



deeplearning.ai

5 Sequence models

Andrew Ng

Note by:
Hossein Mamaghanian
amirhm@gmail.com



"Sequence Models"

"Notations"

Imagine an example of input sequence (a sentence here), where we want to find the "names"

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

$x^{(1)} \quad x^{(2)} \quad x^{(3)}$ \dots $x^{(9)}$
 $T_x = 9$

$y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(9)}$
 $T_y = 9$

$\langle \cdot \rangle$ to refer the index in sequence ,

$x^{(i)(t)}$: index t in training example i

T_x : length of sequence, $T_x^{(i)}$: length of i^{th} training example.

here T_x and T_y are equal but not necessarily always same

"Representing words"

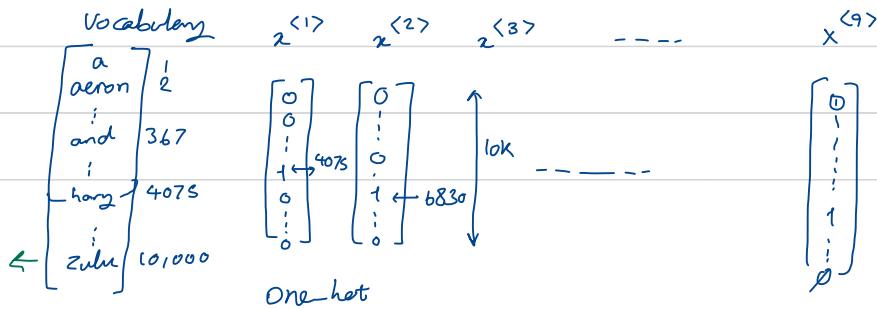
we use a dictionary

in commercial NLP

Dic. size 30w50k

Maybe even +M

x: harry Potter and Hermione...



One-hot : because only one index is non zero
(1) $\in \{0, 1\}^{K \times 9}$

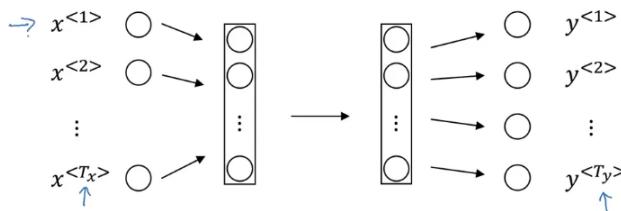
it's solved as a supervised $x \rightarrow y$ learning , input $X \in \mathbb{R}^n$

if a word doesn't exist in our vocabulary we could assign by a word wnv ,

→ how to learn such mapping ?

↓

Why not a standard network?



such naive net.

does n't work well

Problems:

- Inputs, outputs can be different lengths in different examples. → even Padding - is not a good idea
- Doesn't share features learned across different positions of text. idea

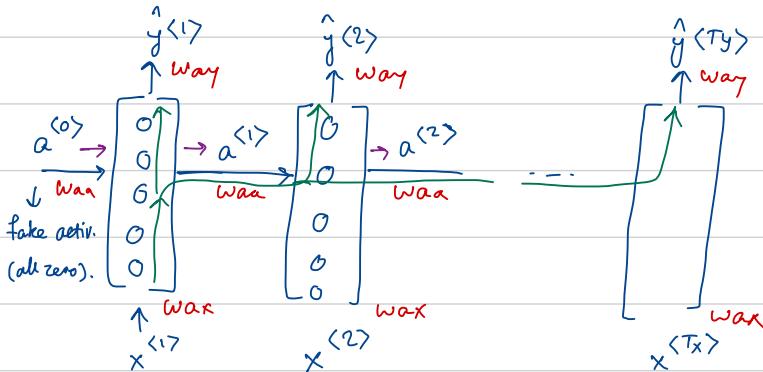
↳ the main problem, with such model for sequence models,

→ for example in such Net , it is not easy to generalize that some word in other positions also could appear, (similar to Conv. Nets,).

→ Input size also huge, weight for layers could be enormous.

→ To overcome these problems, we use an architecture which is called "RNN" Recurrent neural networks"

"Recurrent Neural Networks"



→ where inputs are given in sequential order.

→ This is for case $T_x = T_y$
slightly different if not.

→ every output $\hat{y}^{(t)}$ is calculated

/activation are feed to steps after.

based on info from previous sequence.

weight are same for all layers.

→ This architecture does not take into account the words after which is important (for this we could have BRNN : Bi-directional RNN).

for example:

he said "Teddy Roosevelt was a great President

"Teddy Bears are on sale"

To decide if Teddy is part of name we need info from later

Sometime RNN are represented as:



"forward Propagation"

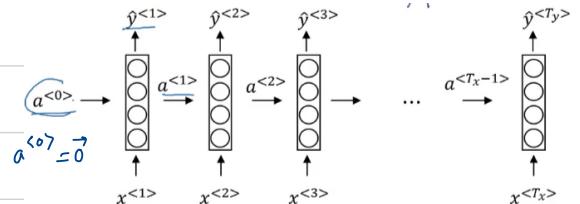
$$a^{(1)} = g_1(W_{aa} a^{(0)} + W_{ax} x^{(1)} + b_a)$$

$$y^{(1)} = g_2(W_y a^{(1)} + b_y)$$

⋮

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

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



activation function Tanh / ReLU

→ sigmoid (if binary)
More common

softmax ...

Notations: $a \leftarrow \underbrace{W}_{\text{R}} \underbrace{a \times X}_{\text{C}}$

in case of more common Tanh other methods are taken to overcome problems like vanishing gradients.

"Simplified Notations"

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

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

More common representation.

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

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

→ W_a is the concatenation of W_{aa} , W_{ax}

if a size is 100 → $W_{aa} \in 100 \times 100$

$W_{ax} \in 100 \times 10,000$

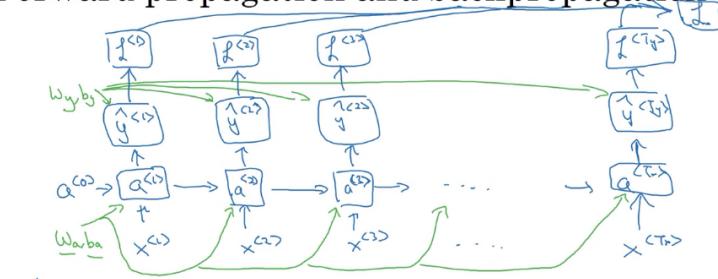
$$W_a = \begin{bmatrix} W_{aa} & W_{ax} \\ \vdots & \vdots \end{bmatrix} = \begin{bmatrix} a^{(t-1)} & x^{(t)} \end{bmatrix} \in \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix}$$

"Back propagation"

logistic-loss or

cross entropy loss

Forward propagation and backpropagation



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

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

← logistic loss (cross entropy loss)

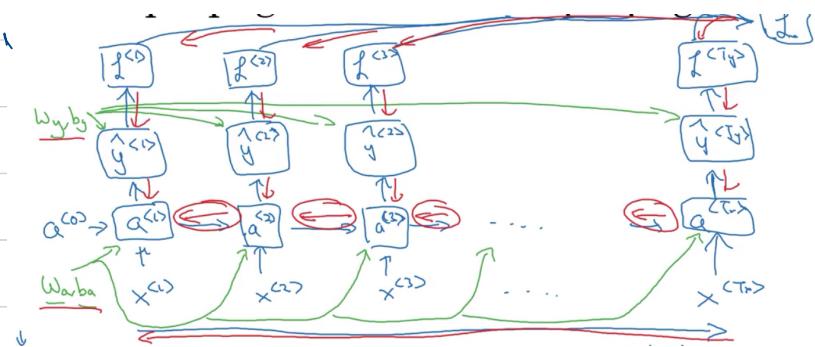
Andrew Ni

in forward computation for every timestamp t , we want to calculate

$$a^{(t)}, x^{(t)} \xrightarrow{\text{calc}} a^{(t)} \quad \rightarrow \text{using } w_a, b_a$$

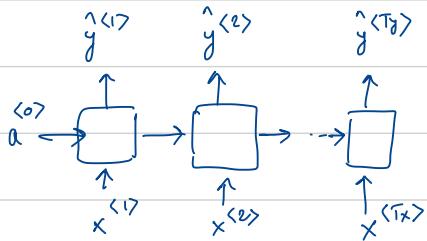
$$\text{and } a^{(t)} \xrightarrow{\text{calc}} y^{(t)} \quad \text{using } w_y, b_y$$

In Back propagation time, error is propagated backward, since it is in the reverse direction of timestamp, it is called "B.P. through time"



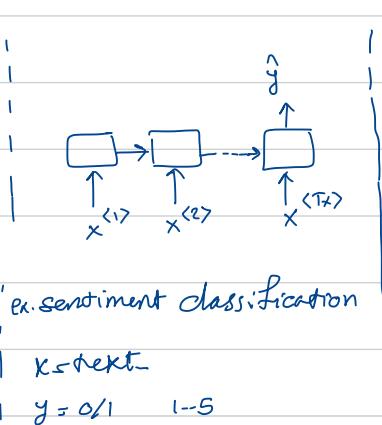
"RNN types"

Many-to-Many



$T_x = T_y$

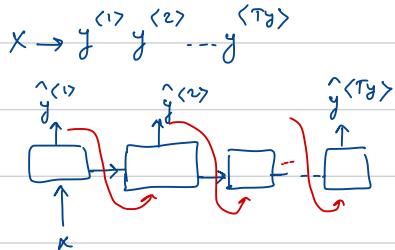
Many-to-One



One-to-One



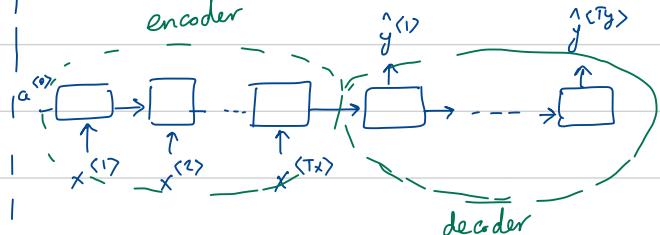
One-to-many



in reality, we use output of previous as input for next.

