

Why sequence models

Example of sequence data.

- Speech recognition. → Machine translation.
- Music generation. → Video activity recog.
- Sentiment classification. → Name entity recog.
- DNA seq. analysis.

Notation:

e.g. Name Entity Recog.

- x : Harry Potter and Hermione Granger invented a new spell.
 $x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \dots \quad x^{<t>} \quad \dots \quad x^{<9>}$
- y :
 $y^{<1>} \quad 1 \quad 1 \quad 0 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0$
 $y^{<2>} \quad y^{<3>} \quad \dots \quad y^{<t>} \quad \dots \quad y^{<9>}$

$T_x = T_y = 9$. length of x and y .

$x^{(i)<t>}$: for the t^{th} index of the i^{th} training example

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

Same for $y^{(i)<t>}$ and $T_y^{(i)}$

Representing Words.

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

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \dots \quad x^{<t>} \quad \dots \quad x^{<9>}$

e.g. Vocabulary: We have 10,000 words in it.

a
aaron
:

One way to represent words: One-hot encoder:

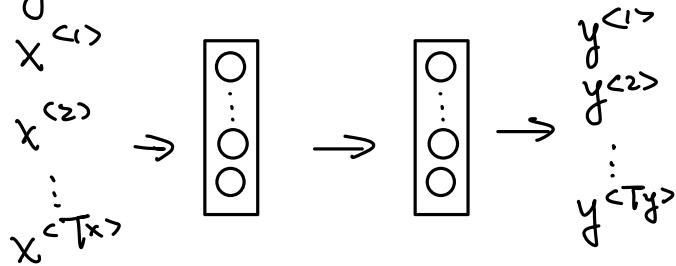
TOP ↑

and	367	$x^{<1>}:$	$\begin{bmatrix} \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\leftarrow 4075$	$10,000$	for all the words.
harry	7075					
Potter	:					
Zulu.	10,000					

For words never appear in the vocab. may use a new token $\langle \text{UNK} \rangle$

Recurrent NN model.

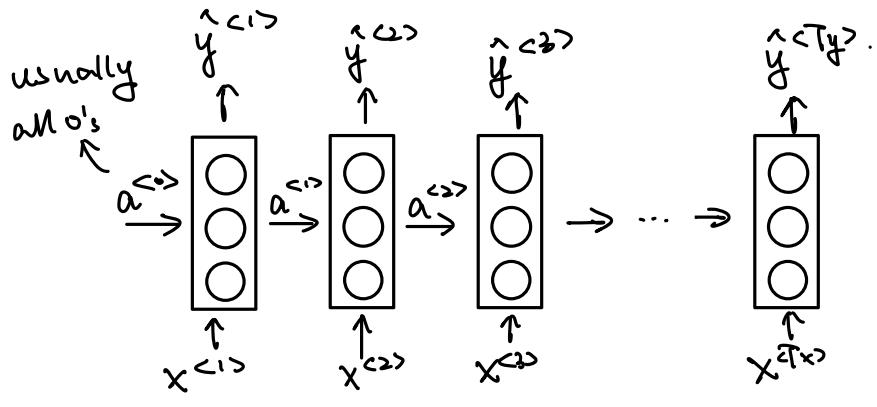
Why not a standard NN.



Problem :

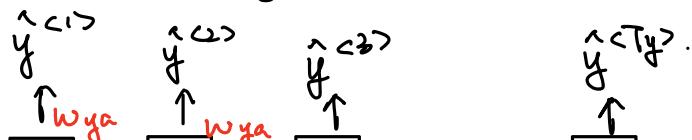
- Inputs outputs can be different length in diff. examples.
- Doesn't share features learned across diff pos. of text.

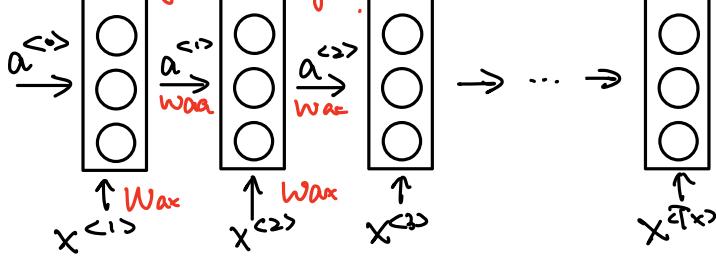
Recurrent NN.



e.g. for the prediction of $y^{<3>}$
it does use information from
 $y^{<1>}$ and $y^{<2>}$.
(But it doesn't have info later)

Forward Propagation.





Computation Step :

$$a^{<0>} = \vec{0}, \quad a^{<1>} = g(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a) \leftarrow \text{Tanh, ReLU} \\ \hat{y}^{<1>} = g(W_y a^{<1>} + b_y). \quad \leftarrow \text{Sigmoid.}$$

...

$$a^{<t>} = g(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a) \\ \hat{y}^{<t>} = g(W_y a^{<t>} + b_y).$$

Choice of activation g :

Simplified Notation for RNN.

$$a^{<t>} = g(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a) \Rightarrow a^{<t>} = g(W_a [a^{<t-1>} \ x^{<t>}]) + b_a \\ \hat{y}^{<t>} = g(W_y a^{<t>} + b_y).$$

Suppose W_{aa} is $(100, 100)$ matrix. W_{ax} is $(100, 10, 100)$ matrix.

