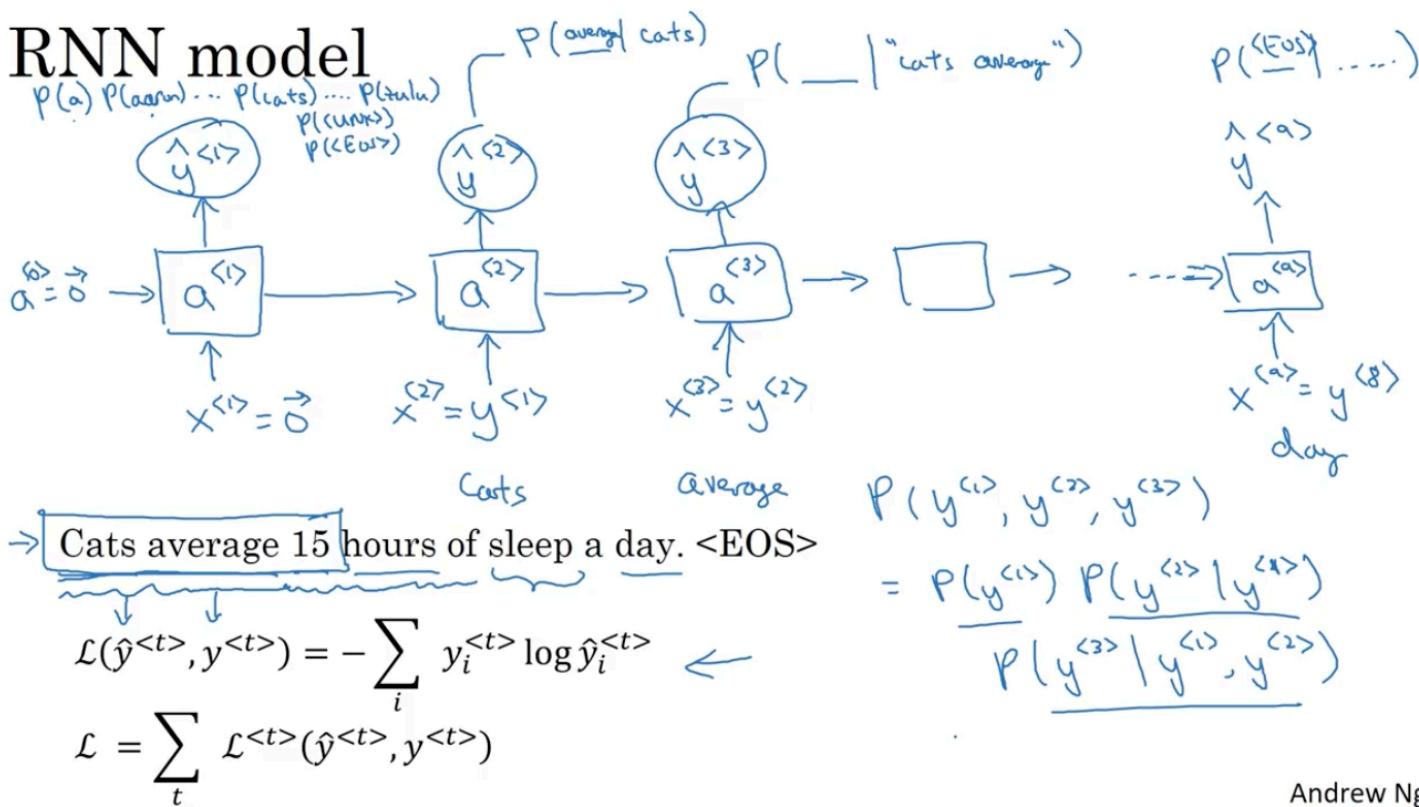
$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad \dots \quad y^{(8)} \quad y^{(9)}$

The Egyptian ~~Mau~~ is a bread of cat. $<\text{EOS}>$

$10,000 \quad <\text{UNK}>$

$<\text{EOS}>$: end of sentence token - punctuation - $<\text{UNK}>$ unique token: names etc

RNN model



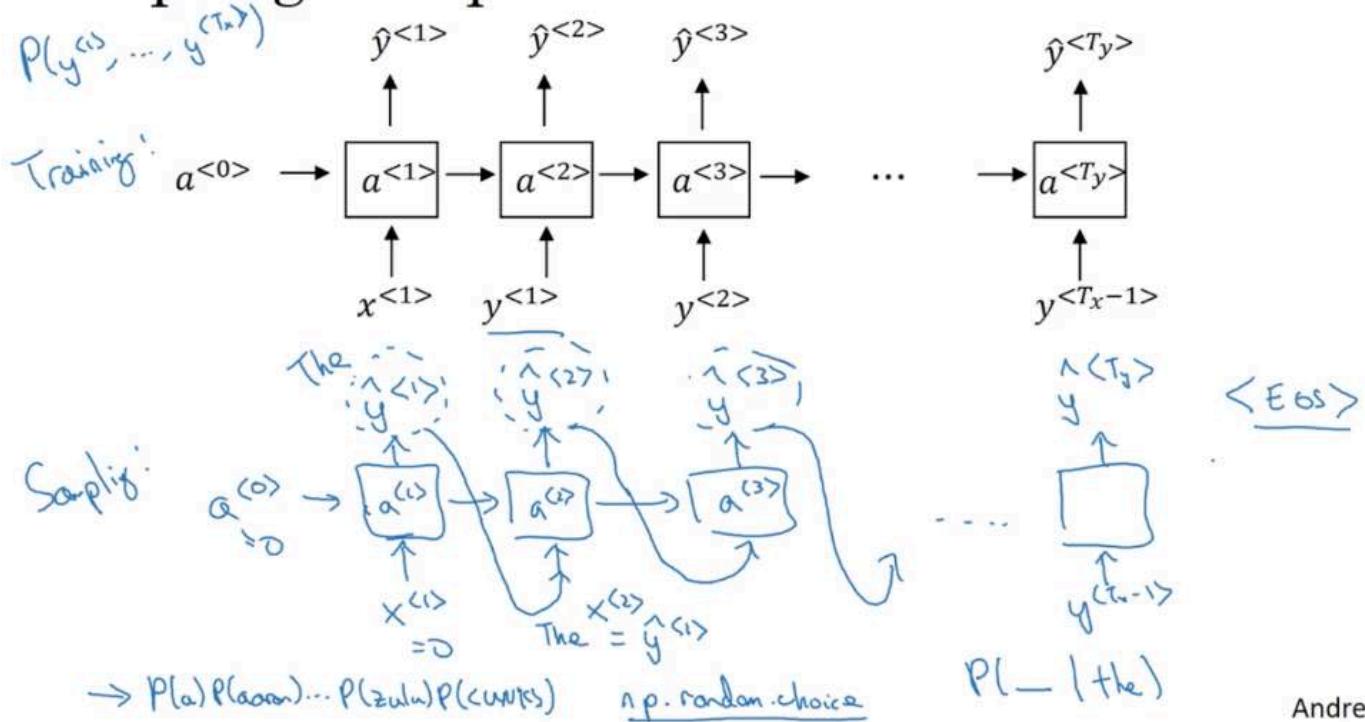
Andrew Ng

softmax loss function

tries to predict next token.

Sampling Sequences from a Trained RNN Language Model

Sampling a sequence from a trained RNN



Andrew Ng

Process:

1. Initialization:

- $a[0] = \theta$ (initial hidden state, often a zero vector).
 - $x[1] = \theta$ or <SOS> (start-of-sequence input).

2. Iterative Generation (for each time step t):

- Input: Previous hidden state ($a[t-1]$) and previous *predicted* output ($y[t-1]$). Note: During training, $y[t-1]$ is the correct previous word; during sampling, it's the *sampled* word from the

previous step._

- RNN computes: $P(y[t] \mid y[1], \dots, y[t-1])$ (probability distribution over vocabulary).
- Sample a word from this distribution (e.g., using `np.random.choice`). This sampled word becomes $y[t]$.
- $x[t+1] = y[t]$ (The sampled word becomes the input for the next time step).

3. **Termination:** Stop when `<EOS>` is sampled or max length is reached.

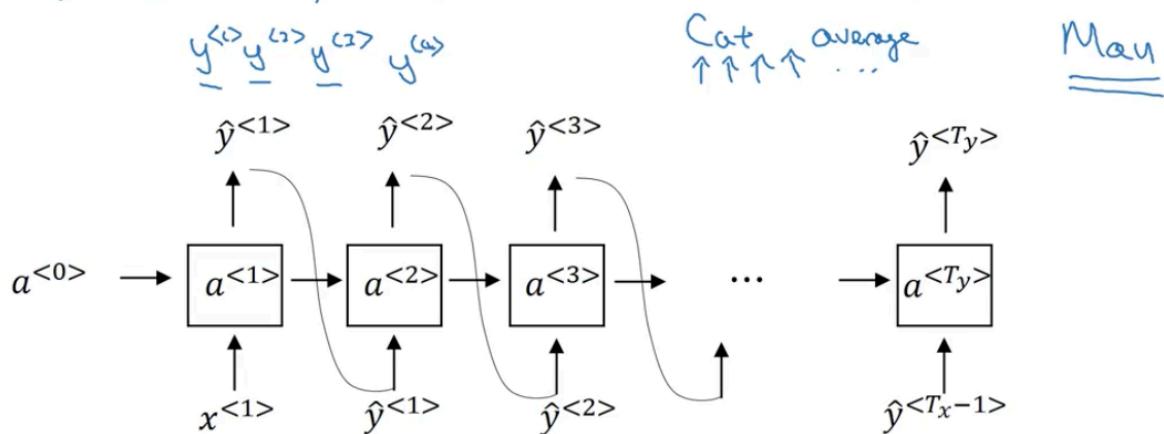
Key Difference: Training vs. Sampling

- **Training:** Uses *correct* previous word ($y[t-1]$) as input. Learns $P(y[t] \mid y[1], \dots, y[t-1])$.
- **Sampling:** Uses *predicted* previous word ($\hat{y}[t-1]$, initially sampled, then fed back) as input. Generates new sequences using the learned probability distribution.

Character-level language model

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

→ Vocabulary = [a, b, c, ..., z, , ., , ;, , o, ..., q, A, ..., Z]



Andrew Ng

Sequence generation

News

President enrique peña nieto, announced
sench's sulk former coming football langston
paring.

“I was not at all surprised,” said hich langston.

“Concussion epidemic”, to be examined. ←

The gray football the told some and this has on
the uefa icon, should money as.

Shakespeare

The mortal moon hath her eclipse in love.

And subject of this thou art another this fold.

When lesser be my love to me see sabl's.

For whose are ruse of mine eyes heaves.

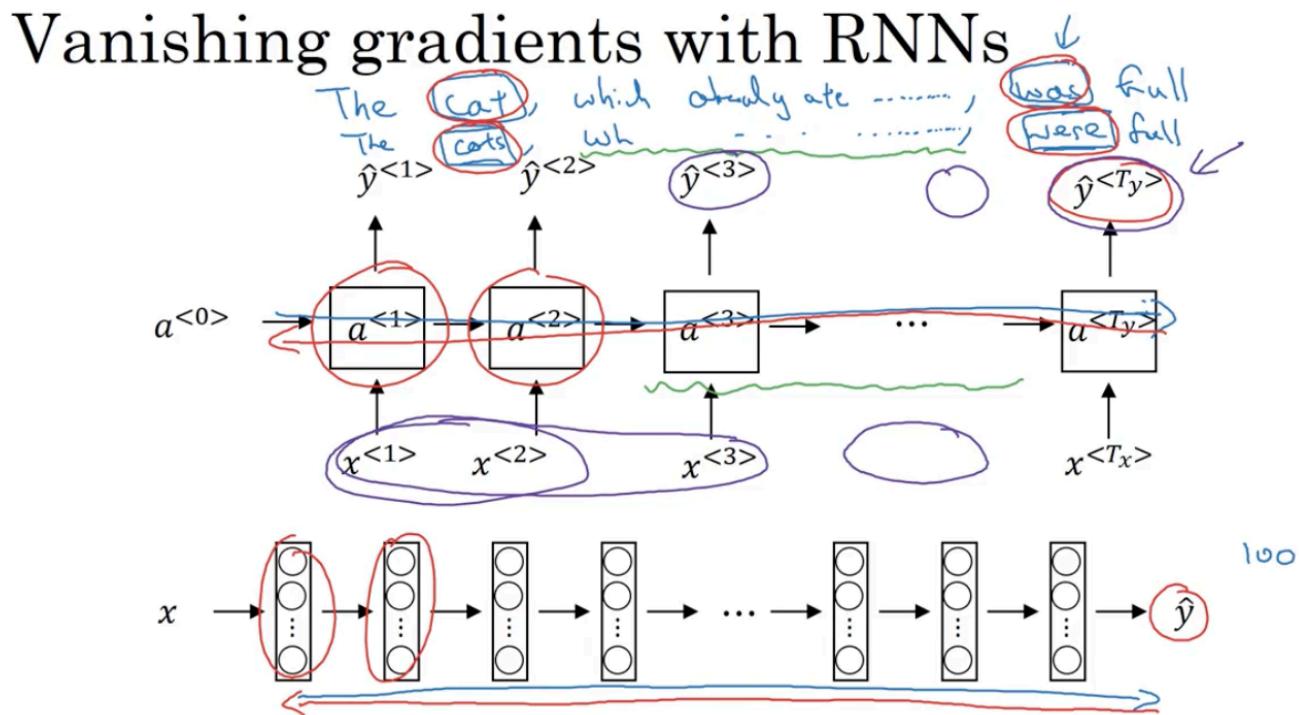
whatever is the input, is the output

Vanishing Gradients with RNNs

RNNs tend to have vanishing gradients problem

very early tokens, may have effect on later tokens, base RNN architecture is not good at preserving such effects to later 'layers'.