Many-to-Many (2nd version $T_x \neq T_y$)

ex: Machine translation



encoder reads the input and decoder generates the output.

"language Model"

→ Model of language which gives a probability for each sentence.

$$\left\{ \begin{array}{l} \text{The apple and pair salad} \\ \text{v . . . pear salad } \checkmark \text{ (More probable)} \end{array} \right. \quad P(\dots) = 3.2 \times 10^{-13}$$
$$P(y^{<1>} , y^{<2>} , \dots , y^{<Ty>})$$

language modelling with an RNN.

→ Training set: large corpus of english text.

→ Step 1: Tokenize.

Cats average 15 hours of sleep a day. <EOS>

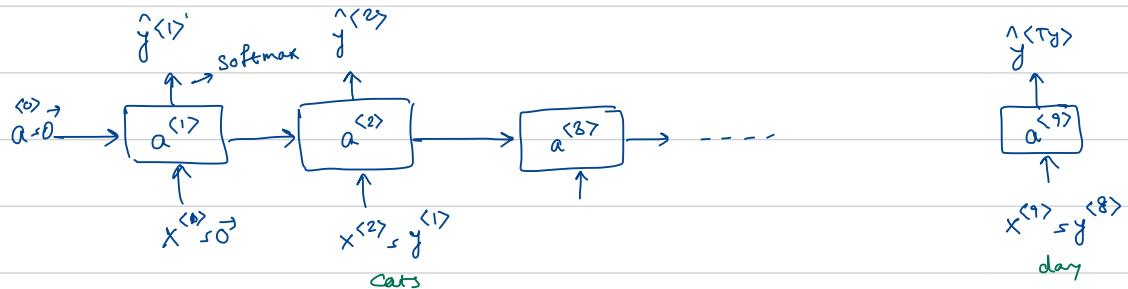
from dic → $y^{<1>} \quad y^{<2>} \quad \dots \dots$

The egyptian Man is a bread of cat <EOS>

→ if word does not exist in vocabulary, replaced by <UNK>

→ another word <EOS> is sometimes used to mark end of sentence.

"RNN model"



$\hat{y}^{(1)}$ is actually the probability of every word in dic to be first word

$$P(a) \cdot P(\text{aeron}) \quad \dots \quad P(\text{cats}) \quad \dots \quad P(\text{zulu}) \cdot P(\text{EOS})$$

$\hat{y}^{(2)}$ $P(\underline{\hspace{2cm}} / \text{cats})$: average

$\hat{y}^{(3)}$ $P(\underline{\hspace{2cm}} / \text{"cats average"})$: 15

: This RNN in every timestep

learn to guess the next word.

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

(softmax cost)

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

→

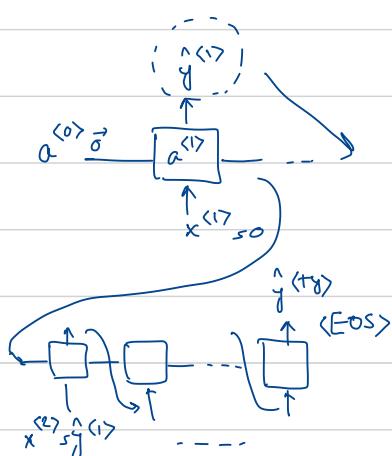
if network correctly learned, given previous words could predict the next words.

$$P(y^{(1)}, y^{(2)}, y^{(3)}) = \underbrace{P(y^{(1)})}_{\text{1st softmax}} \cdot \underbrace{P(y^{(2)} | y^{(1)})}_{\text{2nd softmax}} \cdot \underbrace{P(y^{(3)} | y^{(1)}, y^{(2)})}_{\text{3rd Timestamp softmax}}$$

Sampling a sequence from a trained RNN

→ "to verify or check the model"

for sampling



first step $y^{(1)}$ is the distribution of every word as first word in sentence.

then we can randomly sample based on this distribution. (np.random.choice)

Then this word, will be given to next time-stamp and so on to generate sentence. until we see $\langle \text{EOS} \rangle$.

Character-level language model

→ similarly, it is possible to have character-level language model.

→ vocabulary: [a, b, --, , , , ;, o, --, A, -- Z], $y^{(1)}, y^{(2)}, \dots$ are characters!

cong { → word-level is superior due to better capturing the relation
- slip between word in the sentence.
- computationally heavy

example generated texts

Sequence generation

pros: { No unknown word
{ Not wide-use.

News

President enrique pena 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.

"Vanishing gradient with RNN"

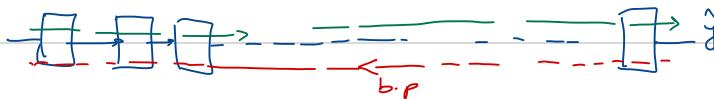
Take example of sentence in NLP, sometimes the words has very long relationship.

→ The cat, which already ate----, was full

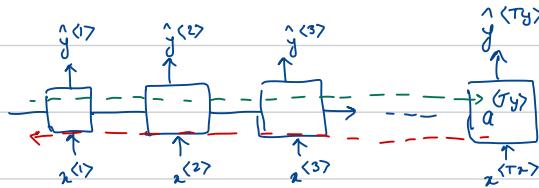
The cats, a c e t .. . w e r l l u full

The basic RNN structure that seen so far is not very good at capturing these far dependencies.

similar to problem in very deep neural networks, vanishing gradient is problem here



similar in RNNs.



The gradient in b.p has
very hard time to affect

The main weakness

The effect of each output will be

earlier layers.

of basic RNNs

very limited to close neighbours. and will not
travel much far to affect earlier layers.

exploding gradients: another problem which could exist in RNNs, similarly

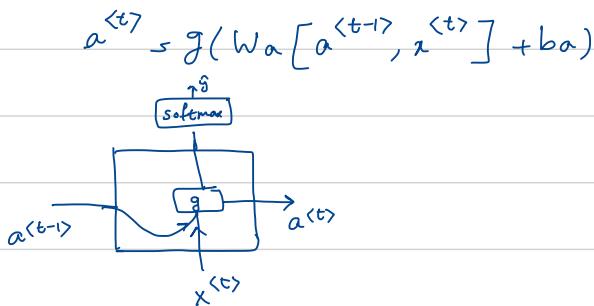
but this is not the main challenges in training RNNs.

→ in this case we see "NaN" in computations.

we could use gradient clipping: put some threshold on gradient
if bigger, rescale or clip. (relatively robust).

"gated recurrent unit" (simplified)

→ modification to RNN hidden unit that make it much better in capturing long range relationship and help with "vanishing gradient problem"



c = memory cell.

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

↳ for gru, c is output a value, this notation is used

more for later describing LSTM

candidate for update $c^{<t>}$

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

activation

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

gate unit. sigmoid

↳ gate unit will decide whether we update $\hat{c}^{<t>}$ or not.

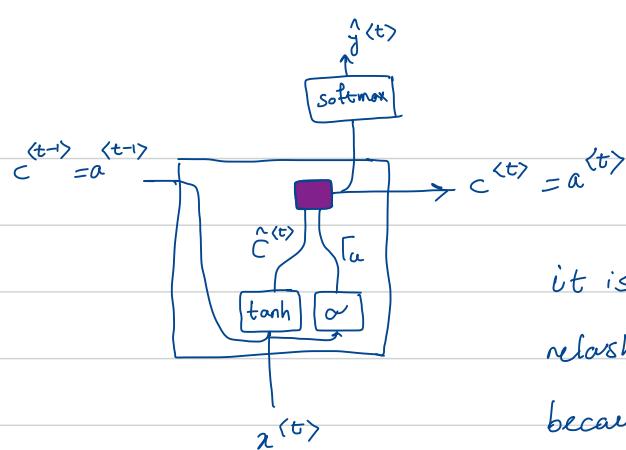
$$c^{<t>} = \Gamma_u * \hat{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

↑ element wise-

$$\begin{matrix} \Gamma_u = 1 & \Gamma_u = 0 & \Gamma_u = 0 \\ c^{<t>} = 1 & \Rightarrow & \end{matrix} \quad \dots \quad \dots$$

↓ since memorize $c^{<t>} = 1$, then select "also"

The cat, which —————, was full.



it is very good at memorizing long range relationships.

because Γ_u is very easy to set to zero
then $c^{(t)} = c^{(t-1)}$,

and since Γ_u is very close to zero, it does not suffer from vanishing gradients.