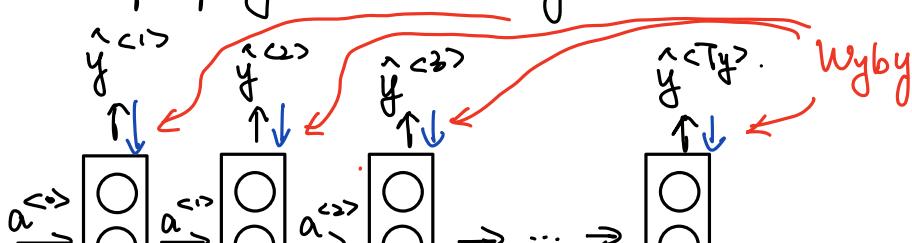
$$W_a = \begin{bmatrix} W_{aa} & | & W_{ax} \end{bmatrix} \stackrel{\uparrow}{\substack{100 \\ \downarrow \\ 100 \rightarrow 10000 \rightarrow }} \quad W_a \text{ is } (100, 10100) \text{ matrix.}$$

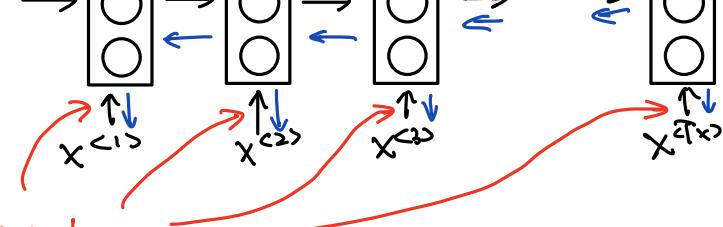
$$[a^{<t-1>} \ x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ \hline x^{<t>} \end{bmatrix} \stackrel{\uparrow}{\substack{100 \\ \downarrow \\ 10000}} \quad \stackrel{\uparrow}{\substack{10 \\ \downarrow \\ 100}} \quad \text{So } W_a [a^{<t-1>} \ x^{<t>}] \\ = W_{aa} a^{<t-1>} + W_{ax} x^{<t>}$$

Rather than use two weight matrices W_{aa} and W_{ax} . we only use W_a .

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

Buckpropagation Through time.





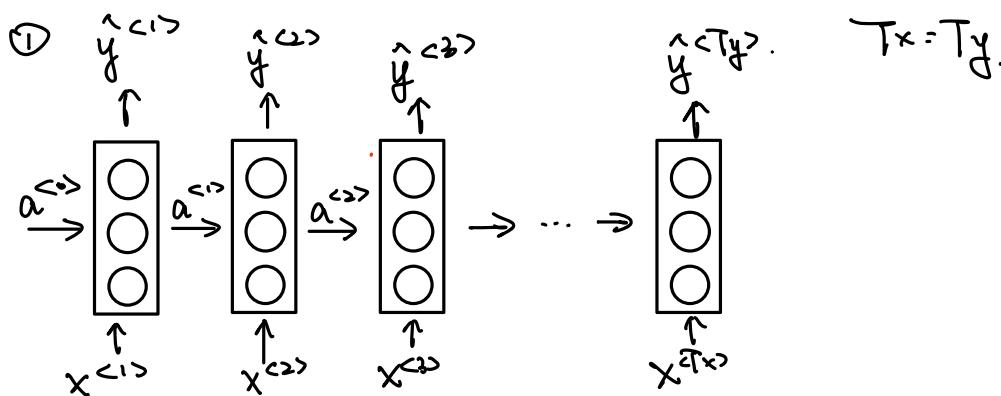
Waba

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

$\mathcal{L}(\hat{y}^{<t>} | y^{<t>}) = \sum_{t=1}^T \mathcal{L}^{<t>}(\hat{y}^{<t>} | y^{<t>})$. Backprop through time.

Different Types of RNNs. (Previously $T_x = T_y$, but $T_x \neq T_y$ in general).

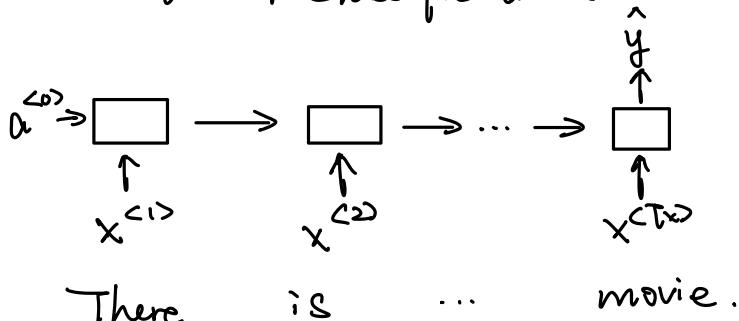
e.g. of RNN Architecture.



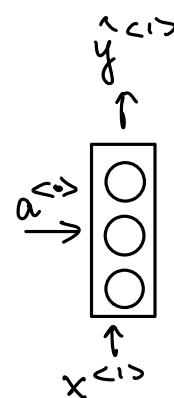
Many-to-many (Equal).

② Sentiment Classification.

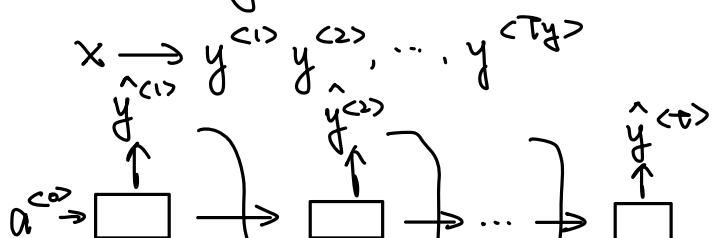
③ Of course you may have one-to-one.