- the cat was
- the cats were



highly local, influenced mainly by inputs near that output

weakness of basic RNN

- if not addressed, not good at capturing long term dependencies.

also this can cause exploding gradients

- easier to spot, parameters blow up - NaN
- soln: gradient clipping - look gradient vector, if bigger than X, clip. Robust solution

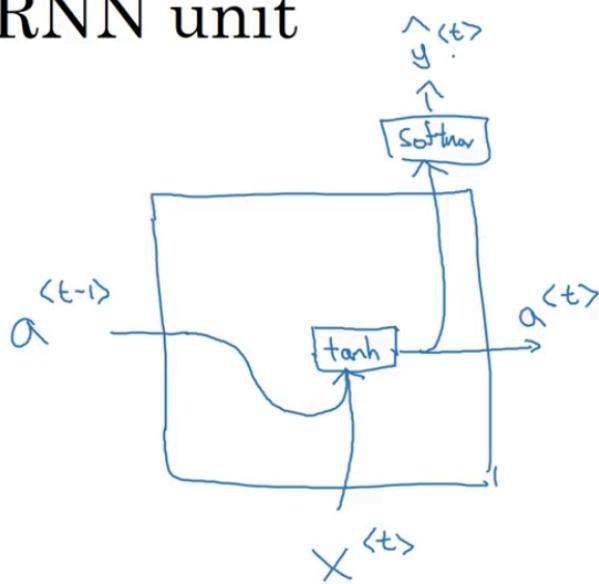
vanishing gradients are much harder to solve

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

Gated Recurrent Unit (GRU)

RNN Unit:

RNN unit

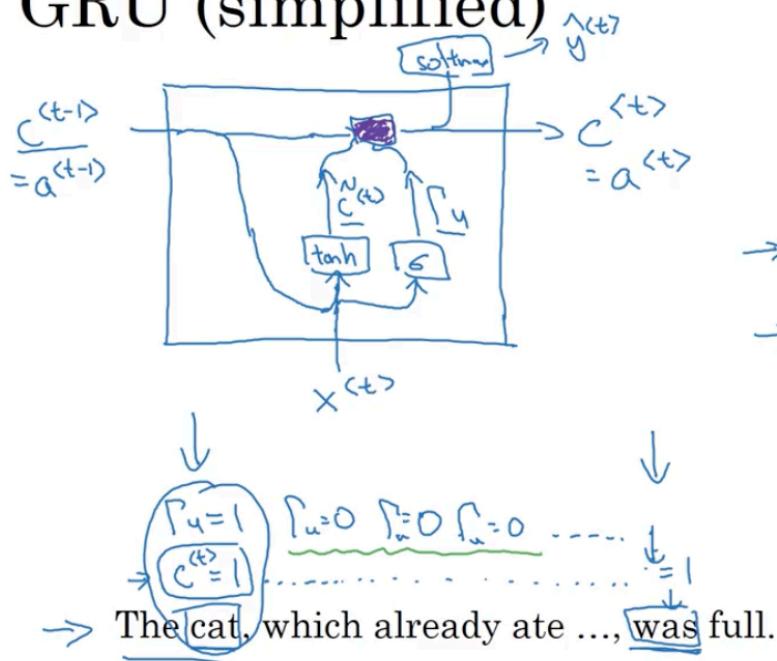


$$a^{<t>} = \tanh(W_a[a^{<t-1>}, x^{<t>}] + b_a)$$

GRU (simplified):

for LSTM, c of t and a of t are not equal, though that is for later

GRU (simplified)



[Cho et al., 2014. On the properties of neural machine translation: Encoder-decoder approaches] ←
 [Chung et al., 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling] ←

Andrew Ng

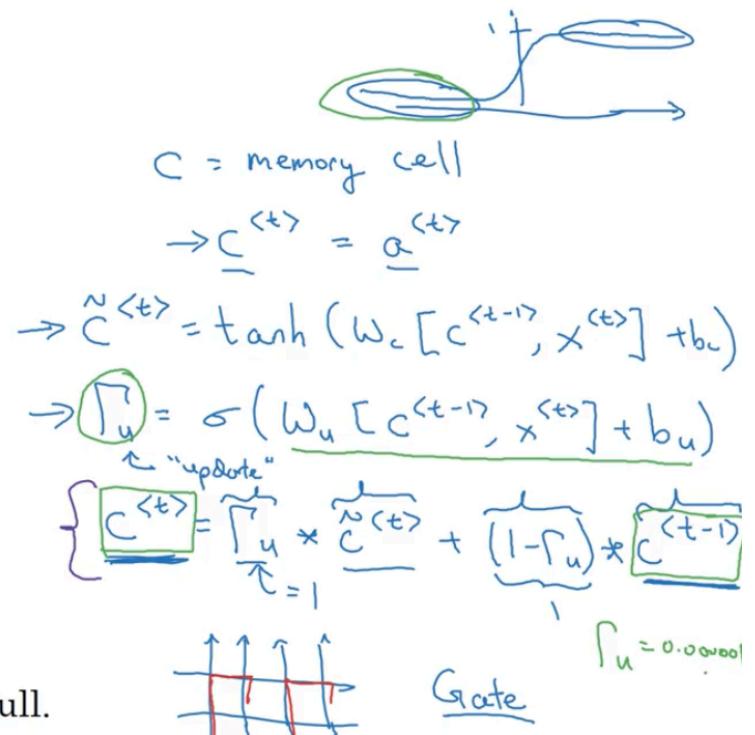
!!!

very good at maintaining the value of the cell - not changing the memory value

c_t can be a x -dimensional-vector (generally)

gamma won't be exactly 0 and 1 mostly

you can choose some bits of the vector to keep them constant and choose some to update - with adjusting what (?)



Full GRU

$$\tilde{h} \quad \underbrace{\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)}_{\text{u}}$$

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

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

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

LSTM

The cat, which ate already, was full.

Andrew Ng

Long Short Term Memory (LSTM)

Even more powerful than the GRU Unit

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

instead of one update gate controlling, for LSTM there are two gamma's, update & forgot also new output gate

GRU and LSTM

GRU

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

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

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

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

$$a^{<t>} = c^{<t>}$$

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

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

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

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

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

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

[Hochreiter & Schmidhuber 1997. Long short-term memory] ←

Andrew Ng

LSTM in pictures

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

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

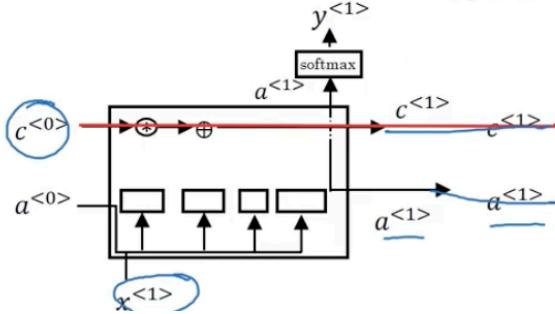
$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

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

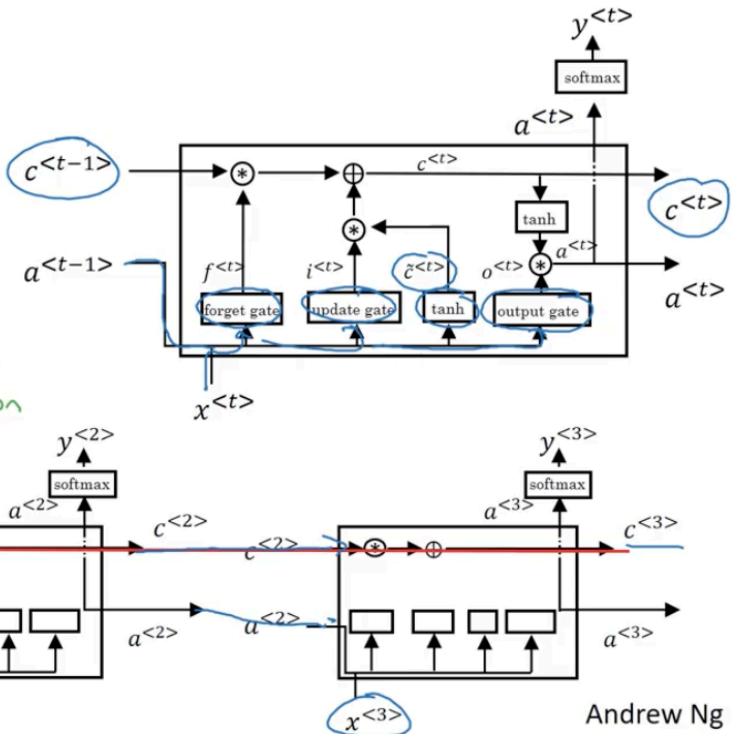
$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

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

peephole
connection



Bidirectional RNN

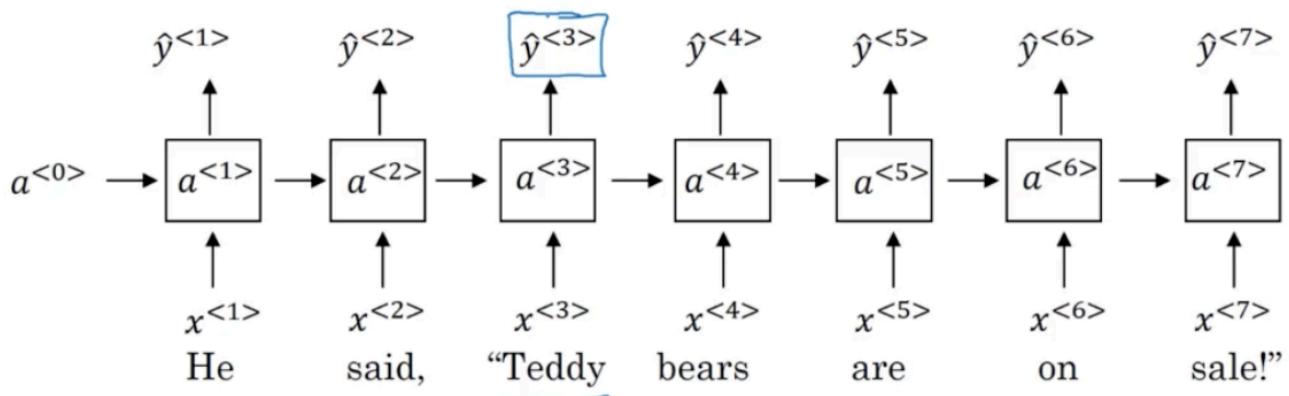


Andrew Ng

Getting information from the future

He said, "Teddy bears are on sale!"

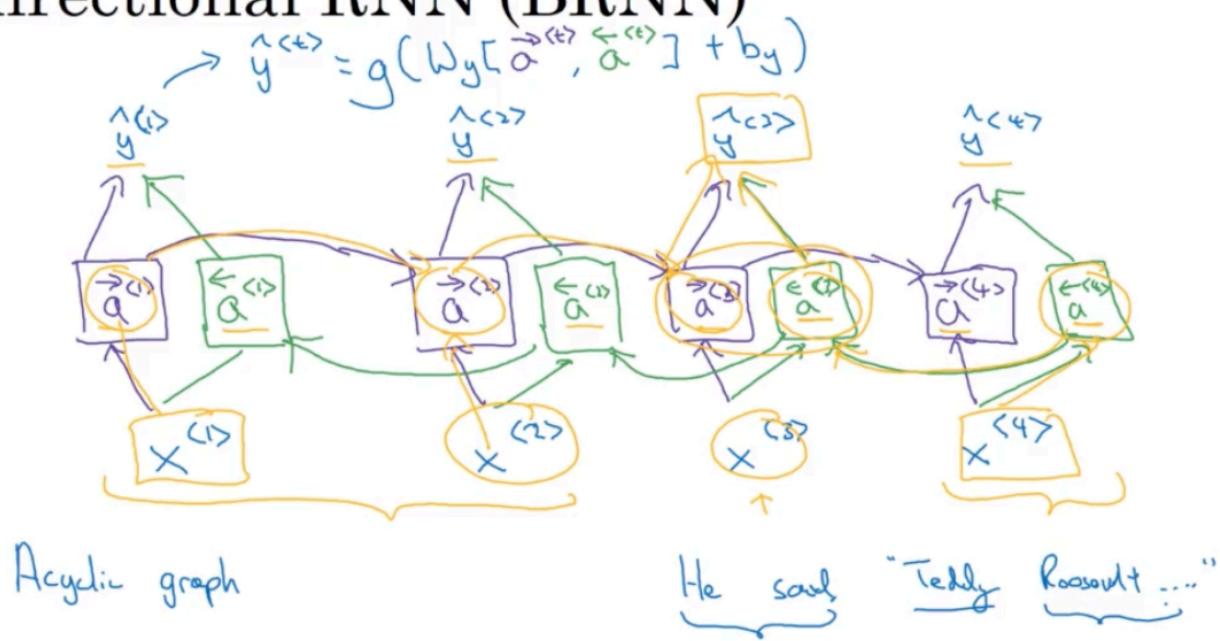
He said, "Teddy Roosevelt was a great President!"



Andrew Ng

these blocks can be either basic RNN, GRU, or LSTM units.

Bidirectional RNN (BRNN)



Andrew Ng

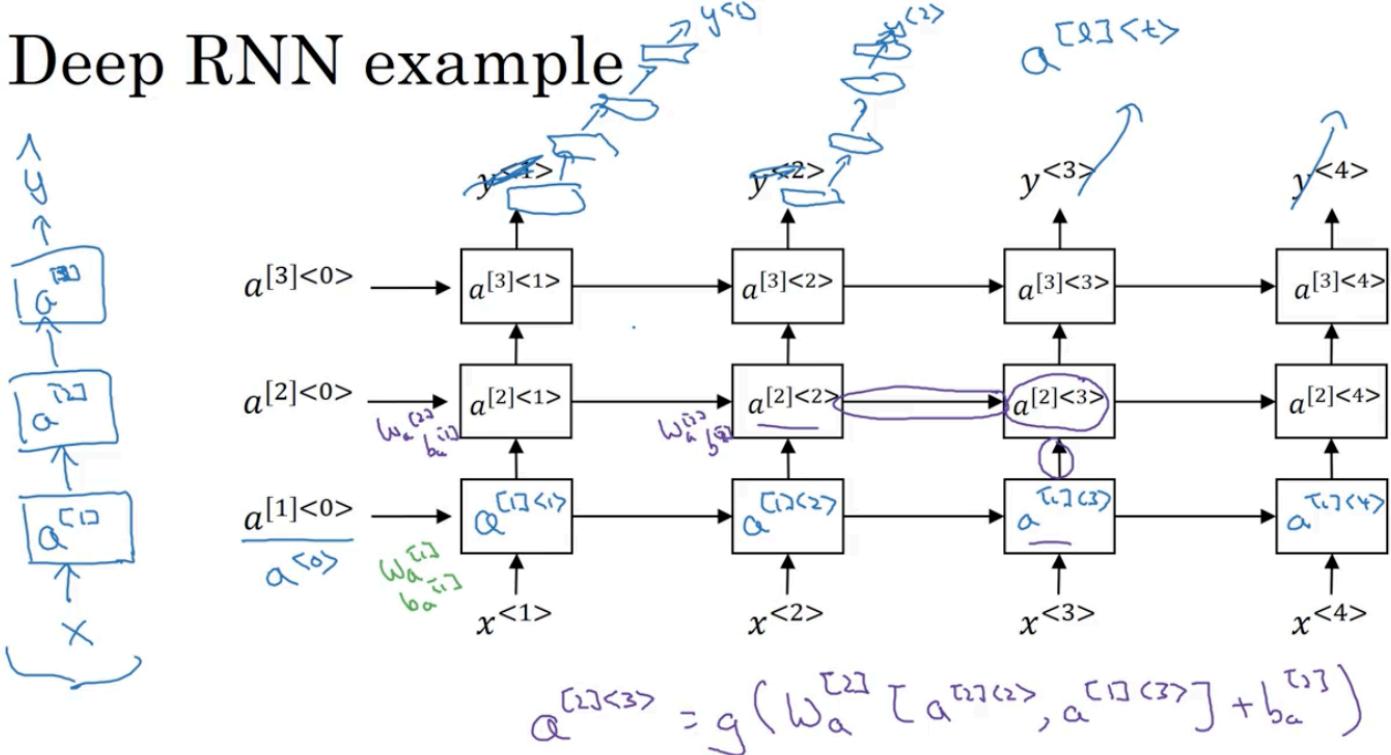
forward propagation, but it goes both to forward and 'backward'
information gets carried in both directions

- BRNN w\ LSTM blocks - used frequently for NLP problems
Disadvantage: Need entire sequence of data to start processing / make prediction

Deep RNNs

RNNs work as they are already, deeper models make them better.

Deep RNN example



Andrew Ng

3 hidden layers

blocks does not have to be regular RNN blocks

there can be only horizontal connections - there are other versions of this network

What you should remember:

- The recurrent neural network, or RNN, is essentially the repeated use of a single cell.
- A basic RNN reads inputs one at a time, and remembers information through the hidden layer activations (hidden states) that are passed from one time step to the next.
 - The time step dimension determines how many times to re-use the RNN cell
- Each cell takes two inputs at each time step:
 - The hidden state from the previous cell
 - The current time step's input data
- Each cell has two outputs at each time step:
 - A hidden state
 - A prediction
- An LSTM is similar to an RNN in that they both use hidden states to pass along information, but an LSTM also uses a cell state, which is like a long-term memory, to help deal with the issue of vanishing gradients
- An LSTM cell consists of a cell state, or long-term memory, a hidden state, or short-term memory, along with 3 gates that constantly update the relevancy of its inputs:
 - A **forget** gate, which decides which input units should be remembered and passed along. It's a tensor with values between 0 and 1.
 - If a unit has a value close to 0, the LSTM will "forget" the stored state in the previous cell state.
 - If it has a value close to 1, the LSTM will mostly remember the corresponding value.