$c^{(t)}$ is vector R^N

$\tilde{c}^{(t)} \in R^N$

$\Gamma_u \in R^N$, intuition is that every element in Γ_u memorize part of information in the sentence.

Full GRU:

$$\begin{aligned} \tilde{c}^{(t)} &= \tanh(W_c [\Gamma_r * c^{(t-1)}, x^{(t)}] + b_c) && , \Gamma_r \text{ weight one added} \\ \Gamma_u &= \sigma(W_u [c^{(t-1)}, x^{(t)}] + b_u) && \text{to simplified version} \\ r &= \sigma(W_r [c^{(t-1)}, x^{(t)}] + b_r) && \text{to decide which elements} \\ h &= \underline{c^{(t)}} = \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)} && \text{of } c^{(t-1)} \text{ are relevant} \end{aligned}$$

some literature use this notation.

for calculating the $\tilde{c}^{(t)}$

LSTM (long, short term memory)

Paper: Hochreiter & Schmidhuber 1997, long short-term memory

other type of units which allows to learn longer dependencies are LSTM units.
(even more powerful than GRU).

Full GRU:

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

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

$$r_f = \sigma(W_f [c^{(t-1)}, x^{(t)}] + b_f)$$

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

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

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

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

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

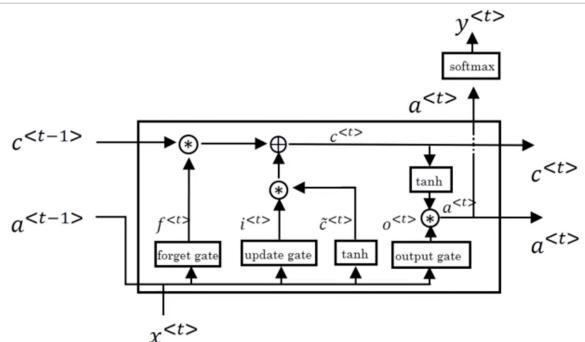
$$\text{output } r_o = \sigma(W_o [a^{(t-1)}, x^{(t)}] + b_o)$$

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

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

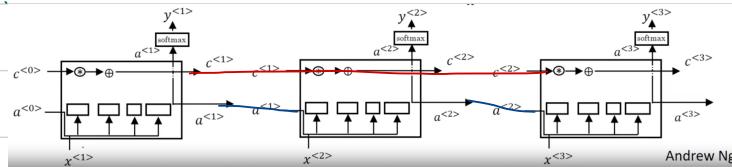
$$\downarrow \\ \text{is no longer directly } = c^{(t)}$$

in LSTM, we have 3 different gates compared to 2 in GRU.



if we connect these units, we build our net

by setting right Γ_u, Γ_f control value, it is easy to keep $c^{(3)} \leftarrow c^{(1)}$ (memorize for) true for gru as well.



Sometime gate value could

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

made to also on previous

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

memory cell values, ($c^{(t-1)}$)

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

this is called peephole connection.

↑ peephole connection.

→ the relationship between $c^{(t-1)}$, Γ would be

one-to-one, (element i^{th} only affect i^{th} element of Γ)

→ LSTM is more strong, much complex. more powerful.

historically more proven choice.

gru is more simpler and easier to build bigger models, train faster.
(maybe just as well).

→ Gru is more recent than LSTM.

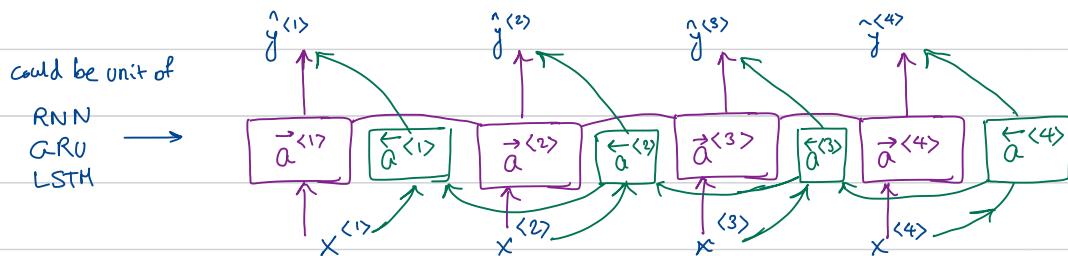
Bidirectional RNN

forward/backward RNNs has a problem that don't take into account the next words

he said, "Teddy Bear one on sale"

"no, "Teddy roosevelt -"

→ This is true whether standard RNN / GRU / LSTM units.

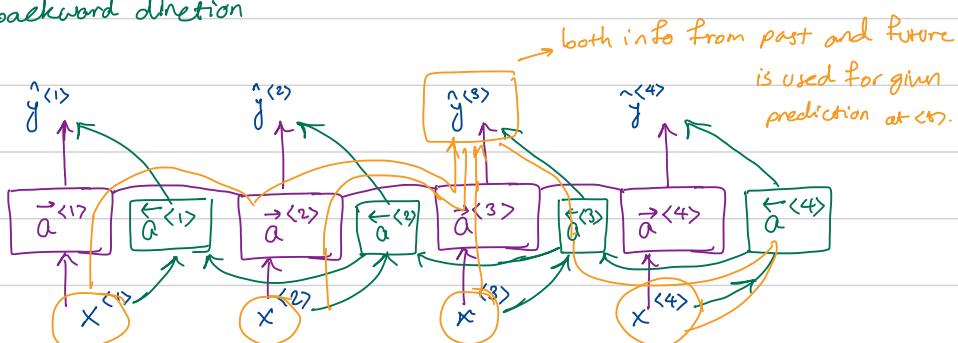


Acyclic graph

first → forward direction

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

second → backward direction

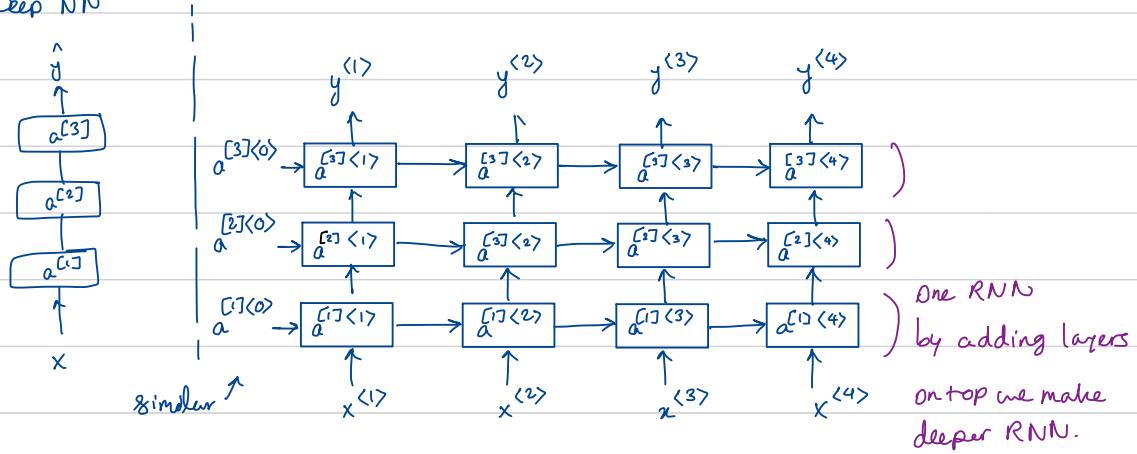


in text (NLP), BRNN w/ LSTM block reasonable choice.

→ disadvantage: entire sequence is needed. for example in speech, we have to wait for full sentence to start predicting. (realtime is not possible).

"Deep RNN"

Deep NN

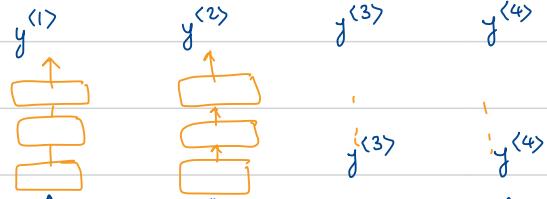


$$a^{[2]<3>} = g(W_a^{[2]} [a^{[2]<2>}, a^{[2]<3>}] + b_a^{[2]})$$

unlike NNs, usually for RNN, having 3 layers is quite a lot and usually deeper versions are not used.

→ instead we may see these versions

which after RNN layers we have

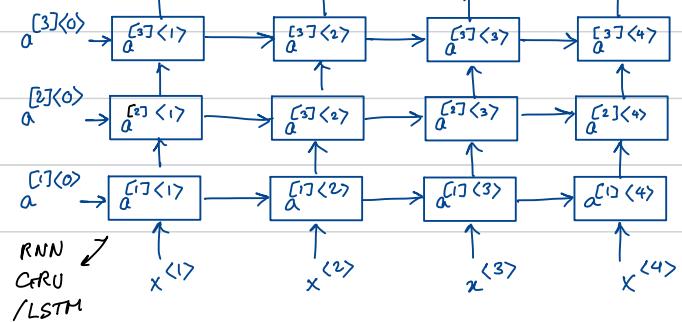


→ RNN layers, which are not connected horizontally.

→ blocks as well could be GRU, or

LSTM blocks,

→ also we could build Deep BRNN as well.



"Back propagation"

NLP and word embedding.

→ "word embedding"

let's model to learn analogies like: Men ~ Women is
King ~ Queen.

weakness of one-hot representation is that treat each word
until itself, and it is difficult to generalize.