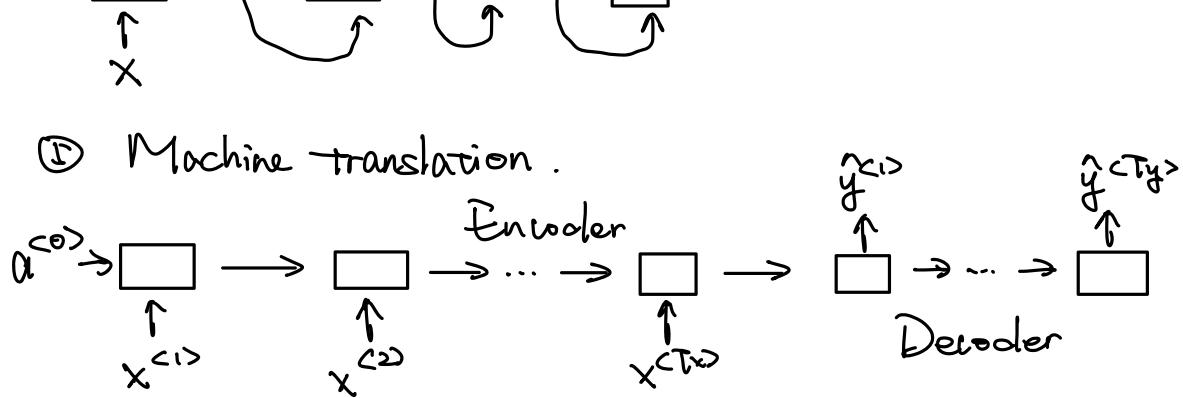


Many-to-one.



④ Music generation.





Many-to-many, (where T_x may not equal T_y)

Language Model and Sequence Generation.

What is language modeling?

Speech recog.

- ① The apple and pear salad.
 - ② The apple and pear salad.

Maybe the second one is much more likely.

$$P(\text{Sentence } \textcircled{1}) = 3.2 \times 10^{-13} \quad P(\text{Sentence } \textcircled{2}) = 15.7 \times 10^{-10} \quad \checkmark$$

Basically, the language model. Take in a sequence $(y^{<1>} , y^{<2>} , \dots , y^{<Ty>})$ and compute: $P(y^{<1>} , y^{<2>} , \dots , y^{<Ty>})$

Language modeling with an RNN.

Training set: large corpus of English text. Vocab size 10000.

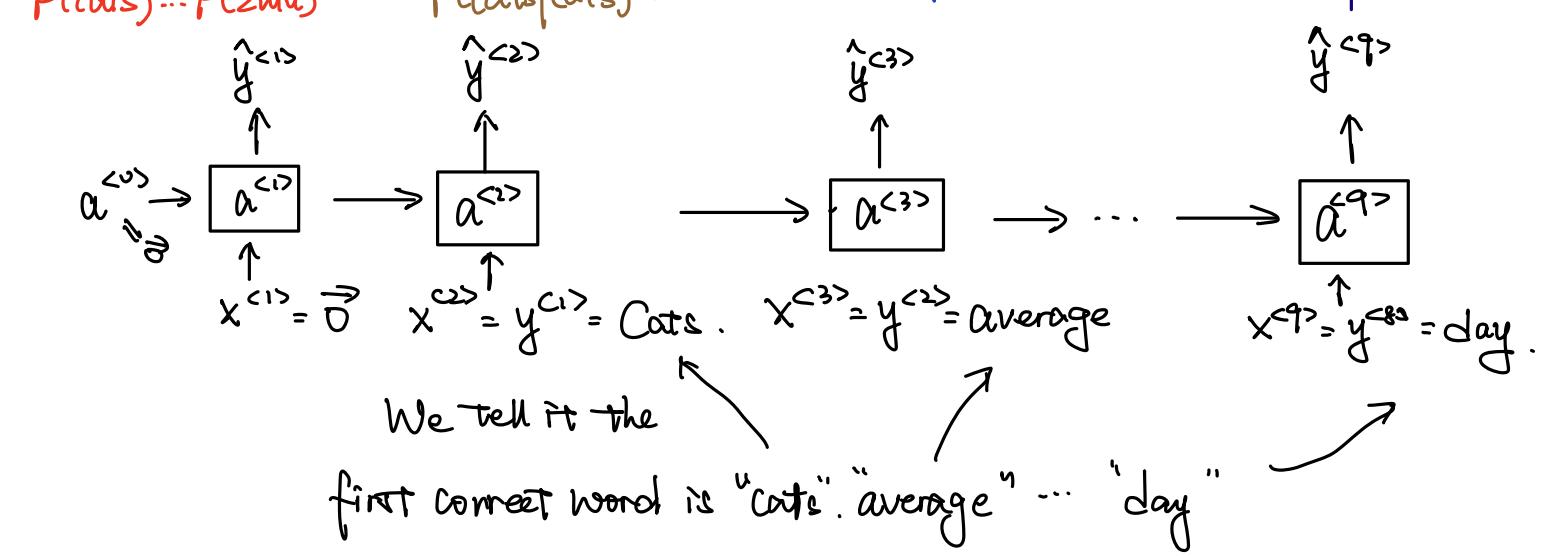
e.g. Cats average 15 hours of sleep a day. <EOS>.