- An **update** gate, again a tensor containing values between 0 and 1. It decides on what information to throw away, and what new information to add.
 - When a unit in the update gate is close to 1, the value of its candidate is passed on to the hidden state.
 - When a unit in the update gate is close to 0, it's prevented from being passed onto the hidden state.
 - And an **output** gate, which decides what gets sent as the output of the time step
 - Very large, or "exploding" gradients updates can be so large that they "overshoot" the optimal values during back prop -- making training difficult
 - Clip gradients before updating the parameters to avoid exploding gradients
 - Sampling is a technique you can use to pick the index of the next character according to a probability distribution.
 - To begin character-level sampling:
 - Input a "dummy" vector of zeros as a default input
 - Run one step of forward propagation to get $a(1)$ (your first character) and $\hat{y}(1)$ (probability distribution for the following character)
 - When sampling, avoid generating the same result each time given the starting letter (and make your names more interesting!) by using `np.random.choice`
 - A sequence model can be used to generate musical values, which are then post-processed into midi music.
 - You can use a fairly similar model for tasks ranging from generating dinosaur names to generating original music, with the only major difference being the input fed to the model.
 - In Keras, sequence generation involves defining layers with shared weights, which are then repeated for the different time steps $1, \dots, Tx$.
 - A sequence model can be used to generate musical values, which are then post-processed into midi music.
 - You can use a fairly similar model for tasks ranging from generating dinosaur names to generating original music, with the only major difference being the input fed to the model.
 - In Keras, sequence generation involves defining layers with shared weights, which are then repeated for the different time steps $1, \dots, Tx$.
-

NLP and Word Embeddings

Word Representations

Word representation

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

$$|V| = 10,000$$

1-hot representation

Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$
\hookrightarrow	\hookrightarrow	\hookrightarrow	\hookrightarrow	\hookrightarrow	\hookrightarrow
e_{5391}	e_{9853}				

I want a glass of orange juice.
I want a glass of apple ?.

Andrew Ng

it makes it harder for algorithm to learn the relations, makes it hard for it to generalize, embeddings are distant - no direct relation in any way

Featurized representation: word embedding

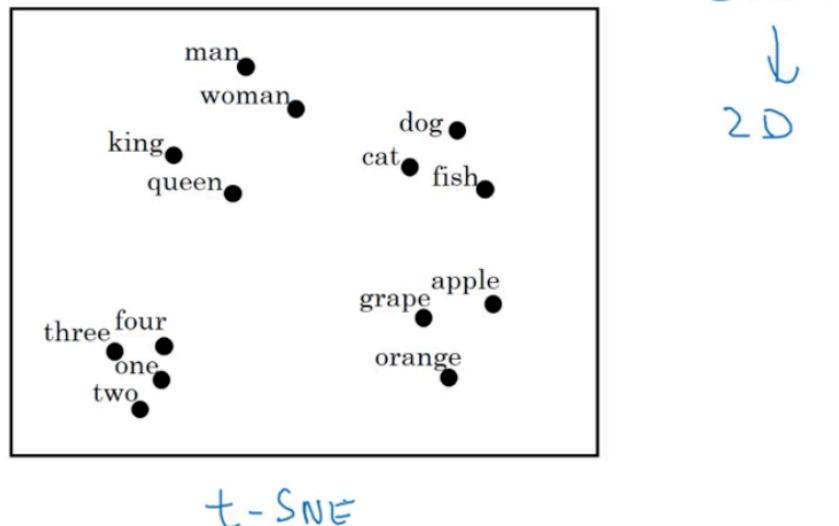
	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.62	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
size						
cost						
alt						
verb						

I want a glass of orange ____.
I want a glass of apple ____.

Andrew Ng

with this approach, representations for orange and apple take on similar values

Visualizing word embeddings



[van der Maaten and Hinton., 2008. Visualizing data using t-SNE]

Andrew Ng

another general approach is to visualize them by embedding these vectors into 2D, a commonly used algorithm for this task: t-SNE

Using Word Embeddings

- Sally Johnson is an orange farmer.
- 1 1 0 0 0 0

New Input:

- Robert Lin is an apple farmer
 - Robert Lin is a durian cultivator
- may not know the durian if small data, but with word embedding, may even understand and match "durian cultivator" to "orange farmer".

1B - 100B words

by reading massive amounts of text

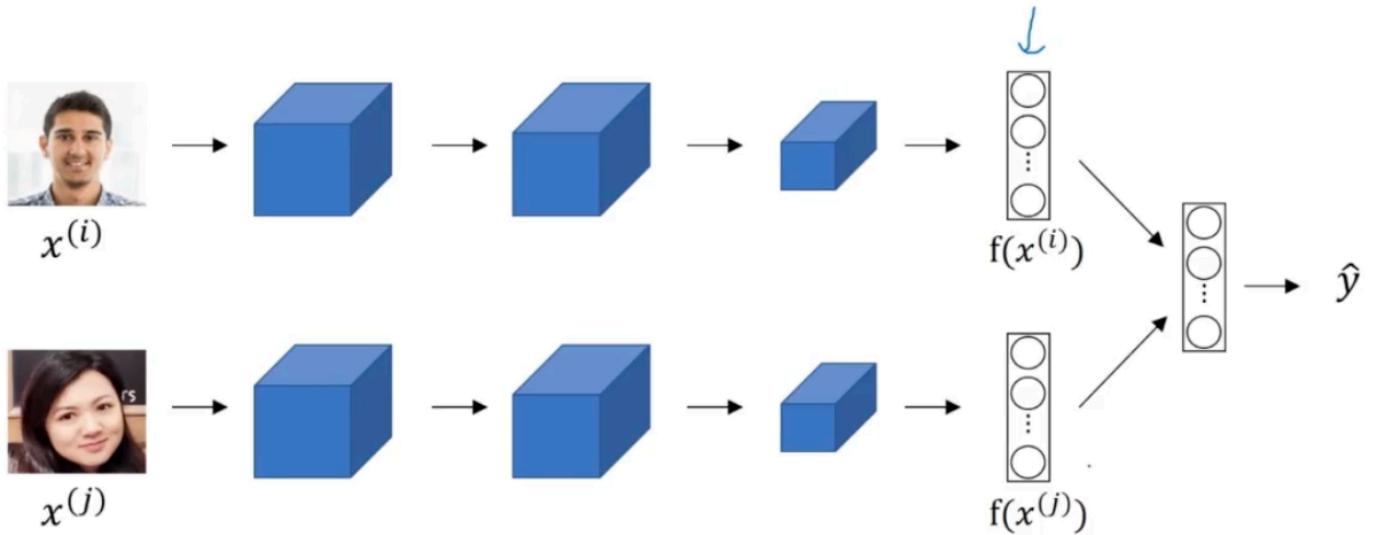
transfer this knowledge to task like name entity recognition

(probably will use BRNN for such task)

Transfer Learning and Word Embeddings

- 1- Learn word embeddings from large text corpus (1-100B words) (Or download pre-trained embedding)
- 2- Transfer embedding to new task with smaller training set (ex: 100k)
- 3- (Opt) Continue to finetune the word embeddings with new data

Relation to face encoding



[Taigman et. al., 2014. DeepFace: Closing the gap to human level performance]

Andrew Ng

encoding = embedding

encoding: even take in a face that has never seen before, and creates a encoding for it

embedding: have a fixed size vocabulary, embeds sentences and meanings

- difference in between the terms

Properties of Word Embeddings

Analogy:

Man -> Woman as King -> ??? - Can the algorithm (embedding) figure this out by itself?

$$\begin{aligned} e_{\text{man}} - e_{\text{woman}} &\approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix} \\ e_{\text{king}} - e_{\text{queen}} &\approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix} \end{aligned}$$

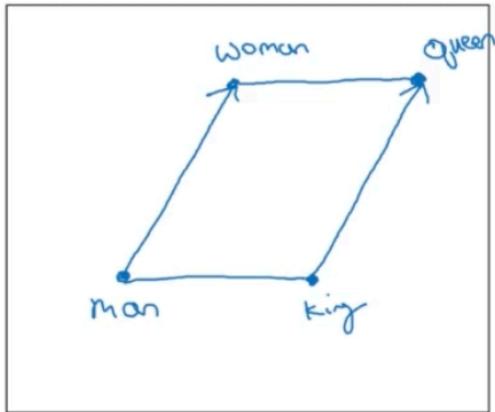
yes

e_{man} - e_{woman} : try to find a vector that is equivalent to this - that is: e_{king} - e_{queen}

- one of the first surprising findings about word embeddings

t-SAE: takes data and plots non-linearly, there is not relationship as such given in the below example (in most cases)

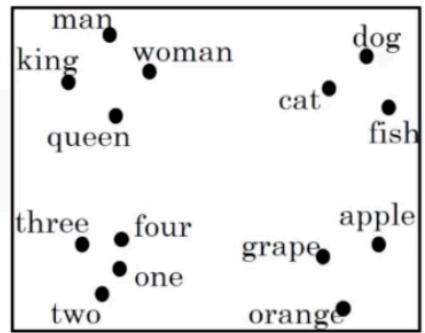
Analogies using word vectors



300D

Find word w : $\arg \max_w$

$300D \rightarrow 2D$



t-SNE

$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\underline{w}}$$

$$\text{Sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$$

30 - 75%

Andrew Ng

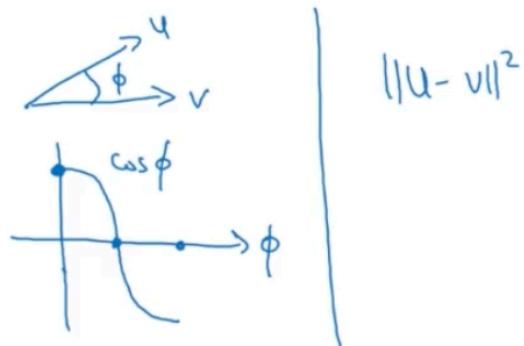
Similarity Function

most commonly used: Cosine similarity

Cosine similarity

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

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



Man:Woman as Boy:Girl

Ottawa:Canada as Nairobi:Kenya

Big:Bigger as Tall:Taller

Yen:Japan as Ruble:Russia

Andrew Ng

Embedding Matrix

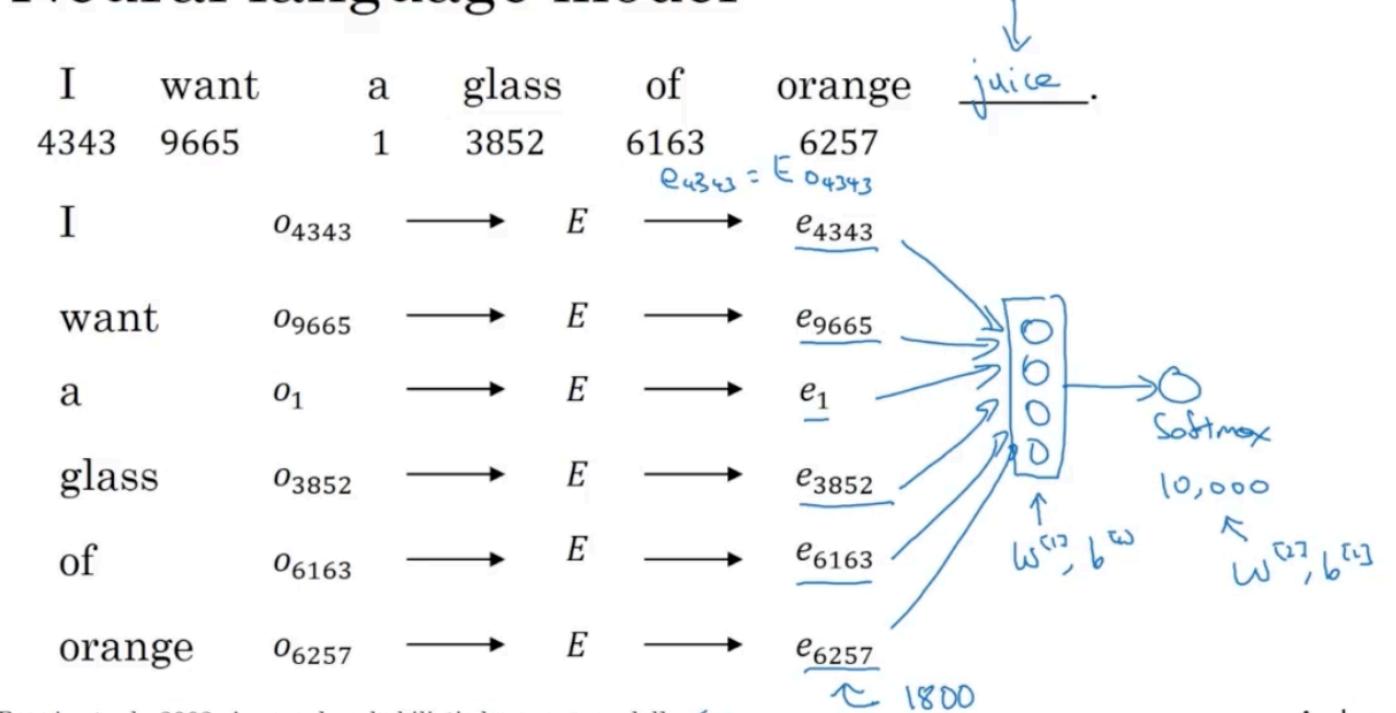
- A matrix learned during word embedding algorithms.
- Dimensions: (embedding dimension) x (vocabulary size). E.g., 300 x 10,000.
- Each column represents the embedding vector for a specific word in the vocabulary.
- E multiplied by a word's one-hot vector (e.g., O_{6257} for "Orange") extracts the corresponding embedding vector (e.g., e_{6257}). Because all elements in one-hot vector are 0 except one, this

multiplication gives a column, that represents that specific word, from the E matrix.

- The result of Multiplying $E O_n$ is equal to E_n .

Word Embeddings

Neural language model

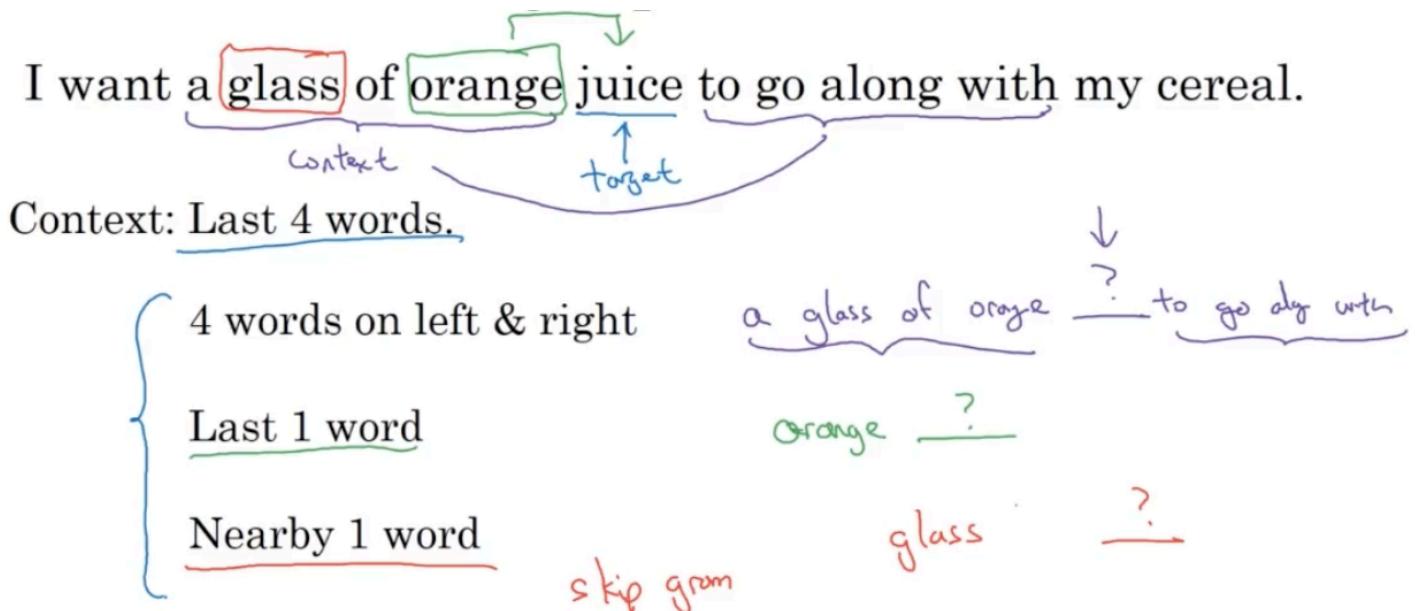


[Bengio et. al., 2003, A neural probabilistic language model]

Andrew Ng

- 6 embedding matrixes stacked together for each word: 1800 dimensions
- another hyperparameter, how many words to back will be looked when predicting the output (ex: always 4 words) - making input size fixed -> 1800 becomes $4 \times 300 = 1200$
- earlier but pretty successful algorithm

Other context/target pairs

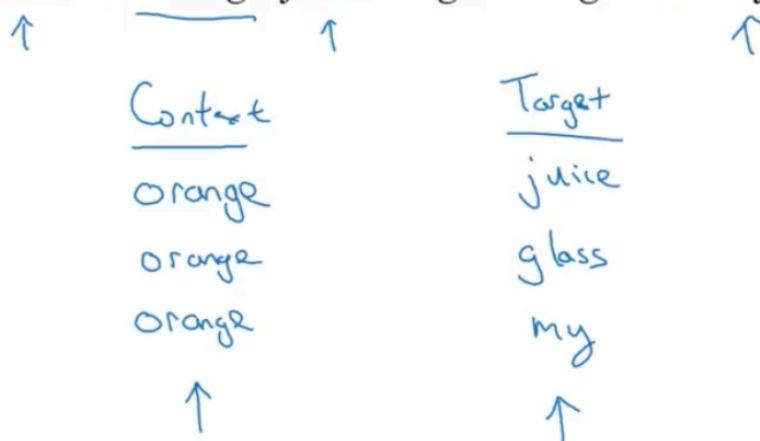


if want to build Language Model, general to give the first type of context - if you just want to learn the embeddings, other choices are relevant too.