Vs[a, aeroplane, --, zulu, <UNK>]

Man (5391)	Women (9853)	King (7157)	...
[]	[0] : 1 0	[]	
↓ 0 5391	↓ 0 9853	↓ 0 7157	

I want a glass of apple juice

orange juice

inner-product of each of these vectors are zero. then really
difficult to generalize that apple/orange, king/queen, --

→ another alternative to one-hot representation is to learn
with feature-representation.

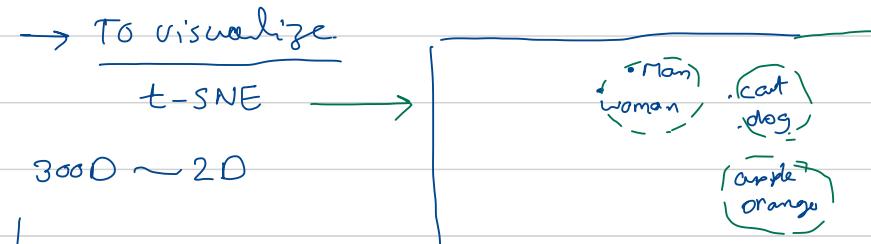
	Man	Women	King	Queen	apple	orange

→ For example 300 dimensional
vector for representing each
word (e^{5391} , e^{9853} , ..., e^{300})

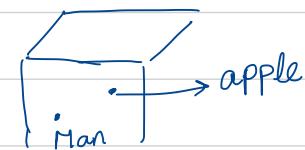
Gender	$\begin{bmatrix} -1 \\ 0.01 \\ 0.03 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 0.02 \\ 0.02 \end{bmatrix}$	-0.95	0.97	0	0
Royal	$\begin{bmatrix} 0.01 \\ 0.03 \end{bmatrix}$	$\begin{bmatrix} 0.02 \\ 0.02 \end{bmatrix}$	-	-	-	-
Age	$\begin{bmatrix} 0.03 \end{bmatrix}$	$\begin{bmatrix} 0.02 \end{bmatrix}$	-	-	-	-
Food	$\begin{bmatrix} 0.02 \end{bmatrix}$	$\begin{bmatrix} 0.02 \end{bmatrix}$	-	-	-	-
:	\vdots	\vdots				
	e^{5391}	e^{9853}				

if we use such embeddings, then representation for apple and orange are close thing, then it would be easier to generalize

In practice we learn these featured representation, and because of this is not easy to neatly interpret each feature.

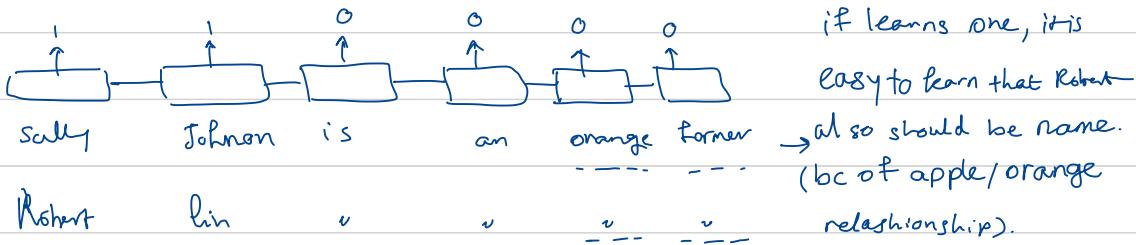


These vectors are called embeddings.



"Using word embeddings."

→ "Named entity recognition example"



a dutch cultivator → how about this

(Normally for such task BRNN should be used) more rare words?

if we have small training set, we may have not seen "durian cultivator", but if we use word embeddings, we should be able to recognize that it is "fruit farmer."

→ idea → using huge data set for word embeddings. $1B \sim 100B$ word
→ later just smaller data set for the task of name entity will be enough

"Transfer learning"

① learn word embeddings from large text corpus.
(or just download)

2. Transfer embeddings to new task with smaller set is available.
(One Pro is that now instead of one-hot vector size of $10^k \sim 300$

3. Optional : fine-tune embeddings with new data.
generally Transfer learning make sense when for the given task,
relatively smaller data is available.

NLP Task which transfer learning could be useful:

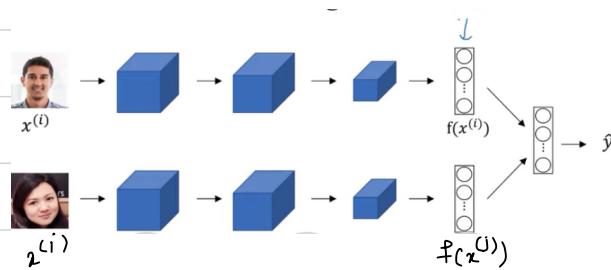
→ Name entity / Text parsing / co-reference recognition.

less useful : language modeling , Machine Translation.

(specially when we have a lot of data dedicated to our Task)

→ Similar concept to face encoding.

→ when we trained a seameese network : face $\xrightarrow{\text{Map}}$ 1280 vector



• encoding/embeddings
one similar concepts.

Difference is that in seameese net we could encode for new faces not seen before.

but in word embedding, we have fixed vocab' set. e_1, \dots, e_{10000}

Properties of word embedding.

linguistic analogies: Man \rightarrow Women as king \sim ?

	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	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

$$\begin{aligned} e_{\text{man}} &= (e_{5391}) \\ e_{\text{woman}} &= (e_{9853}) \end{aligned}$$

Man \rightarrow Women

$$e_{\text{man}} - e_{\text{woman}} \approx$$

$$\begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

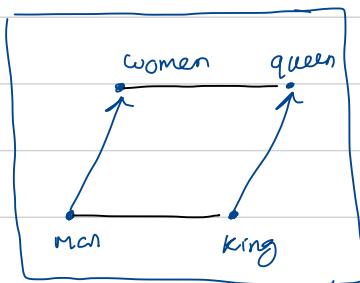
King \rightarrow ?

$$e_{\text{king}} - e_{?} \approx$$

$$\begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

queen

Analogies using word vectors:



300D

$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{?}$

→ find word w :

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

as a performance check for embeddings.

for example 30-75% on analogy reasoning task.

This parallel path is true if visualization done on original 300D space
not in 2D map (ex. by using t-SNE).

Types of similarity functions could be used:

Cosine Similarity:

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

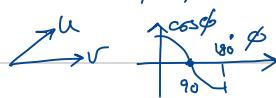
$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

(called cosine, bc it is really
the cos angle between u, v)

→ less popular:

$$\text{euclidian distance: } \|u - v\|^2$$

distance:

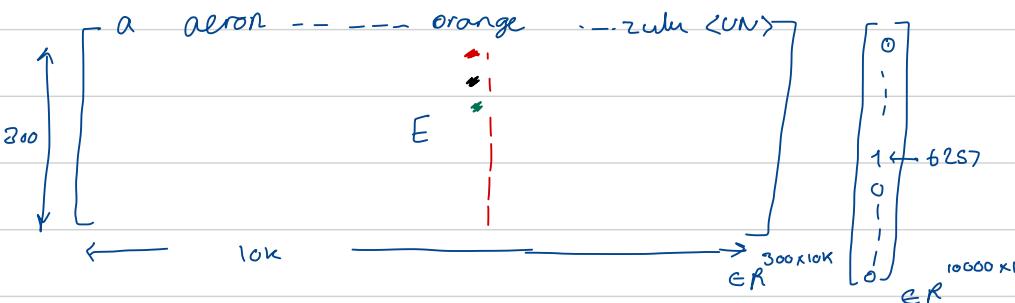


Problems of word embeddings

embedding matrix

Notation for one-hot rep.

O_{6257}



$$E \cdot O_{6257} = \begin{bmatrix} \textcolor{red}{\bullet} \\ \vdots \\ \textcolor{green}{\bullet} \end{bmatrix} = e_{6257}$$

$\rightarrow E \cdot O_j = e_j$ embedding for word j

In practice we really do not use matrix multiplication for this, we use dedicated functions to select the j th column (keras.embedding).

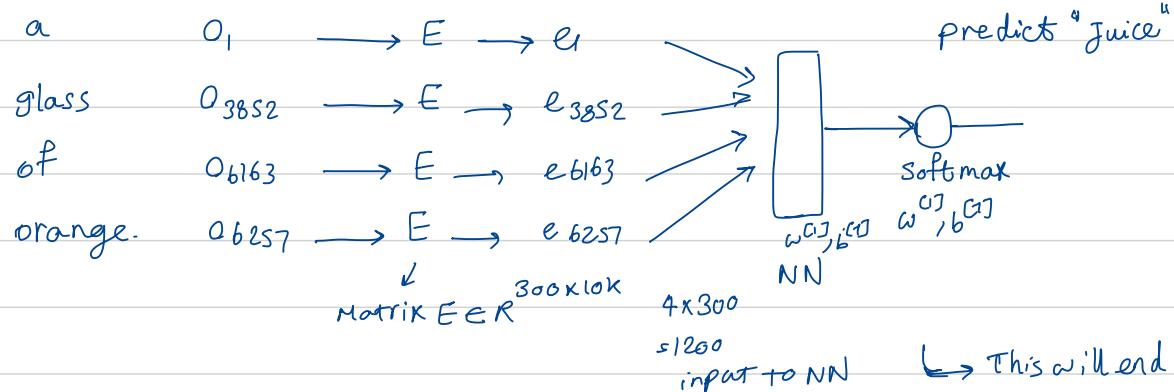
"learning word embeddings"

"Natural language model" → (Paper: Bengio et al 2003, a neural probabilistic language model)

consider ex. I want a glass of orange Juice

4343 9696 1 3852 6163 6257

we build problem like this. that we want to predict the word for example from last 4 word.



other context/target pairs.