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

e.g. The Egyptian Man is a bread of cat. <EOS>.
<LINK>. if not in the vocab.

has a softmax to compute:

$$P(a), P(a|o), \dots, P(a|t) \quad P(a|cats), P(a|o|cats) \quad P(_ | cats, average) \quad P(_ | \dots)$$



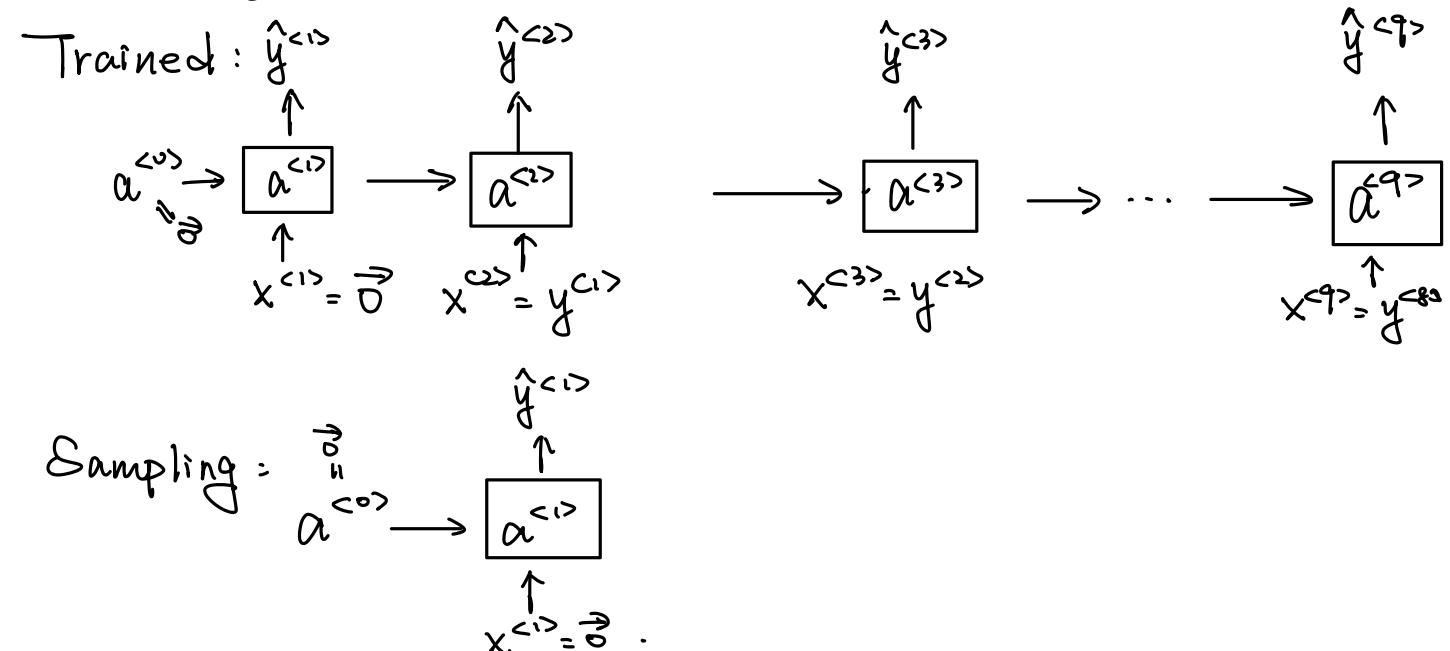
$$\mathcal{L}(\hat{y}^{<t>} \cdot y^{<t>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

$$\mathcal{L} = \sum_t \mathcal{L}(\hat{y}^{<t>} \cdot y^{<t>})$$

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

Sampling Novel Seq.

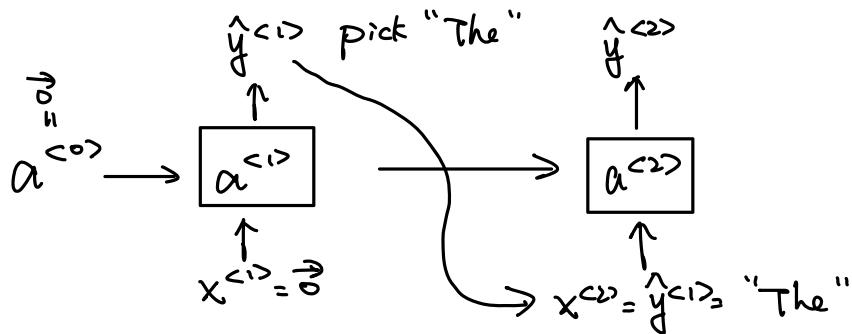
Sampling a seq. from a trained RNN.



$\hat{y}^{<1>}$ will give you all the probability of all the words in the vocab.

Then we just use np.random.choice to pick one.

Step ② . After we picked one . like we pick "The".



And $\hat{y}^{(2)}$ will give you all the probabilities of words given "The". We keep continue this process until:

- ① Sample certain amount of words.
 - ② We sampled the <EOS> token.

Character level language model.

$\text{Vocab} = \{\text{a}, \text{b}, \text{c}, \dots, \text{x}, \text{y}, \text{z}, _, \dots, _, \emptyset, \dots, \emptyset, \text{A}, \text{B}, \dots, \emptyset\}$

Then your $\hat{y}^{(1)}$. $\hat{y}^{(2)}$. will be characters.

e.g. $C \underset{y^{(1)}}{a} \underset{y^{(2)}}{t} \sqcup a v e r a g e$
 $y^{(3)} \dots$

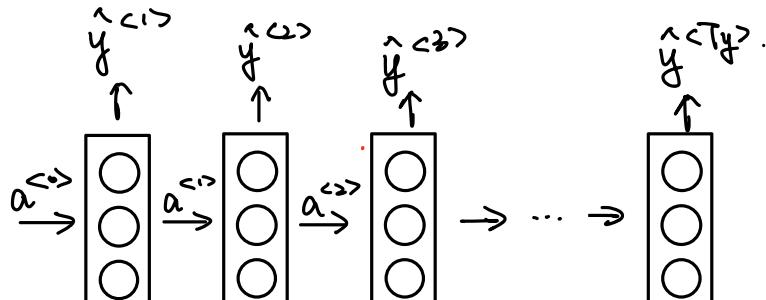
Disadvantage: Too long. Can't catch long range dependencies.