Word2Vec

Skip-grams

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



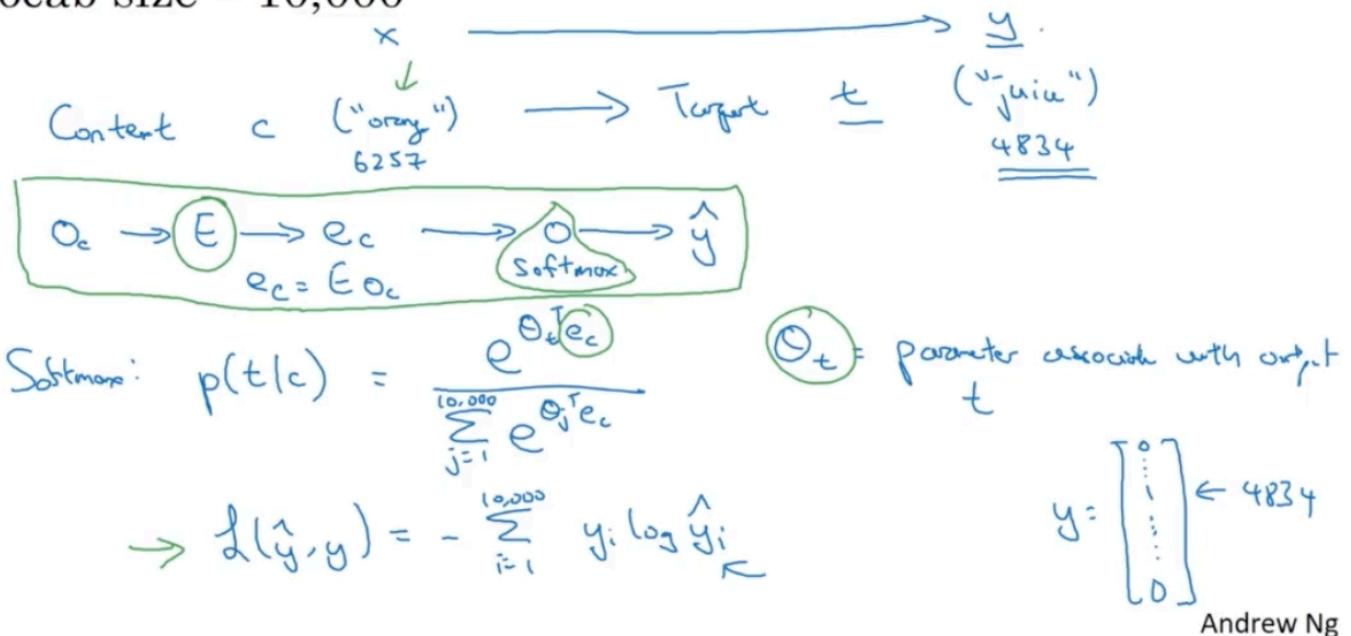
pick a target within some window (ex: +/- 5)

setup a supervised learning problem

not an easy learning problem - there could be a lot of words

Model

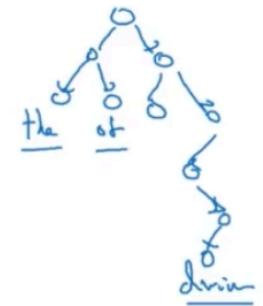
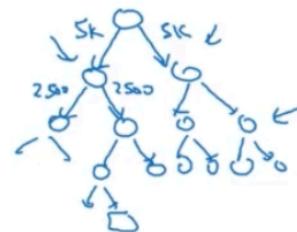
Vocab size = 10,000



Problems with softmax classification

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

log |V|



How to sample the context c ?

→ the, of, a, and, to, ...

t

→ orange, apple, durian

$c \rightarrow t$

Q_{durian}

$P(c)$

hierarchical softmax: does not use equally weighted trees, less commonly words being buried deep in the tree - do not need to go that deep generally

distribution is not taken entirely uniformly, does not want to update values frequently for words that does not come up frequently - different heuristic options for how to choose

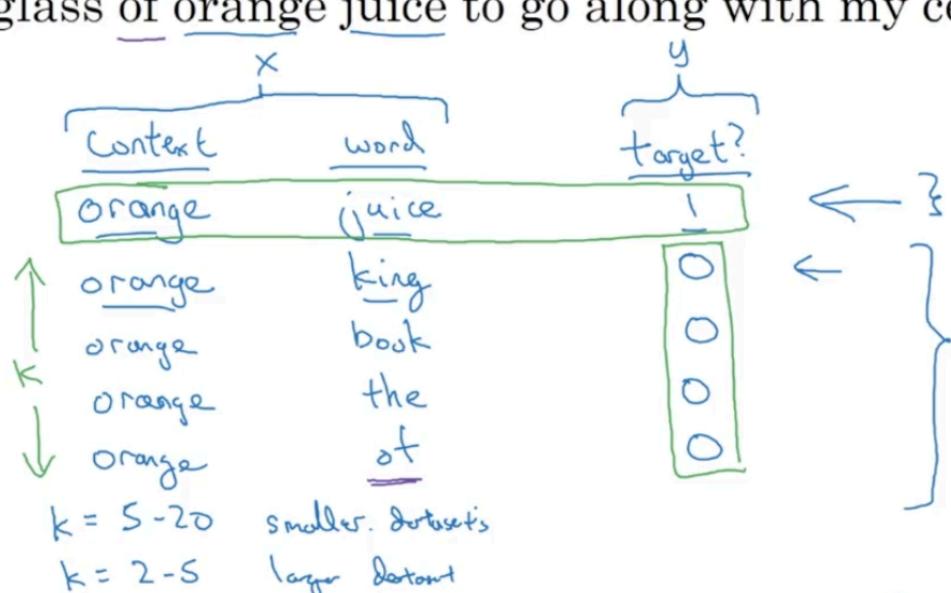
Negative Sampling

something similar to skip-gram model, but much more efficient

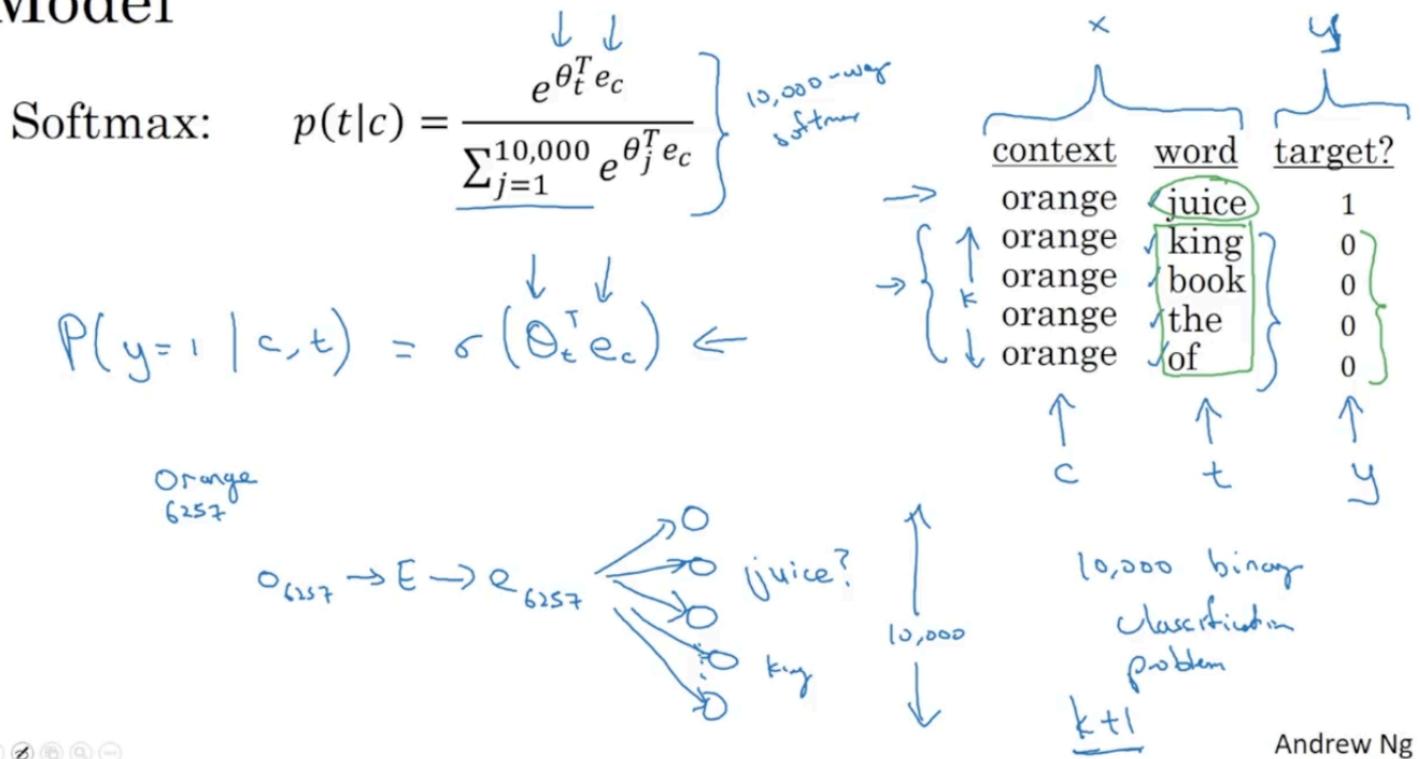
give a positive pair, then for some number of examples, say 'k', put in negative examples. Then create a supervised learning problem that inputs x and predicts the target y.

Defining a new learning problem

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



Model



after you choose the word 'orange', how do you sample negative examples?

random - there are words that comes up frequently - the, of

solt: take an heuristic value, sample proportional to frequency of the word

- random value but is used: $^{3/4}$

Selecting negative examples

the , of, and, ...

context	word	target?
orange	juice	1
orange	king	0
orange	book	0
orange	the	0
orange	of	0

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

$\frac{1}{|V|}$

GloVe Word Vectors

not used as much, but has some enthusiast.

GloVe (global vectors for word representation)

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

c, t

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

$$X_{ij} = X_{ji} \leftarrow$$

[Pennington et. al., 2014. GloVe: Global vectors for word representation]

Andrew Ng

if equal, they always exist together

Model

$$\text{minimize} \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\theta_i^T e_j + b_i + b_j' - \log X_{ij})^2$$

weight term
 $\theta_i^T e_j$
 $f(X_{ij}) = 0 \text{ or } X_{ij} = 0. \quad "0 \log 0" = 0$
 this, is, at, a, ...
 devian

"0?"
 $b_i + b_j' - \log X_{ij}$
 $e_w^{(\text{final})} = \frac{e_w + \theta_w}{2}$
 $\theta_i, e_i \text{ are symmetric}$

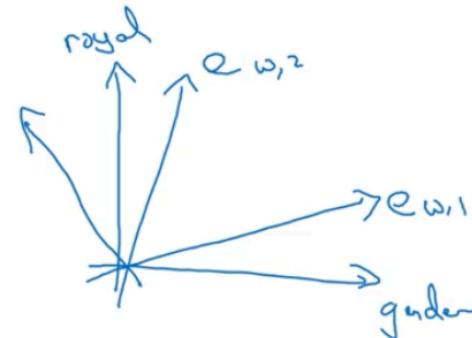
Andrew Ng

simple, but works.

$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\theta_i^T e_j + b_i + b_j' - \log X_{ij})^2$$

A note on the featurization view of word embeddings

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)
Gender	-1	1	-0.95	0.97
Royal	0.01	0.02	0.93	0.95
Age	0.03	0.02	0.70	0.69
Food	0.09	0.01	0.02	0.01



$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\underbrace{\theta_i^T e_j + b_i - b'_j - \log X_{ij}}_0)^2$$

$\leftarrow (\mathbf{A}\theta_i)^T (\mathbf{A}^T e_j) = \underline{\theta_i^T \mathbf{A}^T \mathbf{A} \theta_i}$

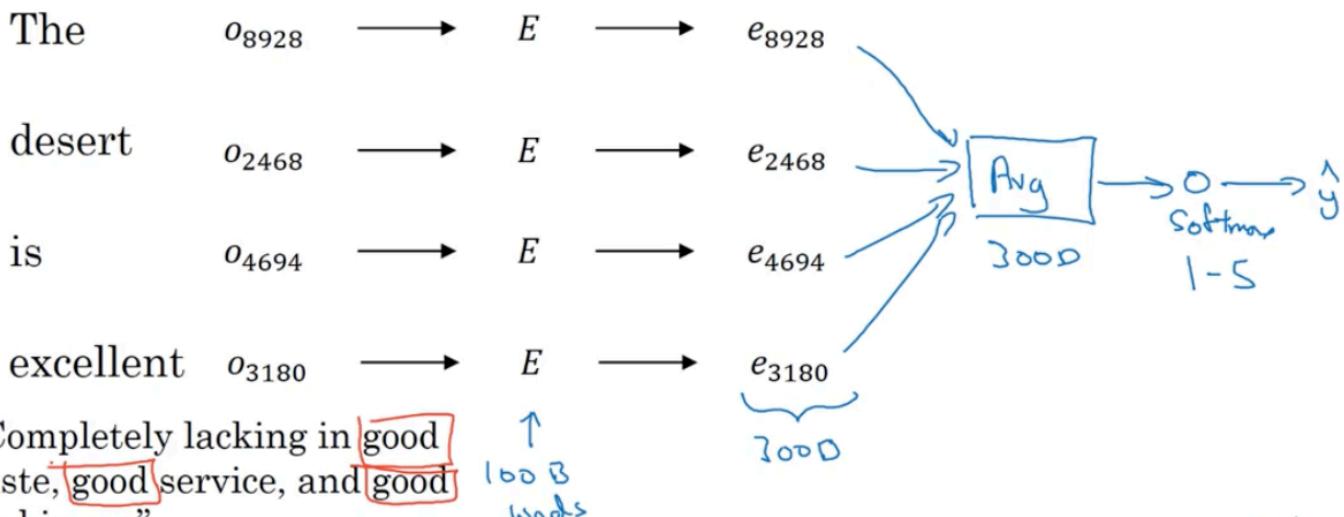
Andrew Ng

brief proof that with an algorithm like this, axis used by these vectors will be easily interpretable (?)