→ last 4 words.

I want a glass of orange juice to go along.

{ 4 word on left & right : a glass of orange --- to go along with

{ last 1 word : orange

Nearby 1 word → (skip gram) ex. for orange: glass (skip of)

word2vec

skip-grams.

I want a glass of orange juice to go along with my ---
context → Target (supervised learning).

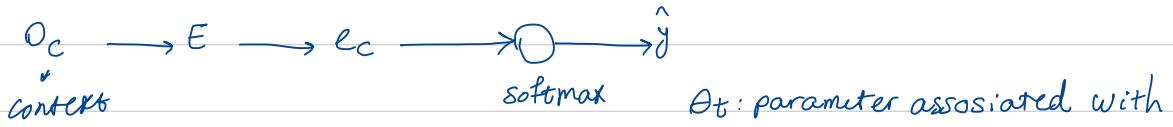
orange	juice	randomly we select the target
orange	glass	word from sentence.
orange.	my	

:

:

→ obviously this would be difficult
supervised learning to really
predict target, but our goal is
just to learn word embeddings.

(No NN between, only $E \rightarrow \text{softmax}$).



$$\text{softmax} = p(t/C) = \frac{e^{\theta_t^T \cdot e_c}}{\sum_{j=1}^{10k} e^{\theta_j^T \cdot e_c}}$$

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

$$y = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \text{one-hot rep.}$$

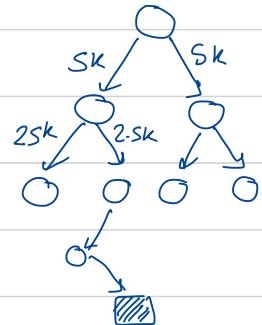
Problems with soft max classification

$$\rightarrow P(t|c) = \frac{e^{\theta_t^T \cdot \theta_c}}{\sum e^{\theta_j^T \cdot \theta_c}}$$

① every time when calculating, we need to sum over all vocab in list (ex. 10k here).

→ we could use "hierarchical softmax".

→ computation complexity would be $\log(v)$
for soft max $\approx |V|$



→ In practice hierarchical softmax does not use perfectly balanced tree.
(More common words be in top, --)
(heuristical approach to build tree--)

how to sample context c ?

if done randomly, you see a lot of the, of, a, --- more frequently.

• use heuristically developed algorithms to sample in balance to cover all words

* Another variation CBOW : Take context from surrounding middle word
and try to predict middle word.

Negative Sampling.

→ Algorithms like W2V has the problem of complex computation on softmax layer (due to large size of the vocabulary set).

Negative sampling is a method, tries to solve a much reduced version of supervised learning problem.

→ Imagine a sentence: I want a glass of orange juice to go along . . .

Training Data are organized like this.

context	word	target	
orange	juice	1	one positive example and
orange	king	0	"K" example for negative which
" "	book	0	randomly found.
" "	of	0	→ even though maybe the randomly
c	x		selected word "of" may really happen
t			on sentence. (ok.)

↳ "supervised learning problem"

Ks 5-20 small datasets
 Ks 2-5 larger datasets

Model:

$$\text{softmax } p(t|c) \leftarrow \frac{e^{\theta_t^T \cdot \theta_c}}{\sum_{i=0}^{l_{\text{tot}}} e^{\theta_i^T \cdot \theta_c}}$$

called negative sampling, since we have one positive example and then randomly sample negative examples.

to model $P(y=1 | c, t)$ we use logistic regression.

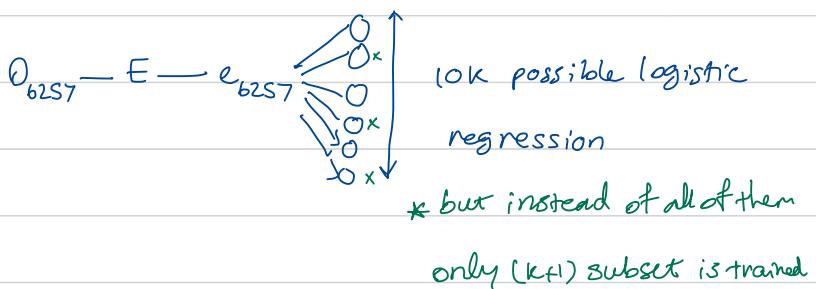
$$P(y=1 | c, t) = \sigma(\theta_t^T \cdot \theta_c)$$

if we have k negative example, the ratio is $k/1$ neg/pos example.

to train this model.

How to train:

example orange
6257



→ training only $(k+1)$ binary logistic regression problem is much cheaper than softmax with 10k units.

→ "selecting negative examples"

One way: candidate for negative could be sample based on empirical frequency of word in dic : a lot of the, of, ...

Other way: uniformly sample $\frac{1}{|V|}$: not representative of real significance of words.

empirically suggested formula $P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{i=1}^{10k} f(w_i)^{3/4}}$

$f(w_i)$: observed

frequency of w_i in lang.

(kind of trade off b/w two extremes of 1&2)

Glove algorithm (global vectors for word representation)

→ simpler but not as popular as W2V or skipgrams.

Paper: Pennington et al, 2014. GloVe: Global Vectors for word representation).

c, t

let x_{ij} : # time i appears in context of j

\downarrow

t

\downarrow

c

→ depends on definition of c, t maybe $x_{ij} \neq x_{ji}$

exp: c, t one the word appear in the proximity of n word x_{ij}, x_{ji}

: c is the word immediately before t : $x_{ij} \neq x_{ji}$

Model:

$$\text{minimize } \sum_{i=1}^{10k} \sum_{j=1}^{10k} f(x_{ij}) (\theta_i^T \theta_j + b_i + b_j - \log x_{ij})$$

\uparrow
similar to " $\theta_t^T \cdot e_c$ "
as before

Solved with GD to minimize:

if $x_{ij} = 0$ then $\log x_{ij} = \text{nan}$, then $f(x_{ij}) = 0$ if $x_{ij} > 0$

another point for f.e.) is to correct the weight for

1) very frequent words This, is, of, a, --

2) less frequent but rather important duration.

→ another point here is that θ_i, e_j are symmetric.

→ (both should initialized randomly → GD → then we can take average

$$e_w^{(\text{final})} = \frac{e_w + \theta_w}{2}$$

GloVe algorithm came later than skip-gram, W2V, ... and simplified a lot the learning embeddings.

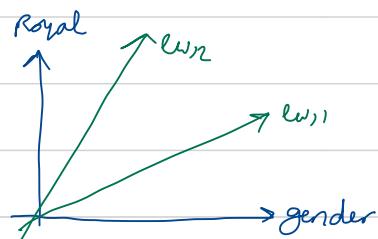
Important Points:

for word embedding learned by all these algorithm, there is no guarantee to be "human interpretable"

	Man	Women	King	Queen
gender	:	-	-	:
Royal	:	-	-	:
Age				

$$\theta_{i \mid j}^T$$

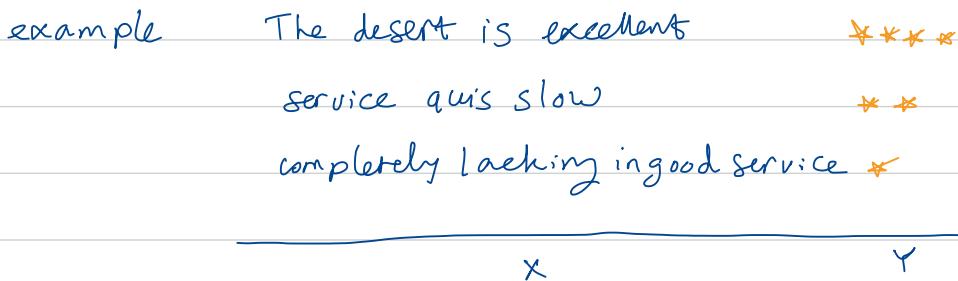
$$\hookrightarrow (A Q_i)^T (A^T \theta_j) = \theta_i^T \theta_j$$



→ best parallelogram for

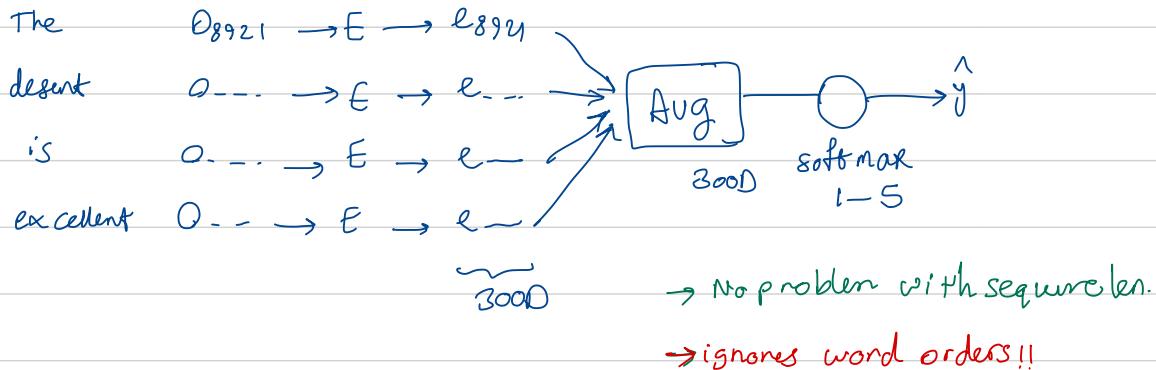
figuring out linguistic analogies holds.