Vanishing gradients with RNNs.

Using RNN, hard to catch this

The cat, which already ate , was full.

The cats, which , were full.



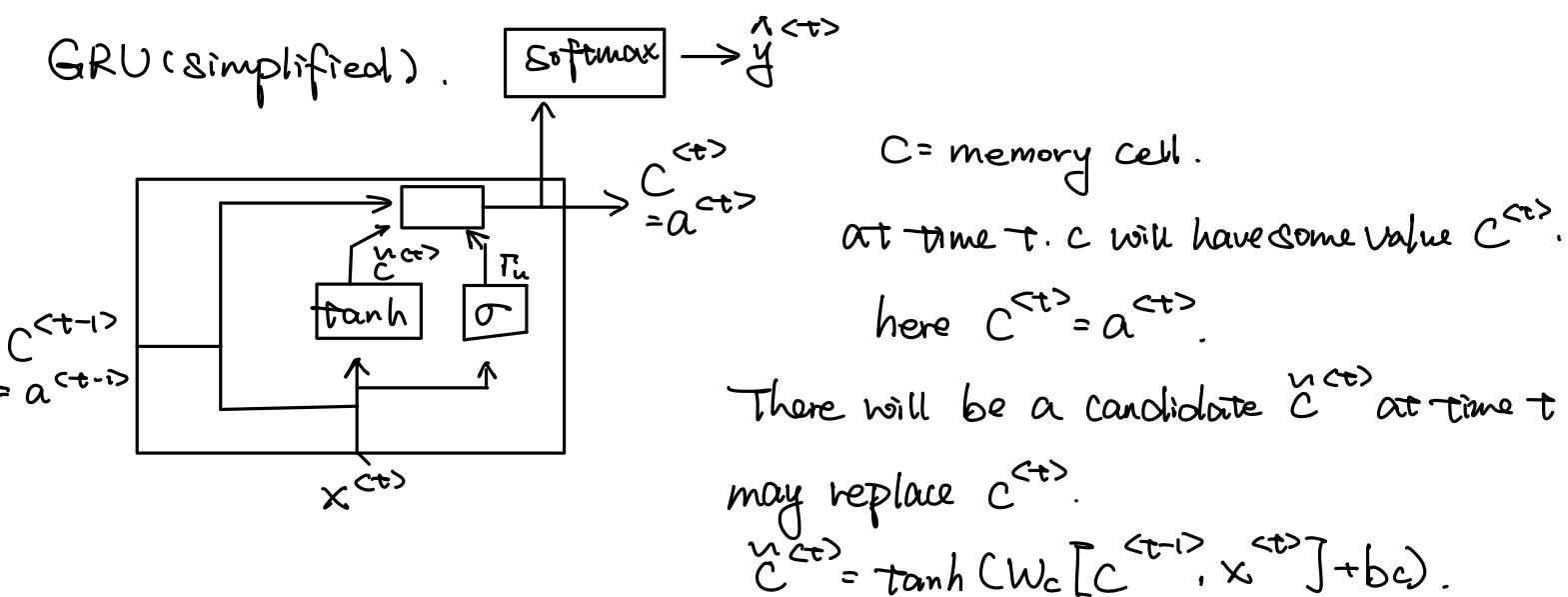
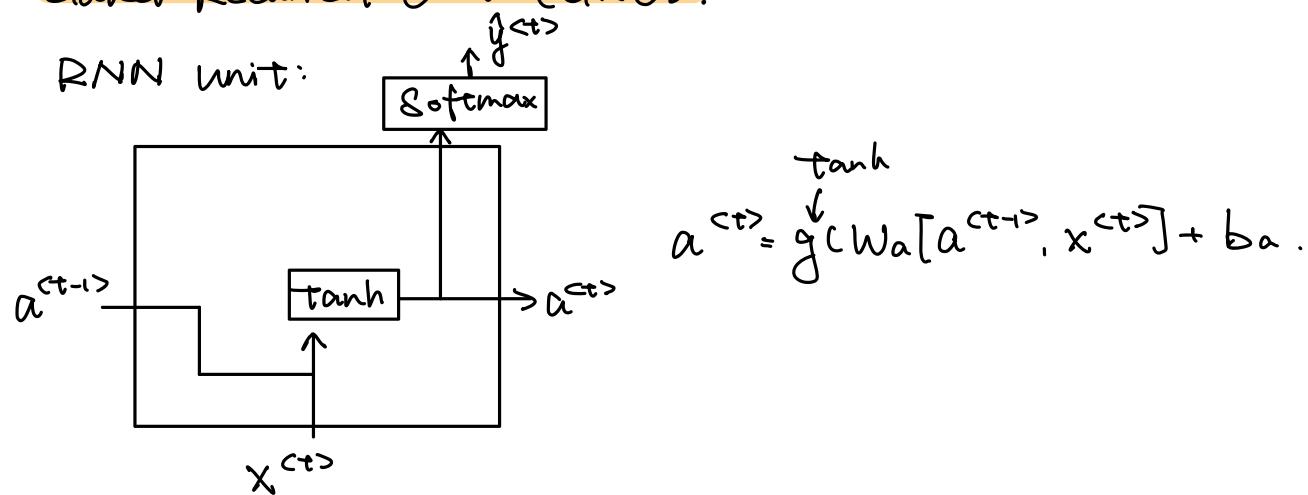
it's hard for the diff. between $y^{(Ty)}$ and $y^{(Ty)}$ to affect early units like $x^{(>)}$.

$$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<t>}$$

By using RNN. $y^{<t>}$ is mainly affected by close items like $x^{<t-2>} \cdot x^{<t-1>} \cdot x^{<t>}$

- * There might also exist exploding gradient problem. We may see NaN values. And we can apply gradient clipping. e.g. if your gradient is greater than some threshold. we just rescale the gradient vec.
- This is easy to solve. Vanishing gradient hard to solve.