Sentiment Classification

Simple sentiment classification model

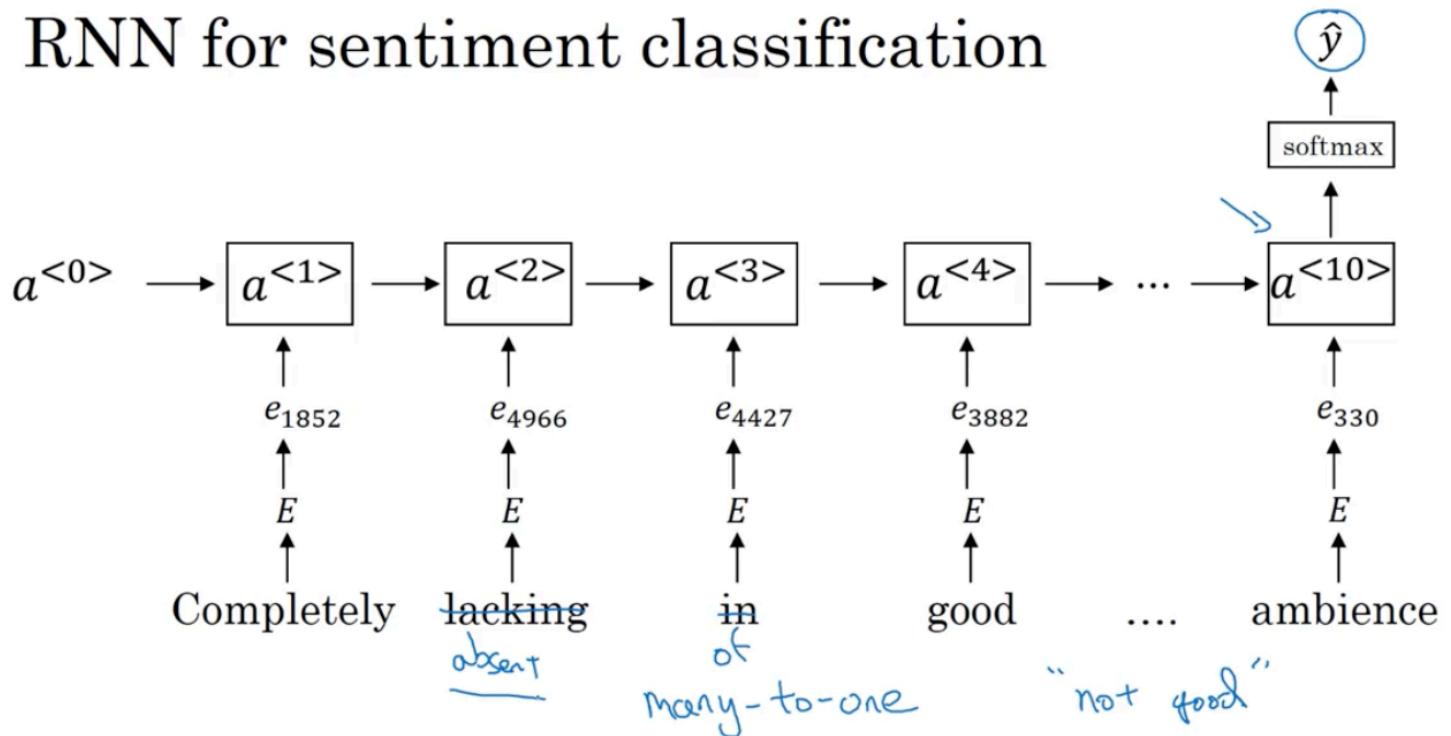
The dessert is excellent
8928 2468 4694 3180



Andrew Ng

below example becomes a good review even though it is not, soln:

RNN for sentiment classification



Andrew Ng

pretty decent sentiment analysis algorithm, better job generalizing, even the words that has not been used specifically

Debiasing Word Embeddings

gender bias, ethnicity bias, sexual orientation bias... - how to diminish

The problem of bias in word embeddings

Man:Woman as King:Queen

Man:Computer_Programmer as Woman:Homemaker X

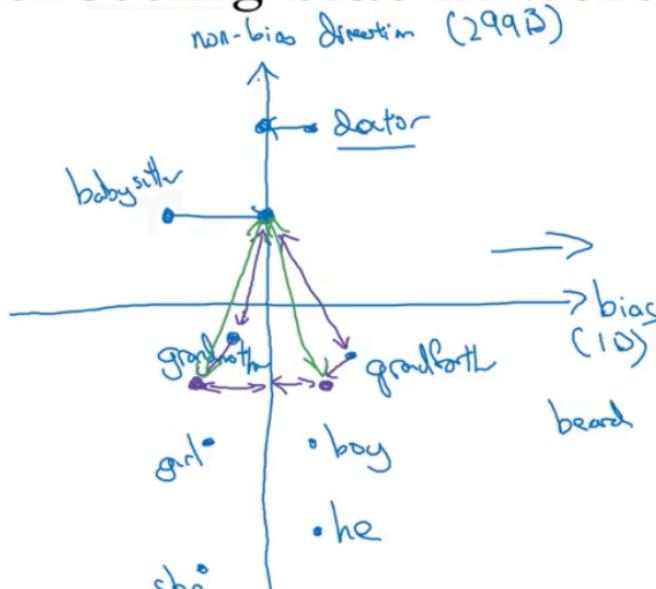
Father:Doctor as Mother:Nurse X

Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of the text used to train the model.

[Bolukbasi et. al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings] L

Andrew Ng

Addressing bias in word embeddings



1. Identify bias direction.

$$\begin{cases} \text{he} - \text{she} \\ \text{male} - \text{female} \end{cases} \rightarrow \text{average}$$

2. Neutralize: For every word that is not definitional, project to get rid of bias.

3. Equalize pairs.

$$\rightarrow \begin{cases} \text{grandmother} - \text{grandfather} \\ \text{girl} - \text{boy} \end{cases}$$

[Bolukbasi et. al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings]

Andrew Ng

What you should remember:

- Cosine similarity is a good way to compare the similarity between pairs of word vectors.
 - Note that L2 (Euclidean) distance also works.
- For NLP applications, using a pre-trained set of word vectors is often a great way to get started.
- Even with a mere 127 training examples, you can get a reasonably good model for Emojifying.
 - This is due to the generalization power word vectors gives you.
- Emojify-V1 will perform poorly on sentences such as "*This movie is not good and not enjoyable*"
 - It doesn't understand combinations of words.
 - It just averages all the words' embedding vectors together, without considering the ordering of words.
- If you have an NLP task where the training set is small, using word embeddings can help your algorithm significantly.
- Word embeddings allow your model to work on words in the test set that may not even appear in the training set.
- Training sequence models in Keras (and in most other deep learning frameworks) requires a few important details:
 - To use mini-batches, the sequences need to be **padded** so that all the examples in a mini-batch have the **same length**.
 - An `Embedding()` layer can be initialized with pretrained values.
 - These values can be either fixed or trained further on your dataset.
 - If however your labeled dataset is small, it's usually not worth trying to train a large pre-trained set of embeddings.
 - `LSTM()` has a flag called `return_sequences` to decide if you would like to return all hidden states or only the last one.

- You can use `Dropout()` right after `LSTM()` to regularize your network.

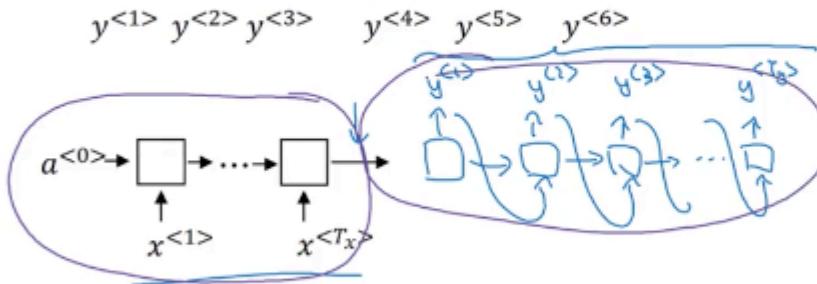
Various Sequence To Sequence Architectures

Basic Models

Sequence to sequence model

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$
 Jane visite l'Afrique en septembre

→ Jane is visiting Africa in September.



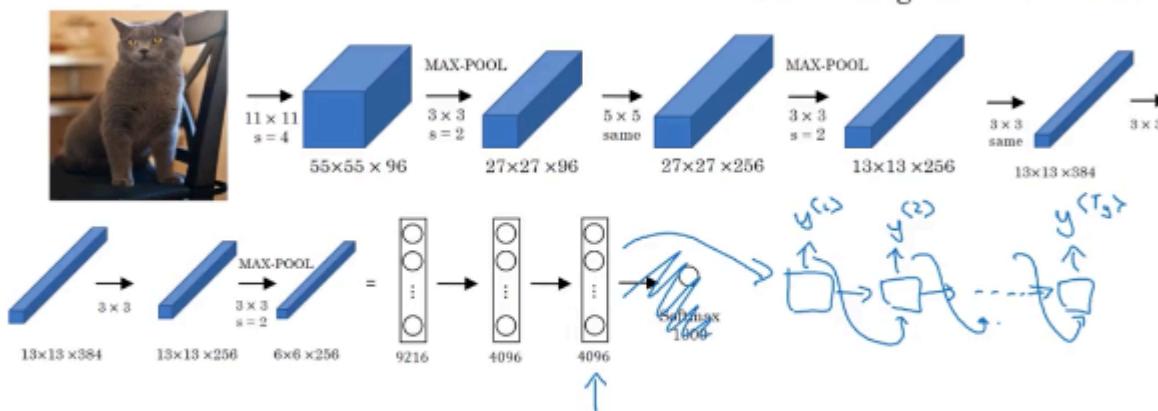
[Sutskever et al., 2014. Sequence to sequence learning with neural networks] ↩

[Cho et al., 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation] ↩

Andrew Ng

Image captioning

$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>} \quad \dots \quad \{$
 A cat sitting on a chair }



[Mao et. al., 2014. Deep captioning with multimodal recurrent neural networks]

[Vinyals et. al., 2014. Show and tell: Neural image caption generator]

[Karpathy and Fei Fei, 2015. Deep visual-semantic alignments for generating image descriptions]

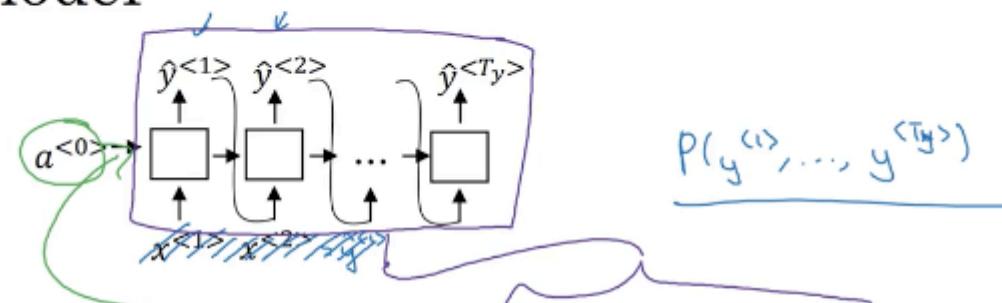
Andrew Ng

works well especially if caption wanted is not that long

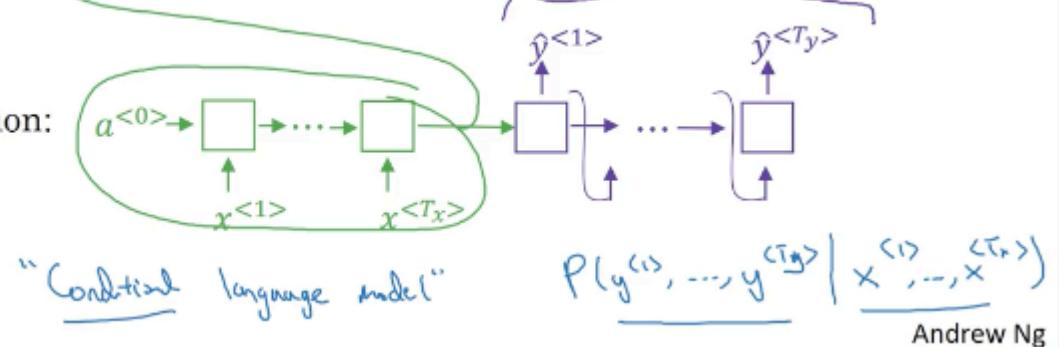
- there were multiple groups that were working on similar things at the time

Machine translation as building a conditional language model

Language model:



Machine translation:



understand the part which makes it 'conditional'

Finding the most likely translation

Jane visite l'Afrique en septembre.

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

- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.
- In September, Jane will visit Africa.
- Her African friend welcomed Jane in September.

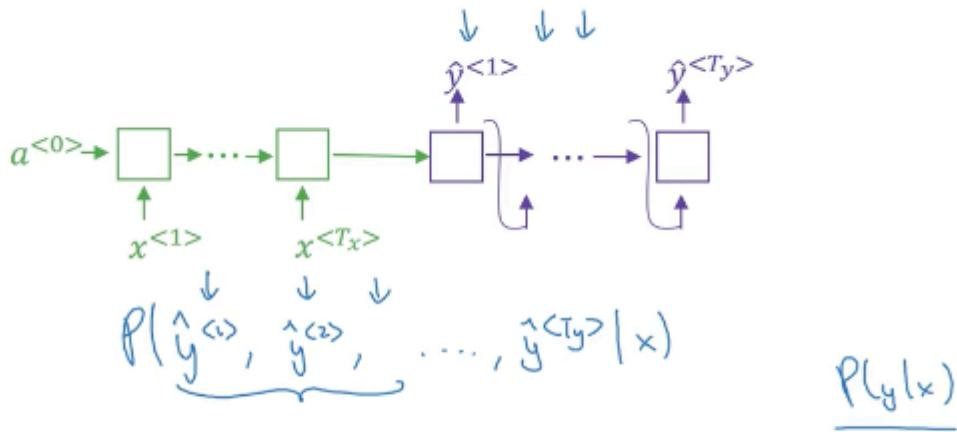
$$\arg \max_{y^{<1>}, \dots, y^{<T_y>}} P(y^{<1>}, \dots, y^{<T_y>} | x)$$

Andrew Ng

find an english sentence y that maximizes, we do not want to sample it.