"Sentiment classification"



challenge: Maybe not of a big training set.

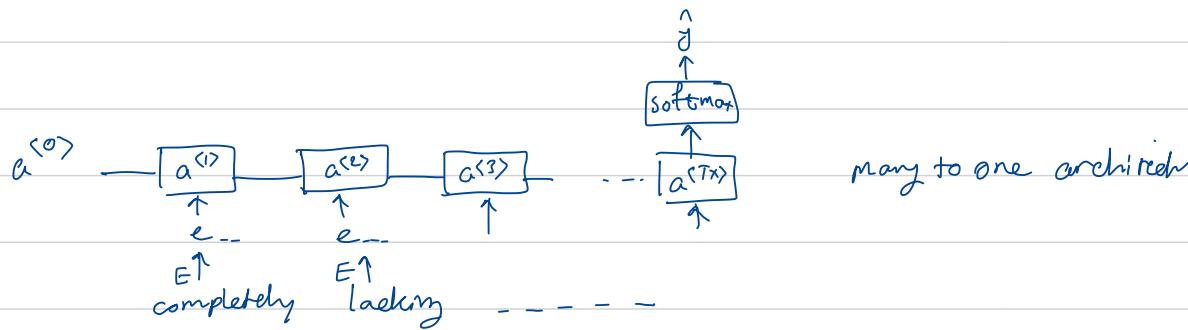
one simple solution, imagine we have word embeddings E



ex. completely lacking in good taste, good..., good...

↳ very likely miss this example due to not considering the order.

"RNN for sentiment classification"



This architecture could learn the orders.

advantages of using pretrained embedding, is that even some word were not in our labed traing set, still would be recognize those and generalize.

Debiasing word embeddings.

example:

bias in terms of gender, age, sexual

social status, orientation, ethnicity, ... which are in texts.

Man: Computer-programmer as Women: Homemaker ~~X~~

Father: Doctor as Mother: Nurse. ~~*~~

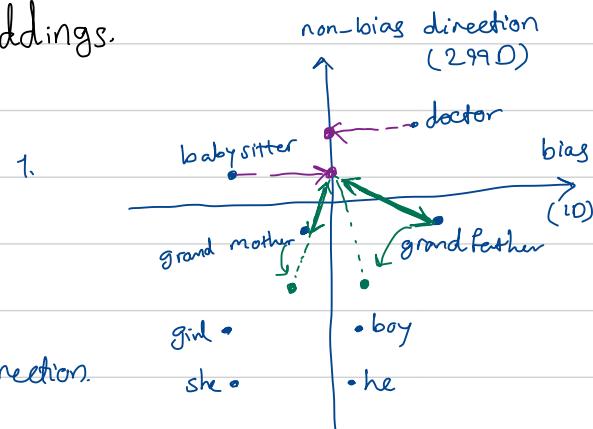
we want to eliminate these type of biases which is learned from texts used for training.

addressing bias in word embeddings.

1. Identify bias direction

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

average. → find the bias direction.



2. Neutralize: for every word that is not definitional, project to get rid of bias. (doctor, babysitter, ...)

3. Equalize pairs: (like the distant of grand Mom/father to baby sitter, doctor transpose to new position to be equally distanced.)

What words has to be neutralized or not?

- a linear classifier could be trained to find which has to be neutralized. (small set like grand Mom/father are definitional).
- same for equalization pairs as well.

Sequence to sequence models.

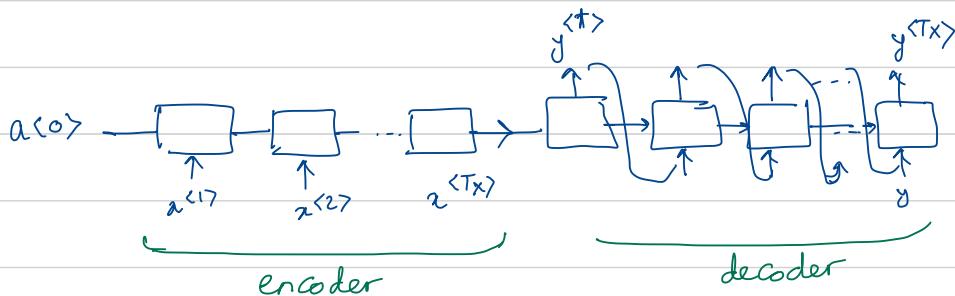
application: machine translation, speech recognition, ..

$x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$ $x^{(5)}$

Jane visite l'Afrique en septembre

↓
Jane is visiting Africa in september

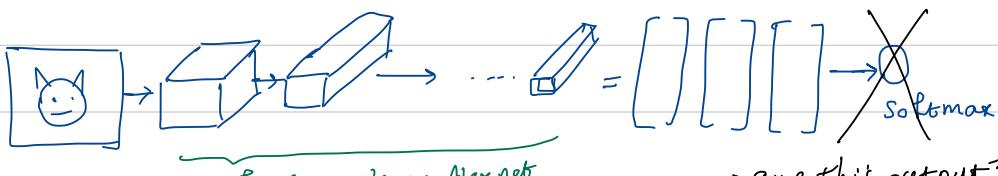
$y^{(1)}$ $y^{(2)}$ $y^{(3)}$ $y^{(4)}$ $y^{(5)}$ $y^{(6)}$



→ This architecture actually works!

(given enough pairs en/Fr for training)

Similar architecture also used for image captioning.



→ give this output to

decoder to generate caption.

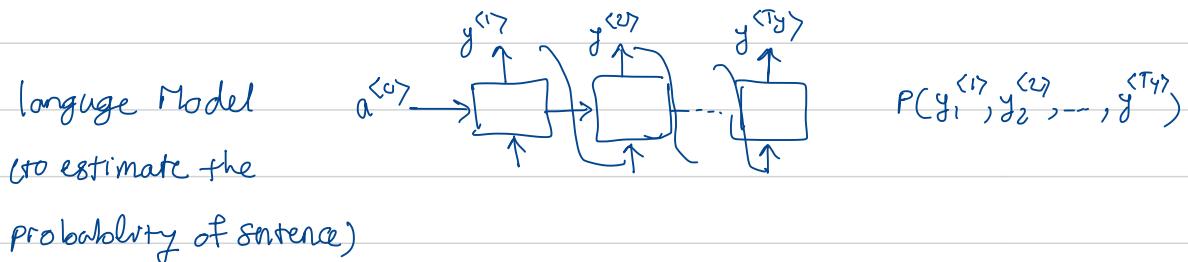
May groups:.. Moos 2014

| Vinyals 2014

| Korpathy and FeiFei 2015

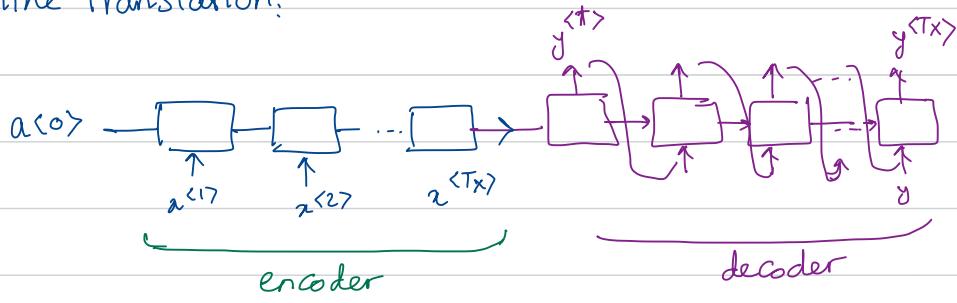
This as well works pretty decent!

"Picking the most likely sentence"



→ such architecture as well, we used to generate Novel sentences.

Machine Translation:



!!

decoder is very similar to language

→ model, just the input is not \emptyset , and actually is generated from other sequences (encoder). because we call this:

"Conditional language model"

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

exp. $P(\text{english sent.} | \text{French sentence})$.

but, what we don't want is to sample randomly (similar to what we did while generating Novel sentences).

Jane visite l'Afrique en septembre

- different samples, not all good!
- ↓ Jane is visiting Africa in september ①
 - | Jane is going to be visiting --- ②
 - | In september, Jane will --- ③

we want an english sentence which maximize:

French sentence.

$$\arg \max_{y^{(1)} y^{(2)} \dots y^{(T_y)}} P(y^{(1)}, \dots, y^{(T_y)} | x)$$

→ why not just use a greedy search?

greedy search: first $P(\hat{y}^{(1)} | x)$ then $P(\hat{y}^{(2)} | x, \hat{y}^{(1)})$, ...

best first → best 2nd, ...

this is not good, indeed we are interested

in best joint probability $P(\hat{y}^{(1)}, \hat{y}^{(2)}, \dots | x)$

for example $P(\text{sentence-1} | x) > P(\text{sentence-2} | x)$

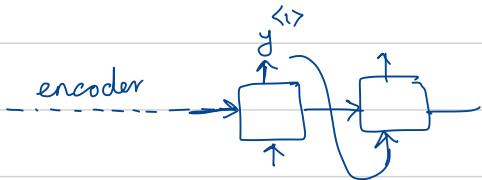
→ if we use greedy $P(\text{Jane is going} | x) > P(\text{Jane is visiting} | x)$

is more frequent in english

but not correct choice, given whole sentence.

also given 10k, vocab list, searching for all combination is also not good!