Gated Recurrent Unit (GRU).



There's a gate. $\Gamma_u = \text{Sigmoid}(W_u[C^{<t-1>}], x^{<t>}]) + b_u$

$$\Gamma_u = 1 \quad \Gamma_u = 0, 20, \dots = 0, \dots$$

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

Γ_u decide when to change $C^{<t>}$

e.g. The cat, which was full..

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

① if $\Gamma_u=1$. Set $C^{<t>} = \tilde{C}^{<t>}$. Go ahead and update.

② if $\Gamma_u=0$. $C^{<t>} = C^{<t>}$. No update.

Of course: $C^{<t>} \cdot \tilde{C}^{<t>}$. Γ_u can be vector. Thus * represent elementwise product. So that we can update some bits without affect others. And diff. bits for diff. use.

Full GRU. r stands for relevance. So Γ_r tells you how relevant is $C^{<t-1>}$

$$\tilde{C}^{<t>} = \tanh(W_c [\Gamma_r * C^{<t-1>}, x^{<t>}] + b_c) \quad \text{to compute next candidate for } C^{<t>}$$

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

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

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

LSTM. (long short term memory) unit.

Previous equations:

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

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

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

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

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

LSTM. No $C^{<t>}$ and Γ_r .

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

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

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

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

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

Gives the memory cell the choice to keep the original value and add new value

Γ_f : forget gate. Γ_o : output gate.

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

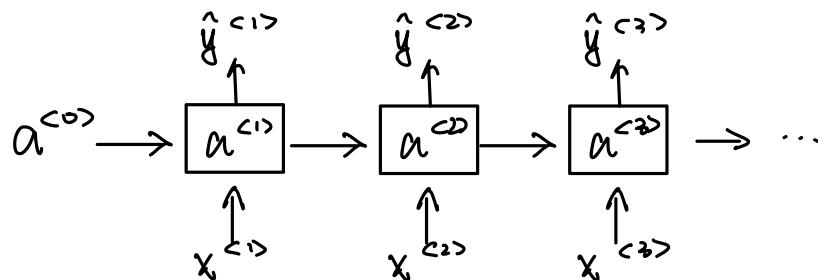
Bidirectional RNN.

Getting information from the future.

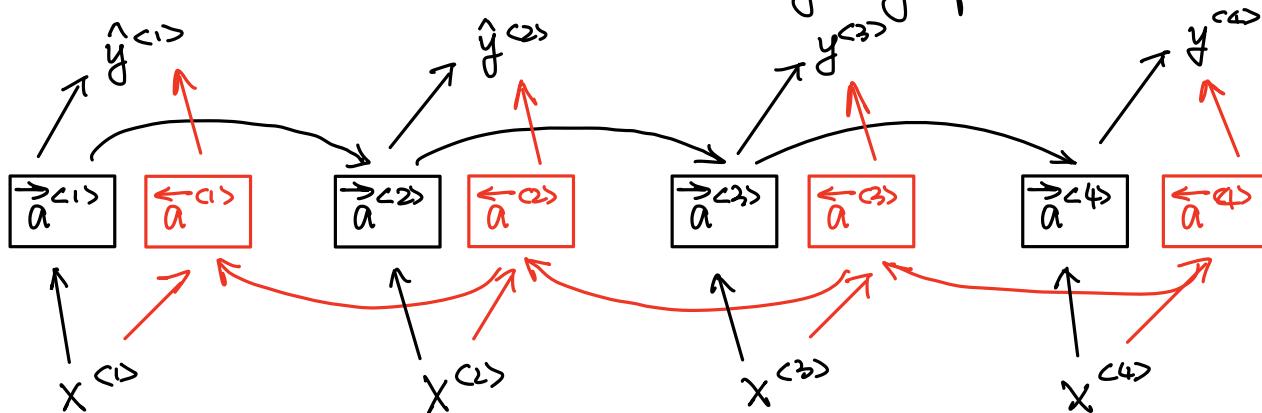
e.g. ① He said : "Teddy bears are on sale!"

② He said : "Teddy Roosevelt was a great president".

We see the word "Teddy". we can't decide it from previous two words : "He said". First two words didn't tell us "Roosevelt" or "bear"



Bidirectional RNN (BRNN).

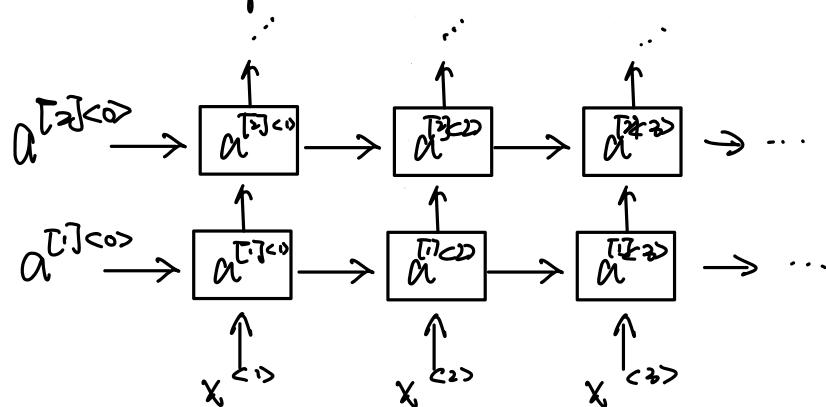


Now we have : $\hat{y}^{<t>} = g(W_y [\rightarrow a^{<t>} \text{, } \leftarrow a^{<t>}] + b_y)$

These blocks can also be GRU and LSTM blocks.