Why not a greedy search?

$$P(\hat{y}^{(1)} | x)$$



- Jane is visiting Africa in September.
 - Jane is going to be visiting Africa in September.
- $P(\text{Jane is going } | x) > P(\text{Jane is visit } | x)$

Andrew Ng

- Pick best word for each unit (highest probability) - greedy search -> does not really work
- going is a more general word, so for example, program may say that second sentence, for the third word, 'going' is the most likely choice - which ultimately results in a worse sentence for translation

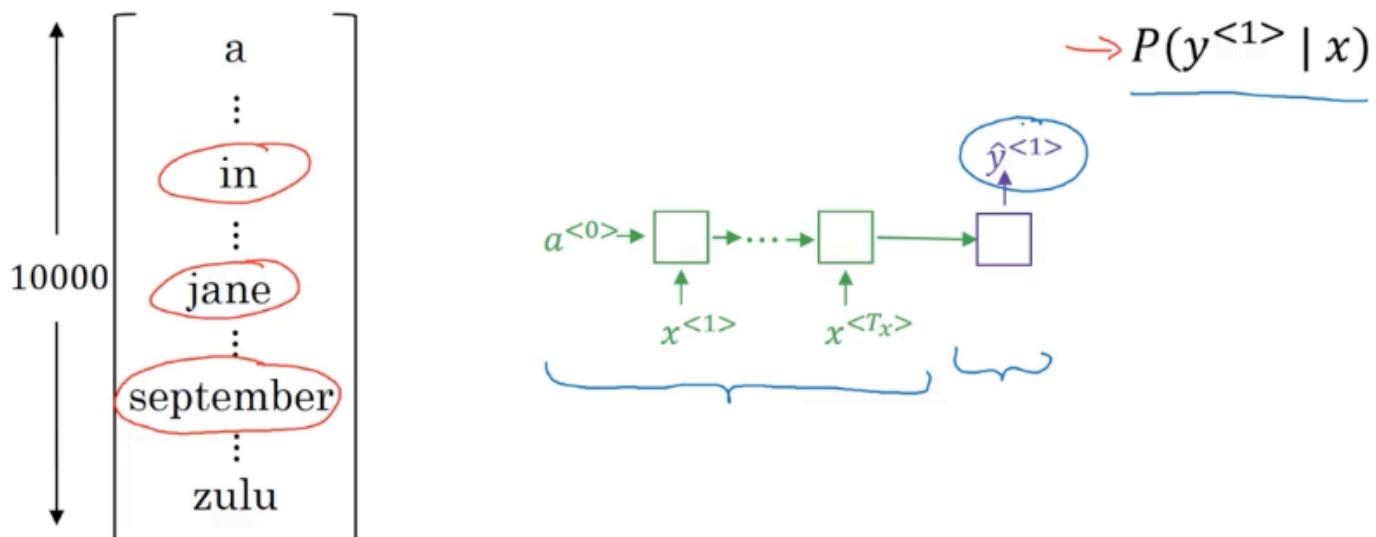
There are numerous possible sentences, for example for 10 word sentences, so we use approximate search algorithms (10k vocab, 10 words = $10^k \times 10$ possible sentences)

Beam Search

Beam search algorithm

$$\underline{B=3} \quad (\text{beam width})$$

Step 1

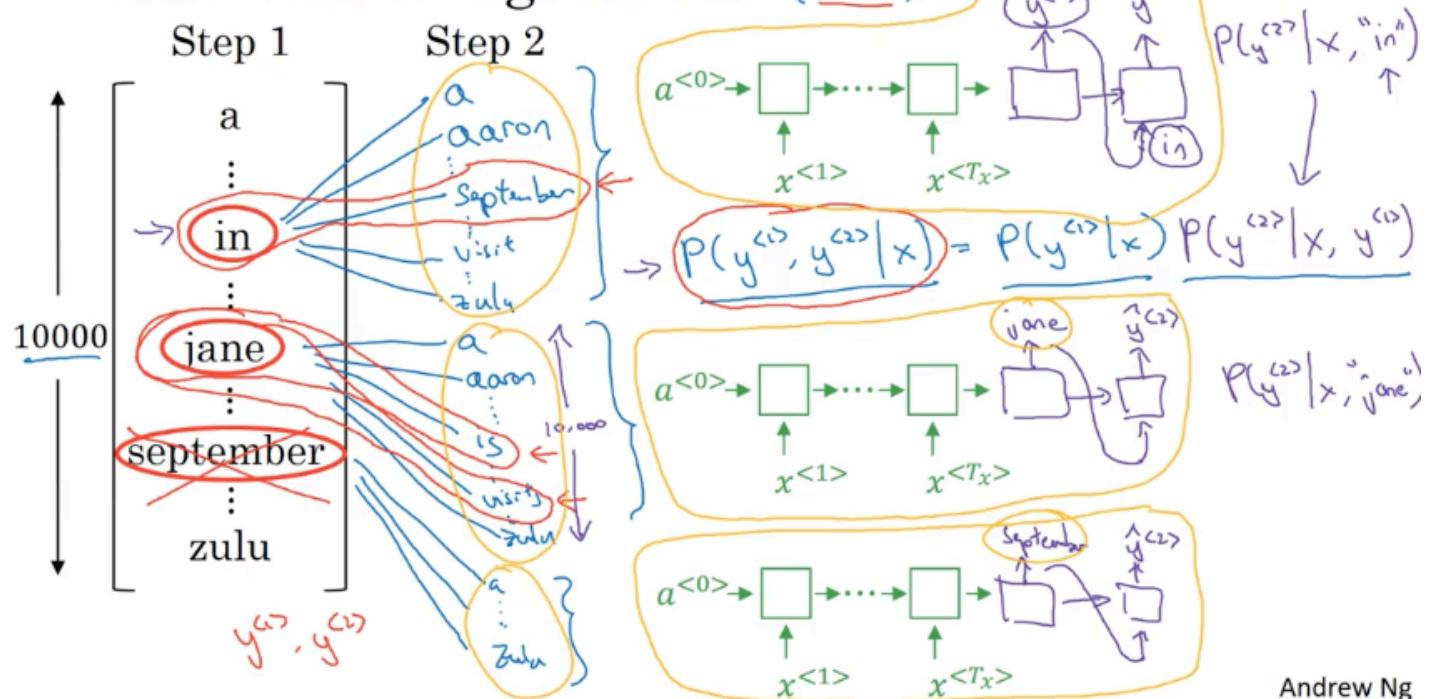


Andrew Ng

- try to pick the first word in translated sentence

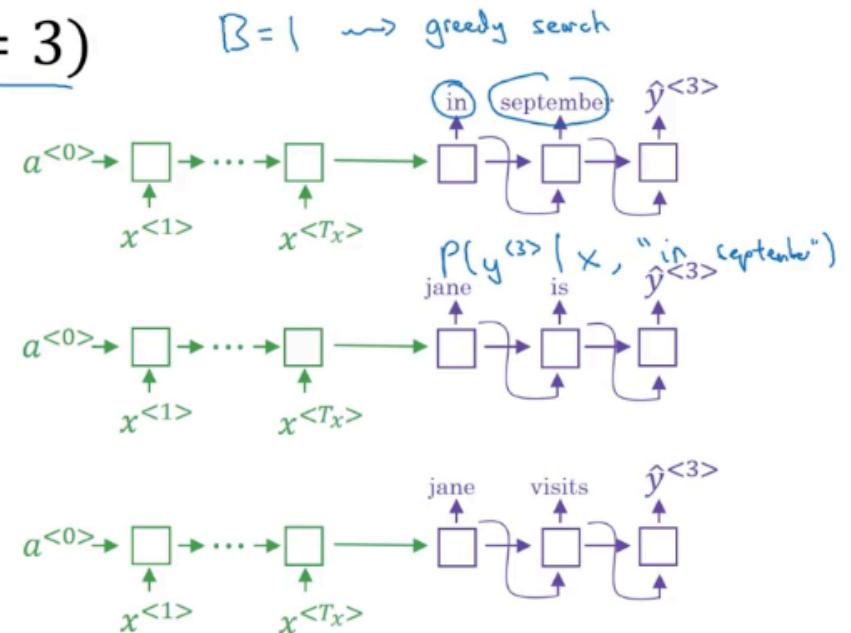
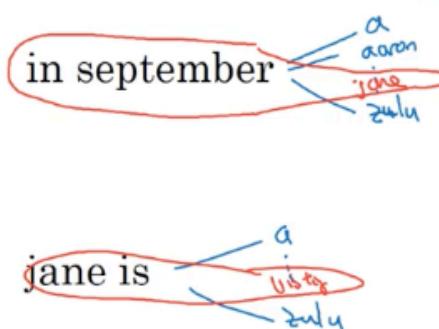
- has parameter B - beam width - can search for several possibilities at once, different than other search algorithms (let's say finds 3 words, chooses them to continue to compute them - hold each for the first word)

Beam search algorithm ($B = 3$)



- now for each of these choices, find the best second choice, dependent on the first choice
- during iteration, eliminates some choices of branches
- don't instantiate all possible outputs, just choose a number of them

Beam search ($B = 3$)



$$P(y^{<1>} | x)$$

jane visits africa in september. <EOS>

Andrew Ng

- and iterate

Refinements to Beam Search

Length Normalization

may get a *numerical underflow* - multiplying many numbers that are very small: 0.0001

- soln: take logs, more stable output values.
- another change to make it work even better: long sentence \rightarrow probability is very low. normalize it by the number of words in the sentence, this significantly reduces the penalty. also there are different/softer adjustments (ex: adding another hyperparameter, alpha).

Length normalization

$$\rho(y^{(1)} \dots y^{(T_y)} | x) = \frac{P(y^{(1)} | x) P(y^{(2)} | x, y^{(1)}) \dots}{P(y^{(T_y)} | x, y^{(1)}, \dots, y^{(T_y-1)})}$$

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

log

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

$T_y = 1, 2, 3, \dots, 30.$

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

$\alpha = 0.7$

$\frac{\alpha}{\alpha-1}$

$\alpha = 1$

$\alpha = 0$

Andrew Ng

How Big Value For Beam Width B?

Beam search discussion

Beam width B?

$1 \rightarrow 3 \rightarrow 10, 100, 1000, \rightarrow 3000$

large B: better result, slower
small B: worse result, faster

Unlike exact search algorithms like BFS (Breadth First Search) or DFS (Depth First Search), Beam Search runs faster but is not guaranteed to find exact maximum for $\arg \max_y P(y|x)$.

Andrew Ng

Error Analysis in Beam Search

This is an heuristic algorithm, so it does not always find the most optimal/accurate result, so how to apply error analysis here?

- got a result (bad), who is to blame? increase the B , never hurts. \rightarrow might not solve. When to decide take time to try to make the RNN part better?
- compute both probabilities, p^* and \hat{p} , both can be bigger than other, dependent on the result you may understand the problem better.

Example

Jane visite l'Afrique en septembre.

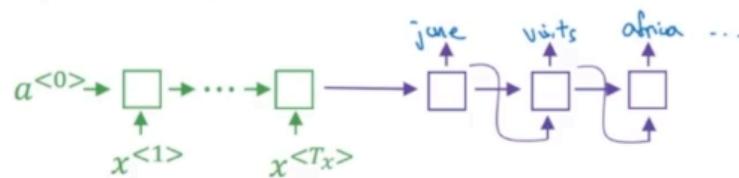
\rightarrow RNN
 \rightarrow Beam Search

BT

Human: Jane visits Africa in September. (y^*)

Algorithm: Jane visited Africa last September. (\hat{y}) \leftarrow

$$\text{RNN computes } P(y^*|x) \geq P(\hat{y}|x)$$



Andrew Ng

Logic:

Error analysis on beam search

$$P(y^*|x)$$

$$P(\hat{y}|x)$$

Human: Jane visits Africa in September. (y^*)

Algorithm: Jane visited Africa last September. (\hat{y})

Case 1: $P(y^*|x) > P(\hat{y}|x)$ \leftarrow arg max_y $P(y|x)$

Beam search chose \hat{y} . But y^* attains higher $P(y|x)$.

Conclusion: Beam search is at fault.

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

y^* is a better translation than \hat{y} . But RNN predicted $P(y^*|x) \leq P(\hat{y}|x)$.

Conclusion: RNN model is at fault.

Andrew Ng

Error analysis process

Human	Algorithm	$P(y^* x)$	$P(\hat{y} x)$	At fault?
Jane visits Africa in September.	Jane visited Africa last September.	$\underline{2 \times 10^{-10}}$	$\underline{1 \times 10^{-10}}$	B
...	...	—	—	R
...	...	—	—	R
			—	R

Figures out what fraction of errors are “due to” beam search vs. RNN model

Andrew Ng

Bleu Score

The BLEU (Bilingual Evaluation Understudy) score is a metric used to evaluate machine translation systems, where multiple good translations can exist for a given input. It compares the machine-generated translation to one or more human-generated reference translations. BLEU calculates a modified precision by counting n-gram matches (unigrams, bigrams, etc.) between the generated and reference translations, clipping counts to the maximum frequency in any reference. These precisions (P_1, P_2, P_3, P_4) are averaged, exponentiated, and adjusted by a brevity penalty (BP) to penalize overly short translations. While not perfect, the BLEU score provides a single real-number metric that has significantly aided machine translation development. It's also used in other text generation tasks like image captioning.

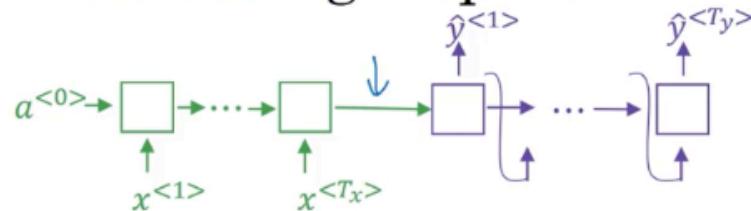
Attention Model

modification of attention,
already been using encoder - decoder pair,
the change:

- the remembering part is not been done properly for the long sequences, the connection between the encoder and decoder is too thin for such job.

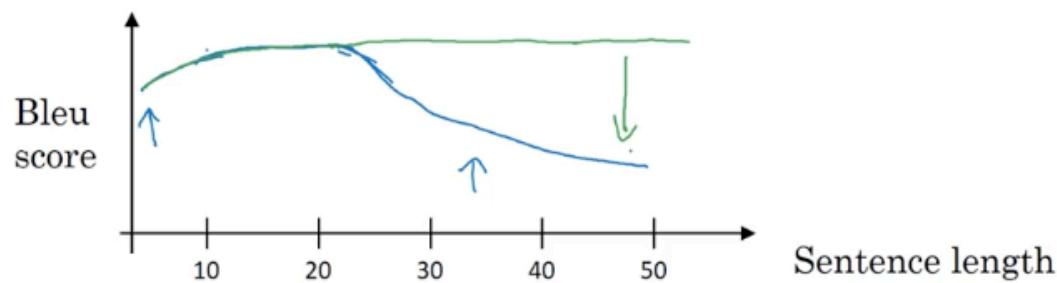
- think how a human does it, processes it part by part, memorizing and processing at once is hard.

The problem of long sequences



Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

Jane went to Africa last September, and enjoyed the culture and met many wonderful people; she came back raving about how wonderful her trip was, and is tempting me to go too.



Andrew Ng

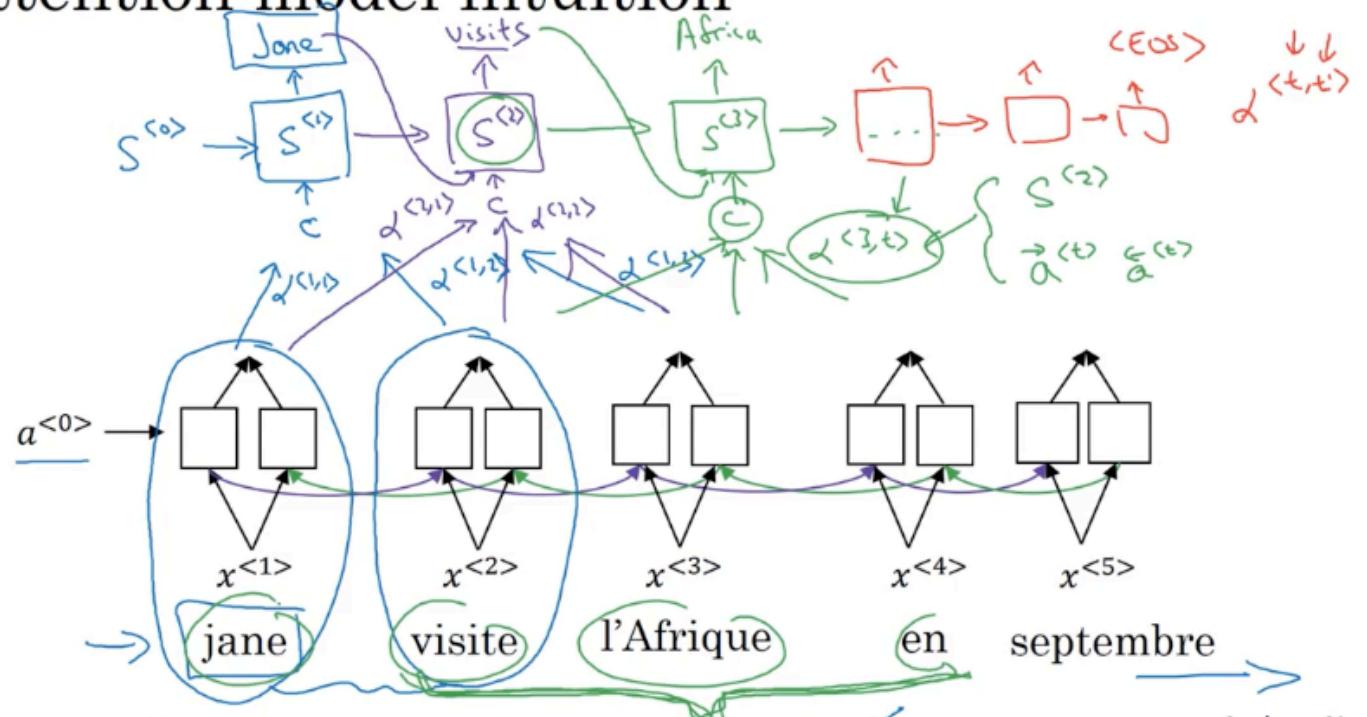
so it is an hard task for the model to remember long sequences

Attention Model Intuition

first developed for machine translation, later use of it generalized (a lot).

how much attention should we pay for each word in the sentence?