→ Beam search Algorithm



① $P(y^{<1>} | x) \rightarrow$ Find n most probable
1. in 2. Jane 3. september.

n : Beam width
 $= 3$ (here).

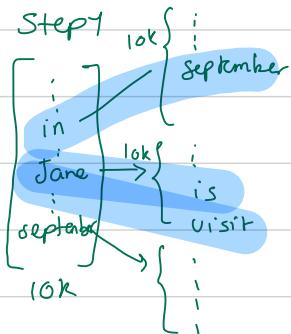
② a. Take First best find $P(\hat{y}^{<2>} | x, \hat{y}^{<1>})$

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

we have this from ① 10k choice

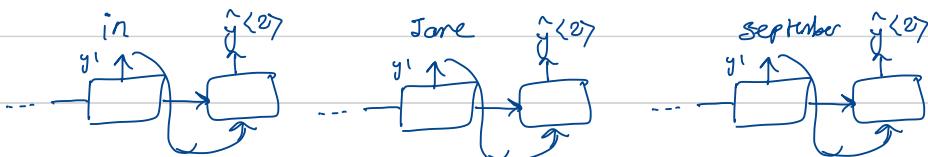
b. Same for second and 3 word as well 10k + 10k choice.

Again Pick 3 (Beam width) most probable.



for example in second step we rejected "september" as the first word.

→ we only need 3 initiation of network to evaluate 30k output.



Same routine continues:...

$y^{(1)}$ $y^{(2)}$
in September

tok
Jane

$P(y^{(3)} | x, y^{(1)}, y^{(2)})$

Jane is

tok
visiting

Jane visits

tok
africa

This routine of search continues, until it finds $\langle \text{Eos} \rangle$ and most likely sentence.

Beam-width = 1 → becomes greedy search.

→ Refinement to beam search

length normalization

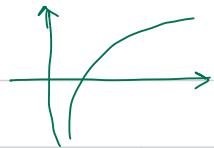
$\arg \max_y \prod_{t=1}^T P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$ → these probabilities are

instead we use \log small number and product

$\arg \max_y \sum \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$ will be very small for float

→ More numerically stable.

log is strictly monotonically increasing function, should give same results.



→ One of the problem with these objective function is that it always prefers smallest sentence. (as $Ty \uparrow P \downarrow$)

then we add Normalization

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

the term helps a lot to decrease the penalty for longer sentences. or heuristically

$$\frac{1}{Ty^\alpha}$$

$\alpha = 0.7$ (good).

$\alpha = 1$ full normalization

$\alpha = 0$ No at all

Summary :

we run the algorithm of Beam-search and we find a lot of possible sentence. (till $Ty \leq 30$)

Then we use the "Normalized log-likelihood" to pick the best.

→ Beam-width

| larger : More probable to find best / More computation

| small : worse results / faster.

10 ~ 100 (For products, 100 is already huge), research 1K~3K

to get the best results. (improvement of course diminish as goes very high)

→ Beam search Unlike (BFS) or (DFS) is not guaranteed to find the best but much faster. (Breadth/ Depth First search).

Error Analysis in Beam search

Beam search is heuristic search algorithm and does not guarantee to find the best possible. suppose:

→ Input: Jane visite l'Afrique en septembre

human: Jane visits Africa in September (y^*)

Algo : Jane visited Africa last september \hat{y} (Pretty bad translation).

who is to blame?

since here we had --- [enc.] --> [dec.] ↑ Beam-search.

two part 1. RNN model

2. Not globally optimal Search

error Analysis: there are two possibility:

case 1: $P(y^*|x) > P(\hat{y}|x)$

Means, our RNN correctly knows that y^* is more probable (B.S. is fault)

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

Beam search was fine, RNN is fault. → Increase Beam-width.

(small complication, may be due to bad normalization at B.S could be responsible as well. → all relevant technique for net (More data or architecture...))

Normally analysis done for a number of examples to find the most impactful of errors.

human	Algo	$P(y^* x)$	$P(\hat{y} x)$	fault
--	---	2×10^{-10}	7×10^{-10}	B
--	--	:	:	R
--	--	:	:	Y

"Bleu score" (bilingual evaluation understudy)

How about if there are multiple good translations?

*understudy is theatre is the person that learn to take the role of senior player if necessary)

Paper: Papineni 2002, Bleu ...

Imagine example
that we have 2
reference human
translation

ref1: The cat is on the mat

ref2: there is a cat on the mat

MT: the the the the the the (pretty awful).

Simple Precision: $\frac{7 \text{ (appear on ref)}}{7 \text{ (word on MT)}}$ $\leq 1 \leftarrow \text{bad!!}$
on MT words (the)

Modified P = $\frac{2 \text{ (count_clip)}}{7} \rightarrow$ Maximum number of appearance
over all reference sentence.

→ 2nd exp: MT: the cat the cat on the mat

bigrams (P on every two-pairs).

the cat	2	1
cat the	1	0
cat on	1	1
:	1	1
count	count	clip

$$P_1 = \frac{\sum_{\text{unigram}}^{} \text{count}_{\text{clip}}(\text{unigram})}{\sum_{\text{unigram}}^{} \text{Count}(\text{unigram})}$$
$$P_n = \frac{\sum_{n\text{-gram}}^{} \text{count}_{\text{clip}}(n\text{-gram})}{\sum_{n\text{-gram}}^{} \text{Count}(n\text{-gram})}$$

if one exact ref is found $P_1, \dots, P_n = 1.0$

More detail:

P_n : Bleu score on n-grams only

$$\text{combined bleu score} \quad \underbrace{B_p \cdot \exp\left(\frac{1}{n} \sum_{n=1}^K P_n\right)}$$

B_p (brevity penalty): to penalize generating short sentences.
(short sentence have higher chance of getting higher score).

B_p (heuristically defined)

$$B_p = \begin{cases} 1 & \text{if } MT_out_len > ref_out_len \\ \exp(1 - MT_out_len / ref_out_len) & \text{otherwise} \end{cases}$$

"Attention Models"

Paper: Bahdanau, et al 2014. Neural machine translation.

one of the most influential idea's in deep learning.

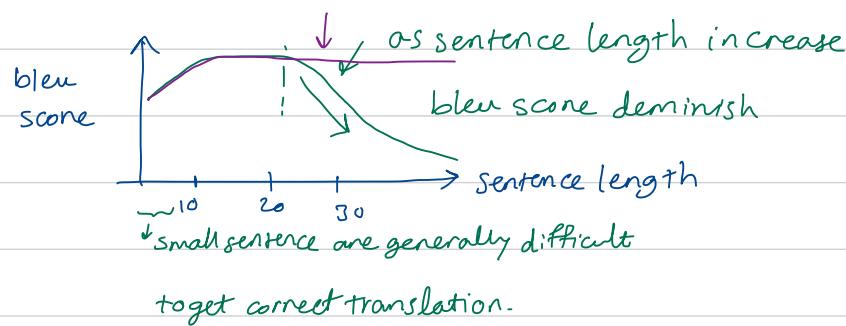
Problem with long sentences.

exp → fr: Jane s'est rendue en Afrique -----, ---, --

en: Jane went to Africa fast -----

since for MT, we use encoder/decoder, having a long sentence, it is very difficult for encoder to memorize whole sentence and then generate the sentence. human as well, does the translation part by part.

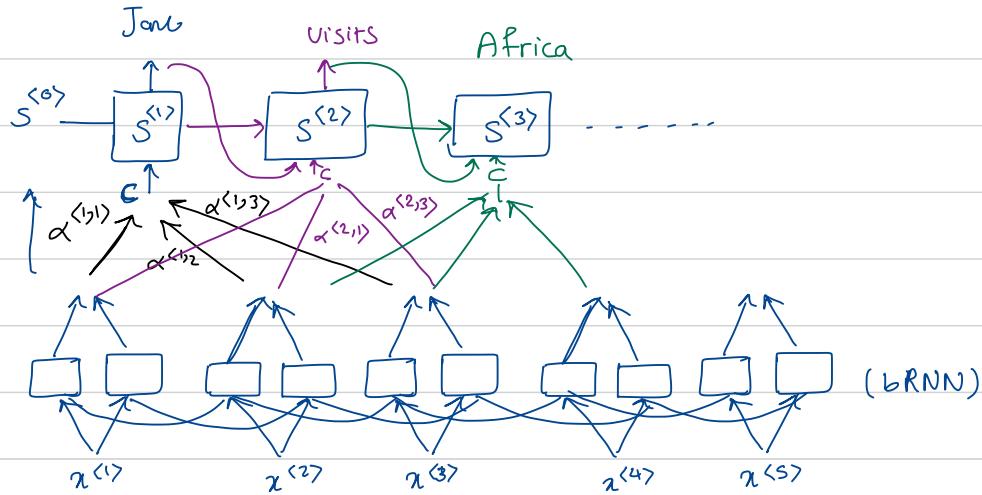
with attention model.



with attention model, Network does the translation more like human in small parts.

originally developed for MT but later extended for many others.