Deep RNNs

Deep RNN example.



Also we can replace blocks use GRU/LSTM blocks

$$a^{T \geq 3} = g(Wa^{T_2} [a^{T_2 \geq 2}, a^{T_1 \geq 2}] + ba^{T_3})$$

Word Representation.

$V = [\text{a, aaron, \dots, zulu, } \in \text{UNKs}]$. $|V| = 10,000$.

1 hot representation. $a = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{10,000}$.

Problem. Taking any two representation. The inner product is always 0.

So talking about similarity. There's no difference between "apple" & "orange" and "queen & orange".

featured representation: word embedding.

e.g. Man Woman King Queen Apple Orange.

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	0.7	0.69	0.03	-0.02
-----	------	------	-----	------	------	-------

Fruit	0.04	0.01	0.02	0.01	0.95	0.97
-------	------	------	------	------	------	------

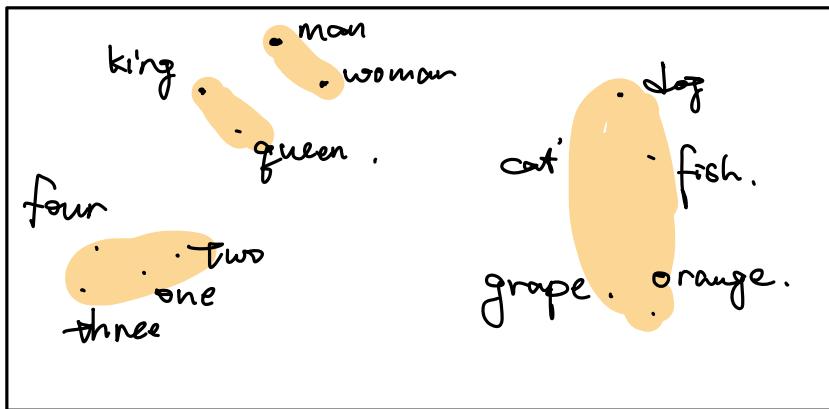
And there are many other types. Allow people to use number to evaluate this thing.

e.g. Maybe there're 300 number of features. And we can use this 300 dim vector to represent words.

So. now it is better to generalize things better.

e.g. Model can learn how to fill words after "apple" by learning from orange.

Visualizing Word Embeddings:

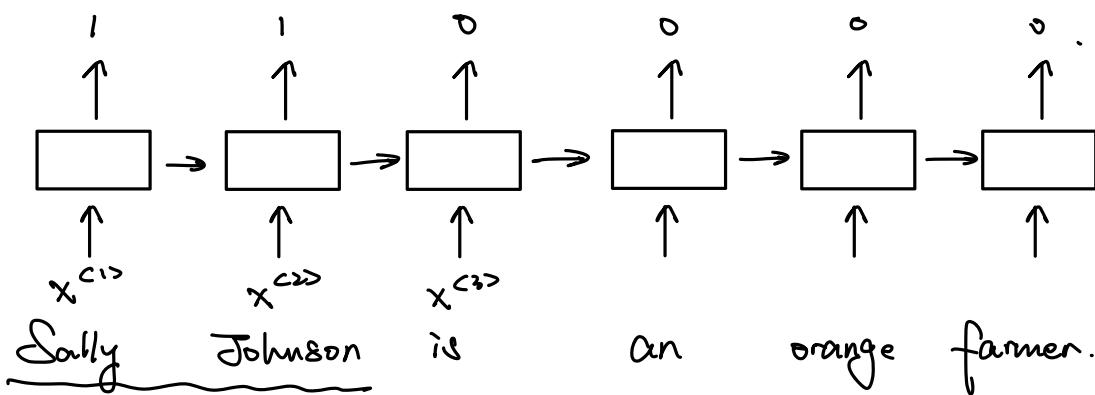


$\rightarrow 3000 \downarrow m \rightarrow 2D$.

t-SNE.

Using word embeddings.

Name entity recognition example: Recog. people's name.



New input. Test.

Robert Lin is an apple farmer.

Using feature representation. it will be easier for the model to recog. the similarity for "apple" and "orange".

In reality. we have tons of text. 1B to 100B words.

You can use them to train. And then apply to your own model.

Transfer the task to you maybe 100k words model.

Transfer learning and word embeddings.

- ① Learn word embeddings from large text corpus. (1B-100B words).
(Or download pre-trained model).

- ② Transfer

③ Optional: Continue to finetune the word embeddings with new data

Properties of word embedding.

Analogies:

know man \rightarrow woman. for king \rightarrow queen. How to know this?

	Man	Woman	King	Queen	Apple	Orange
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	0.7	0.69	0.03	-0.02
Fruit	0.04	0.01	0.02	0.01	0.95	0.97

Recall this:

$$\begin{aligned} e_{\text{man}} &= \begin{bmatrix} -1 \\ 0.01 \\ 0.03 \\ 0.04 \end{bmatrix} & e_{\text{woman}} &= \begin{bmatrix} 1 \\ 0.02 \\ 0.02 \\ 0.01 \end{bmatrix} & e_{\text{man}} - e_{\text{woman}} &\propto \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix} \\ e_{\text{king}} & \quad \quad \quad e_{\text{queen}} & \quad \quad \quad e_{\text{king}} - e_{\text{queen}} &\propto \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}. \end{aligned}$$

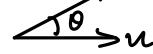
So the diff. are about the same. $e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{queen}}$.