Attention model intuition



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

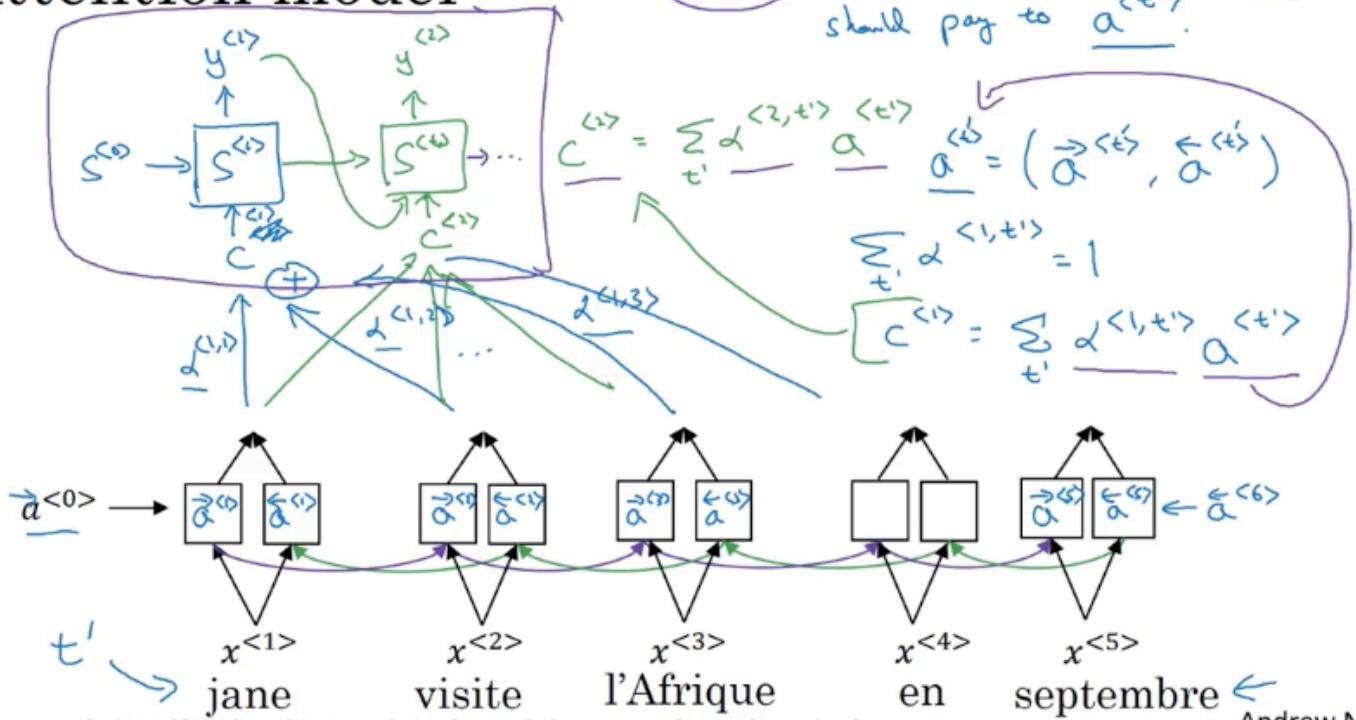
Andrew Ng

for t , how much attention is needed for t' (top right corner)

Attention Model

model pays attention pays attention to part of the input (more to some, less to some)

Attention model



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

Andrew Ng

upper part seems like a regular RNN structure

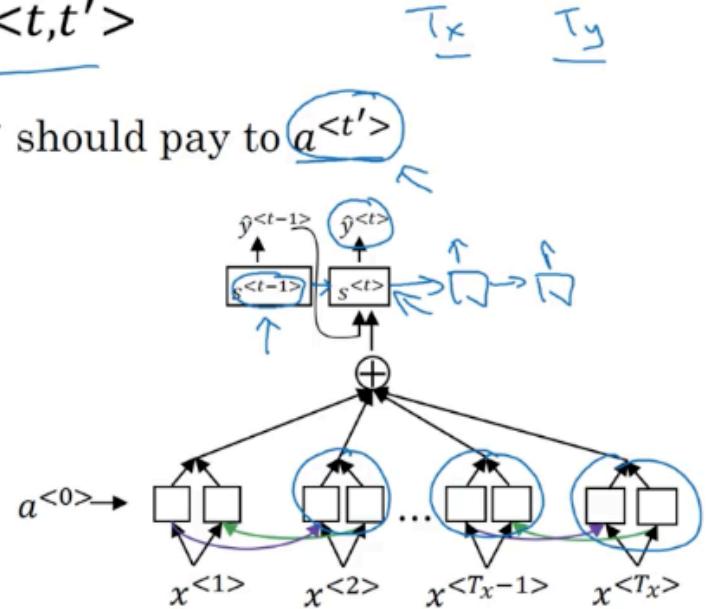
how to define the attention weights (alpha value in purple):

Computing attention $\underline{\alpha^{<t,t'>}}$

$\alpha^{<t,t'>} = \text{amount of attention } y^{<t>} \text{ should pay to } \underline{a^{<t'>}}$

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

$$e^{<t,t'>} = \frac{s^{<t-1>} \cdot a^{<t'>}}{\alpha^{<t'>}}$$



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

[Xu et. al., 2015. Show, attend and tell: Neural image caption generation with visual attention]

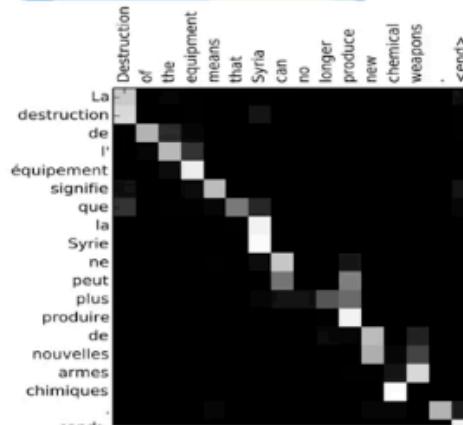
Andrew Ng

Attention examples

July 20th 1969 → 1969 – 07 – 20

23 April, 1564 → 1564 – 04 – 23

Visualization of $\alpha^{t,t'}$:



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

Andrew Ng

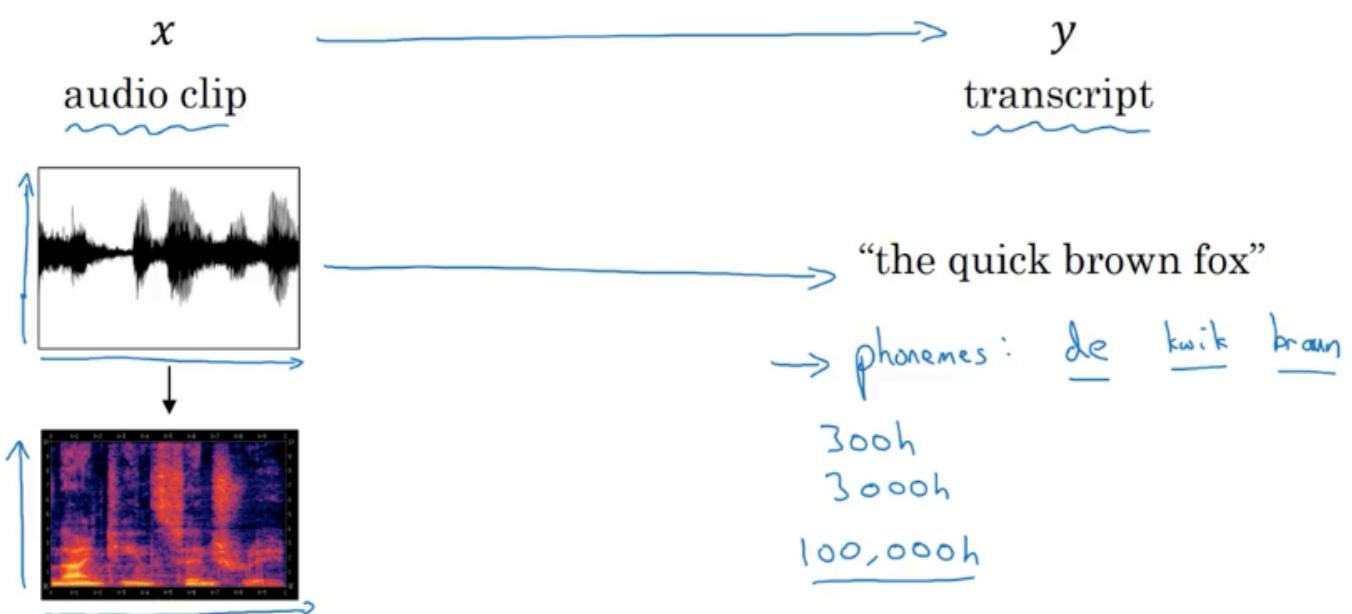
Speech Recognition

One of the best things about sequence to sequence modeling: highly accurate speech recognition

little changes in air pressure - ear, and microphone senses this

- common preprocessing step: process raw audio, generate a spectrogram. how loud is this audio for different frequencies. (human ear also does something that is similar to this)
- at first it was general to transcribe audio with phonemes. now it is not necessary. (or any other hand engineered representations)
- couple thousands → many research papers are written on such amounts of data.
- couple hundred thousands → commercial use.

Speech recognition problem



Andrew Ng

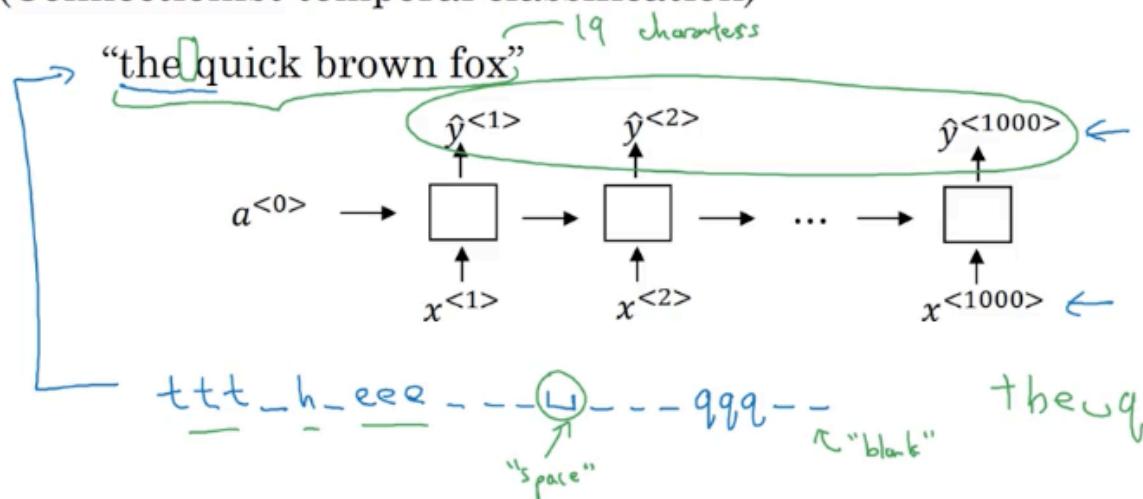
- you may use attention mechanism
- you may use CTC cost

generally # of time steps in the audio is much bigger than the output. (10 hz x 10 seconds, output is not even near it.)

soln:

CTC cost for speech recognition

(Connectionist temporal classification)



Basic rule: collapse repeated characters not separated by “blank”

[Graves et al., 2006. Connectionist Temporal Classification: Labeling unsegmented sequence data with recurrent neural networks] Andrew Ng

Trigger Word Detection

things you can wake up with your voice

Examples:

What is trigger word detection?



Amazon Echo
(Alexa)



Baidu DuerOS
(xiaodunihao)



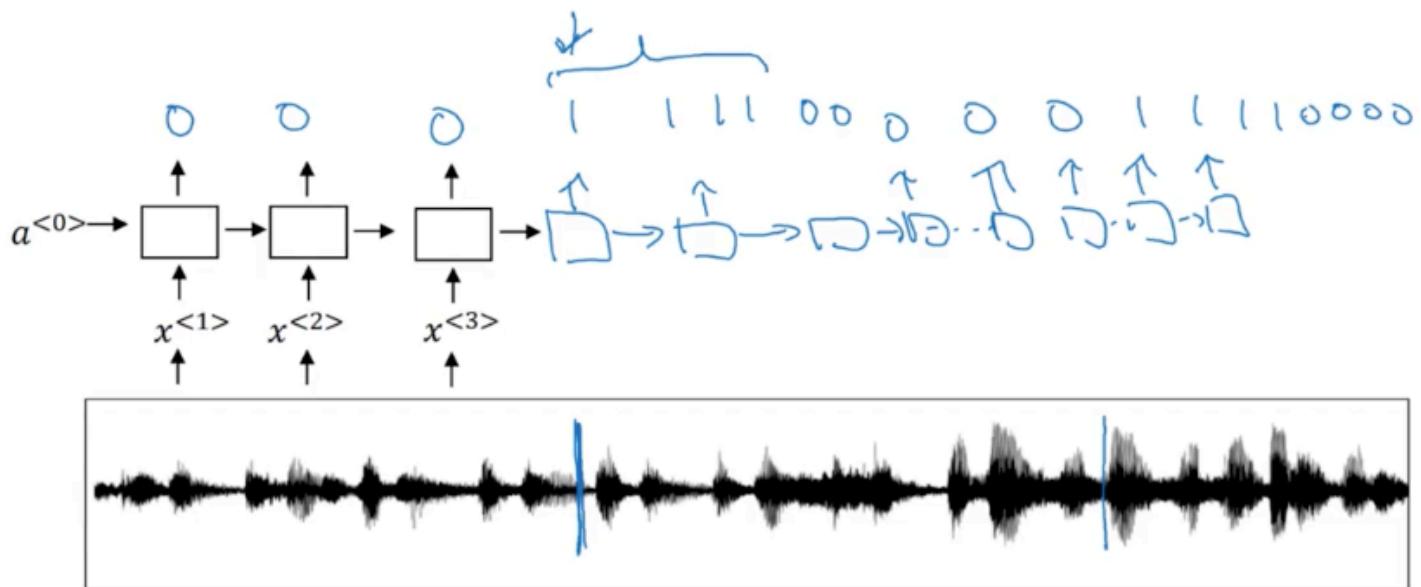
Apple Siri
(Hey Siri)



Google Home
(Okay Google)

Example algorithm:

Trigger word detection algorithm



Andrew Ng

instead of single 1, may put couple of 1's for trigger instances. - a little hack - Andrew Ng

- Data synthesis is an effective way to create a large training set for speech problems, specifically trigger word detection.
 - Using a spectrogram and optionally a 1D conv layer is a common pre-processing step prior to passing audio data to an RNN, GRU or LSTM.
 - An end-to-end deep learning approach can be used to build a very effective trigger word detection system.

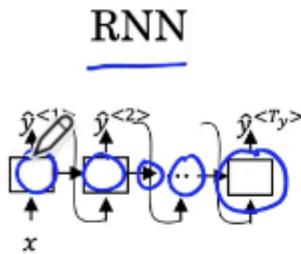
Transformer Network

Best models are often on transformer architecture

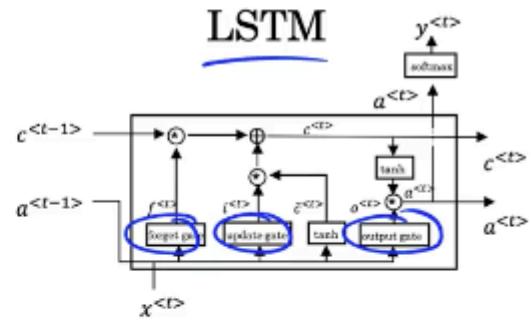
New and SOTA

Transformer Network Motivation

increased complexity
sequential



GRU



Andrew Ng

each unit is in essence a bottleneck by itself

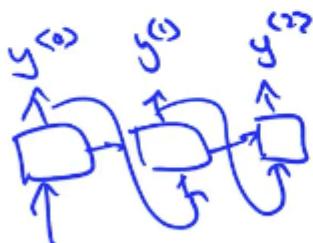
- transformer solves it by enabling to process everything at once.