"Attention model intuition"



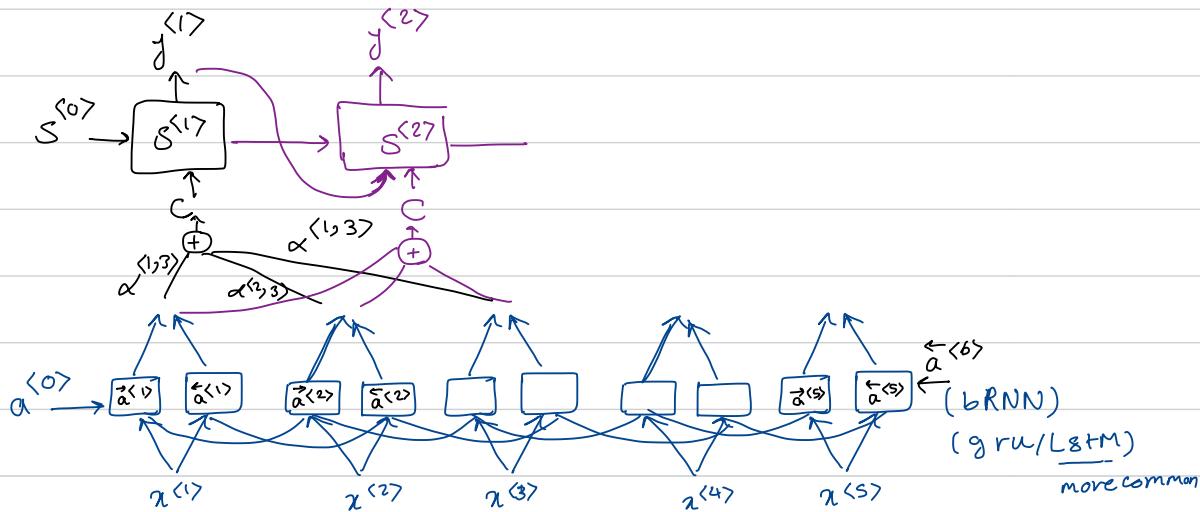
Jane Visite l'Afrique en septembre

the idea is that while generating the translation we also gave a context from similar position of french sentence

$$\alpha^{(t, t')} \rightarrow \begin{cases} S^{(t-1)} \\ \overset{\rightarrow}{\alpha}{}^{(t)}, \overset{\leftarrow}{\alpha}{}^{(t)} \end{cases} \quad t: \text{english}, t': \text{French}$$

this help MT to also look at the local position of original sentence.

Attention model : (detail)



Jane visite l'Afrique en septembre

t' : used for French sentence

$$\alpha^{(t')} = (\tilde{\alpha}^{(t')}, \hat{\alpha}^{(t')})$$

→ context calculated by a sum: $\sum_{t'} \alpha^{(1,t')} = 1 *$

$$c^{(1)} = \sum_{t'} \alpha^{(1,t')} \tilde{\alpha}^{(t')}$$

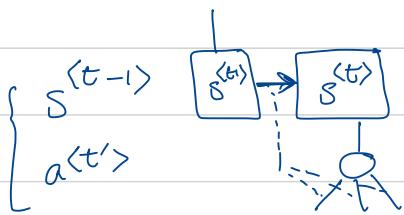
* Indeed $\alpha^{(t,t')}$ is amount of attention $y^{(t)}$ should pay to $\alpha^{(t')}$

$$\alpha^{(t,t')} = \frac{\exp(e^{(t,t')})}{\sum_{t' \neq t} \exp(e^{(t,t')})}$$

(like a softmax)
→ ensure to sum to 1

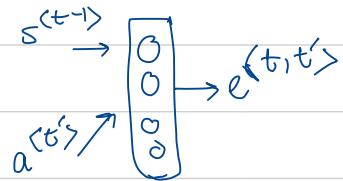
How to compute $e^{(t,t')}$?

$e^{(t,t')}$ should be calculated from



we could train a very simple model

to learn this function and trust it



due to using attention, this algorithm

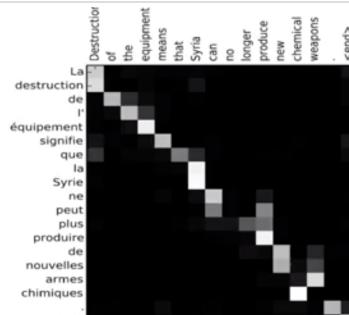
runs at quadratic time/cost but for MT where T_x, T_y are not very big, it could be acceptable. ($T_x \cdot T_y$ times).



same idea of attention used in image captioning and other area's as well.

Paper: Xu et.al. 2015. attend and tell:...

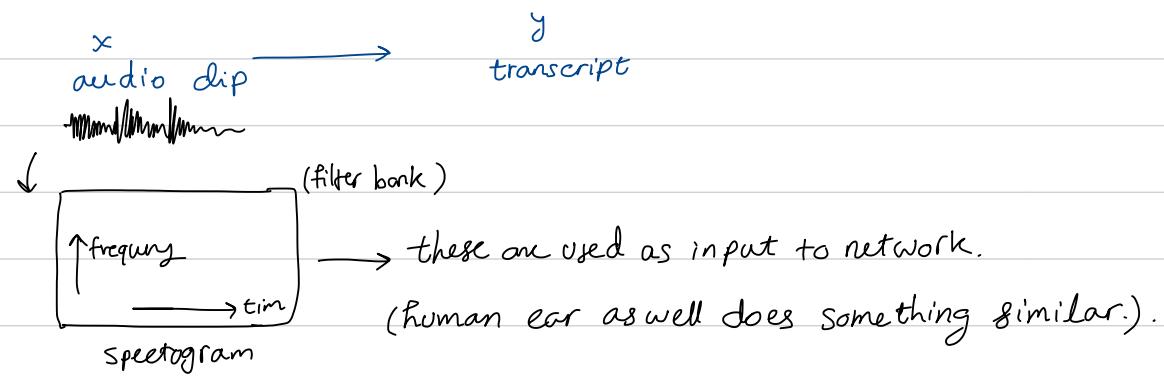
Visualization of $\alpha^{(t,t')}$:



→ plotting the attention

→ usually pays attention to correct word.

"Speech recognition"



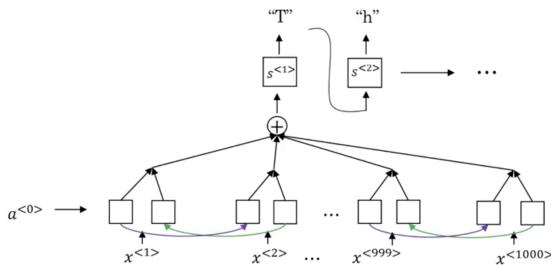
before phonemes concepts used: like.

The quick brown → d e kuik brawn (hand engineered features).

{ 300h / 3000h used in academia

Commercial product 100K h of data for training.

one solution could be to use Attention models for speech recognition

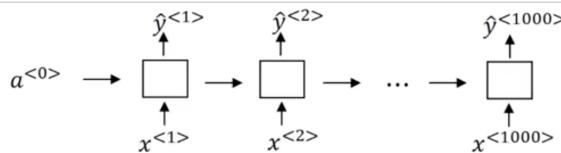


Second option

CTC Cost for speech recognition

(connectionist temporal classification) Graves et. al 2006.

since imagine $x^{<1>} - x^{<1000>}$ are from 'to' second speech, then y is really don't have that many characters.



→ in reality bidirectional
gru/LSTM is used
and deeper models.

output may be like:

ttt - h - eee - - - ^{space} ↕ - - - 999 - - -
↑
CTC collapse → the quick - - - ↘ blank

→ this allows network to have 1000 output by repeating letters.

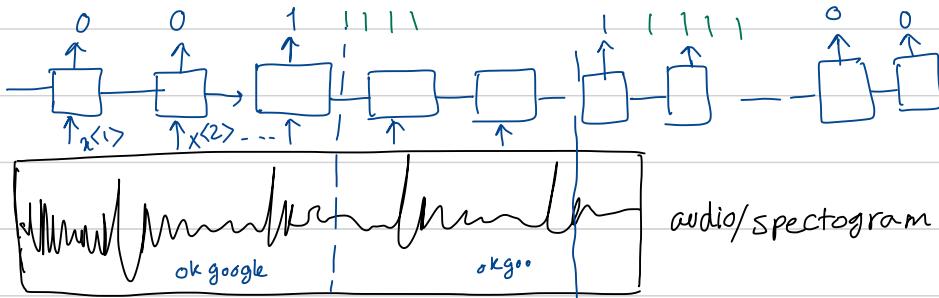
both options (using attention models and CTC models
both works)

Trigger word detection

(or keyword detection).

examples: devices which wake by a keyword; hey siri, ok google
(still evolving area).

one way could be ...



everytime in training data trigger word finishes, output will become 1.

To improve on unbalanced output label (lots of 0s and rare 1s), we could output for few more timestamps.