Transformer Network Intuition

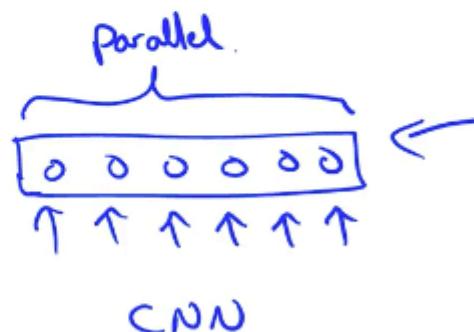
• Attention + CNN

• Self-Attention

• Multi-Head Attention



$A^{<1>} A^{<2>} A^{<3>} A^{<4>} A^{<5>}$



[Vaswani et al. 2017, Attention Is All You Need]

Andrew Ng

Self-Attention

Self-Attention Intuition

$A(q, K, V)$ = attention-based vector representation of a word

calculate for each word $A^{<1>} \dots A^{<s>}$

RNN Attention

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

\downarrow \downarrow \downarrow
 $x^{<1>} \quad x^{<2>} \quad x^{<3>}$
Jane visite l'Afrique

Transformers Attention

$$A(q, K, V) = \sum_i \frac{\exp(q \cdot k^{<i>})}{\sum_j \exp(q \cdot k^{<j>})} v^{<i>}$$

\downarrow \downarrow \downarrow
 $x^{<4>} \quad x^{<5>} \quad v^{<3>}$
en septembre

[Vaswani et al. 2017, Attention Is All You Need]

Andrew Ng

how you use some word may affect how you want to represent it - continent, destination, place both have softmax parts

main difference: for every word, have 3 values called query, key, value = attention value

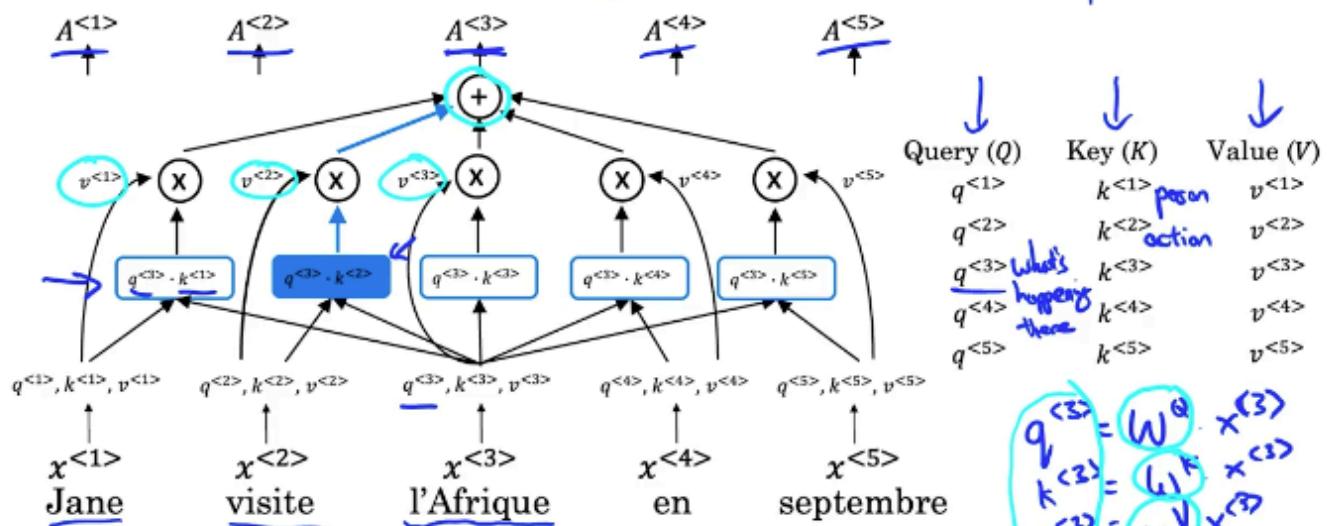
Self Attention Representation $A^{<3>}$

- $q^{<3>}$ question you get to ask - what's happening there?
- $k^{<1>}$ how good is the answer to this question - (this is repeated for each word in the sentence: $k^{<2,3,4..>}$)
- maybe $k^{<2>}$ is most relevant: highest score (let's say)

Self-Attention

$$A(q, K, V) = \sum_i \frac{\exp(q \cdot k^{<i>})}{\sum_j \exp(q \cdot k^{<j>})} v^{<i>}$$

$$\underline{\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V}$$



[Vaswani et al. 2017, Attention Is All You Need]

$$\begin{aligned} q^{<3>} &= W^Q \cdot x^{<3>} \\ k^{<3>} &= W^K \cdot x^{<3>} \\ v^{<3>} &= W^V \cdot x^{<3>} \end{aligned}$$

Andrew Ng

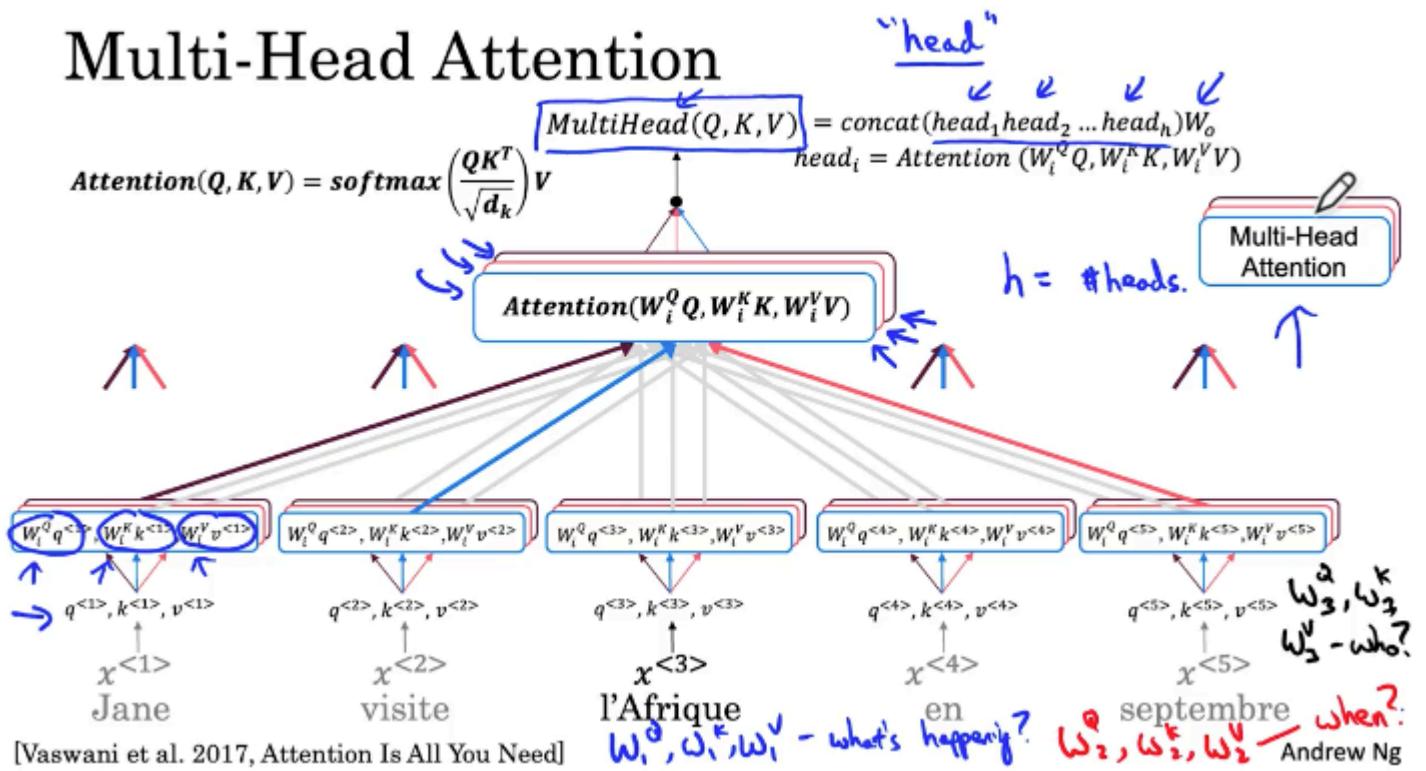
square root dk prevents it from exploding

Multi-Head Attention

A big for loop over the self-attention mechanism

- first question: what's happening?
- now there are other questions - other heads: when's happening? - happening to who? etc. (imaginary)

Multi-Head Attention



Multi-Head Attention

- **Concept:** Computes self-attention multiple times ("heads") in parallel, each with different learned weight matrices (W^Q, W^K, W^V). Each head, can be thought of answering a question, such as who, when, and what's happening.
- **Process:**
 1. For each head i :
 - Linearly project input Queries (Q), Keys (K), and Values (V) using learned matrices W_i^Q , W_i^K , W_i^V .
 - Compute self-attention: $\text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$
 2. Concatenate the outputs of all heads.
 3. Multiply the concatenated output by a final weight matrix W^O .

- **Formula:**

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)W^O$$

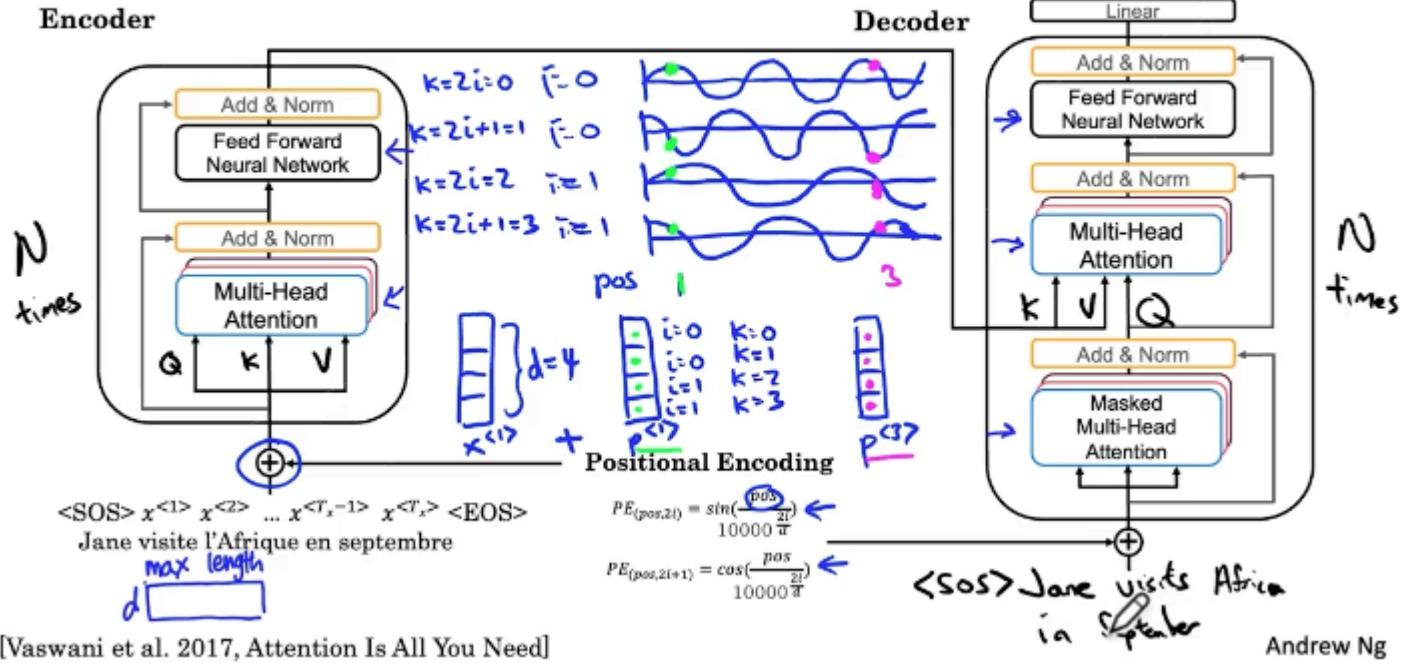
$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

- **Benefit:** Allows the model to attend to different aspects of the input sequence simultaneously (e.g., "what's happening," "when," "who"), capturing richer contextual information.
- The number of heads is represented by h .
- The computations of all heads can occur in parallel.

The multi-head self-attention calculation is identical to standard self-attention; the difference is that it is performed multiple times with distinct, learned linear projections of Q, K, and V.

Transformer Network

Transformer Details



Refer to the attached figure for the overall architecture of the transformer (from “Attention Is All You Need”), which includes both an **encoder** and a **decoder**. Below is a brief walkthrough of how it all fits together:

1. Encoder

- **Embeddings + Positional Encoding:** Each word/token is turned into an embedding and **positional encodings** (sine/cosine signals) are added to capture word order.
- **Multi-Head Self-Attention:** The encoder uses **Q, K, V** matrices (all derived from the embeddings) to allow each word to attend to other words. Having multiple heads enables the model to focus on different relationships in parallel.
- **Feed-Forward Network:** A fully connected layer processes the attention outputs.
- **Add & Norm:** Residual connections and a normalization layer (akin to batch normalization) help speed up learning and preserve positional information.
- The encoder block is typically repeated **N times** (e.g., 6).

2. Decoder

- The goal of the decoder is to produce the target sequence (e.g., an English translation).
- **Masked Multi-Head Attention:** The decoder first attends to its own past outputs (with future positions “masked” in training so the model can’t “peek ahead”).
- **Multi-Head Attention:** Next, the decoder attends to the **encoder’s output** (queries from the decoder and K, V from the encoder).
- **Feed-Forward + Add & Norm:** Similar to the encoder, the attention output is passed through a feed-forward layer and normalization.
- This decoder block is also repeated **N times**, and finally feeds into a **linear + softmax** layer to predict the next token.

3. Training & Masking

- **Masking** is applied in the decoder's first attention layer to simulate predicting one word at a time, even though the entire correct sequence is known during training.
- This encourages the model to learn how to generate tokens step-by-step without cheating by looking ahead.

4. Takeaways

- The **encoder** captures contextual embeddings for the entire input sequence.
- The **decoder** uses its current output tokens (so far) to query the encoder's representation and generate the next token, step by step.
- **Positional encoding** ensures the model knows the relative position of words.
- **Residual connections** and **normalization** help stabilize and accelerate training.

Use these concepts (e.g., attention mechanisms, positional encodings, mask during training) as the basis for understanding transformer variants such as **BERT** and **DistilBERT**.

Q = interesting questions about the words in a sentence

K = qualities of words given a Q

V = specific representations of words given a Q

- The combination of self-attention and convolutional network layers allows for parallelization of training and *faster training*.
 - Self-attention is calculated using the generated query Q , key K , and value V matrices.
 - Adding positional encoding to word embeddings is an effective way to include sequence information in self-attention calculations.
 - Multi-head attention can help detect multiple features in your sentence.
 - Masking stops the model from 'looking ahead' during training, or weighting zeroes too much when processing cropped sentences.
-

References

Week 1:

- [Minimal character-level language model with a Vanilla Recurrent Neural Network, in Python/numpy](#) (GitHub: karpathy)
- [The Unreasonable Effectiveness of Recurrent Neural Networks](#) (Andrej Karpathy blog, 2015)
- [deepjazz](#) (GitHub: jisungk)
- [Learning Jazz Grammars](#) (Gillick, Tang & Keller, 2010)
- [A Grammatical Approach to Automatic Improvisation](#) (Keller & Morrison, 2007)
- [Surprising Harmonies](#) (Pachet, 1999)

Week 2:

- [Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings](#) (Bolukbasi, Chang, Zou, Saligrama & Kalai, 2016)
- [GloVe: Global Vectors for Word Representation](#) (Pennington, Socher & Manning, 2014)
- [Woebot.](#)

Week 4:

- [Natural Language Processing Specialization](#) (by [DeepLearning.AI](#))
- [Attention Is All You Need](#) (Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser & Polosukhin, 2017)