Find word w: $\arg \max_w \text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$.

Cosine similarity.

$$\text{sim}(u, v) = \frac{\vec{u} \cdot \vec{v}}{\|u\|_2 \|v\|_2}$$

This is actually the cosine of the angle between u and v.



Embedding Matrix: How to learn word embedding.

$$\begin{array}{ccccccc} & 6357 & & & & & \\ & \text{orange} & \dots & \text{zulu} & \text{<UNK>} & \dots & \end{array}$$

1000 words in the vocab.



O_{625} : the one-hot encoding vector of "Orange".

Consider $E \cdot O_{625}$. This will select the feature column of "orange" in E .

Embedding for "orange" in the feature space.

Actually this is not efficient. In practice, use specialised function to look up an embedding in E instead of using matrix multiplication.

Learning word embeddings.

Neural Language Model.

I want a glass of orange —

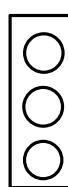
index: 4243 9665 1 3852 6163 6257.

I $O_{4243} \rightarrow E \rightarrow e_{4243}$.

want $O_{9665} \rightarrow E \rightarrow e_{9665}$

a ...

$e_i \rightarrow$



Softmax

glass ...

e_{3852}

of ...

e_{6163}

orange. ...

e_{6257}

Over entire vocab size.
Here we have 10,000 possible outputs.

Of course, since there are too many weights, we may only have a fixed history. E.G: 4 words. We only consider "a glass of orange".

Other text/target pairs.

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

Context . ↓
 target

Context? Last 4 words.

4 words on left & right. ↴

"a glass of orange ? to go alone with"

Last 1 word.

"orange ?".

Nearby 1 word. skip gram model.

Word2Vec.

Skip-grams. (Suppose you have following sentence in your training)

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

Context

orange

orange

orange

Target.

juice

glass.

my.

Model:

Vocab size = 10,000k

$x \xrightarrow{\quad} y$.

Context C ("orange") \longrightarrow Target \hat{y} ("juice").

6257

4834

To represent "orange". We can start from onehot O_{6257} . or $E_{6257} = \vec{O}_{6257}$.

Then we feed E_{6257} into a softmax model. And try to output \hat{y} .

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

θ_t : param associate with output t .

$$L(y, \hat{y}) = -\sum_{i=1}^{10,000} y_i \log \hat{y}_i. \quad (\text{Usual loss for softmax}).$$

$$y = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow \text{index } 4834.$$

So the matrix \hat{t} will have a lot of parameters.

And we keep training it.

Problem with softmax classification.

If $|\text{vocab}| = 10,000$. Computation too large for $p(t|c)$.

Other method: Hierarchical softmax. e.g. 10000 size.



So that don't need to sum over all the 10,000 examples.

Once we choose the context c , we can sample target t in certain window.

How to sample the context c ?

e.g. Sample less usual words. (Not words like "a", "the", "of").

Negative Sampling: Deliberately sample negative examples to train.

Defining a new learning problem. (New supervised learning).

I want x glass of orange juice to go along with my cereal
Context word Target

New learning task:

Pair:

orange	juice	1
--------	-------	---

input: a pair of word.

↑ orange king pick from random

output: will they appear together or not.

↖ Orange book

Normally: $K=5-20$ for small dataset.

↓ orange the

$K=2-5$ for large dataset.

↓ orange of also of appear
 ↑ C ↑ next to "orange" 0

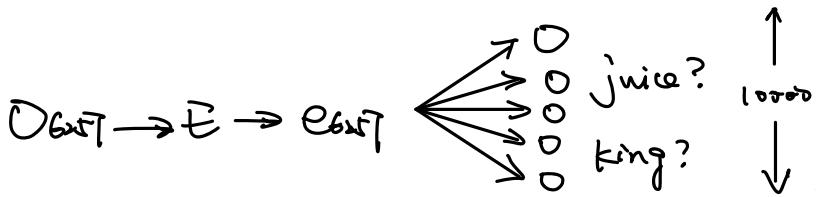
we can think of we have $\frac{1}{K}$ pos&neg sample.

Model:

$$\text{Softmax} = P(t|c) = \frac{e^{\theta_t^\top e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^\top e_c}} \quad (\text{from previous slide}).$$

$$\text{Now: } P(c|y=1|c, t) = \sigma(\theta^\top e_c)$$

Orange
6:st



Turn a 10,000 softmax problem into 10000 binary classification problem
And we only deal k+1 of them.

Selecting Negative example. How.

- ① Sample according to appear frequency. (but sample a lot of "a" "the")
- ② Equal probability. (None representative).
- * ③ $P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$. The frequency of words appear but take power of $\frac{3}{4}$.

GloVe word vectors. (GloVe: global vectors for word representation).

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

x_{ij} : # times i appear in content of j

This captures how often does word i and j appear together.

Model:

$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) (\theta_i^\top e_j + b_i + b_j - (\log x_{ij})^2) + c$$

Weighting term

① $f(x_{ij}) = 0$ if $x_{ij} = 0$

Convention: $\theta \otimes 0 = 0$

② Don't give too much weight for words like "a", "the", "of" ...

Don't give too less weight for infreq. words like "durian"

③ e_i and e_j are symmetric. Need to initialize θ and e at random.

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

A note on the featurization view of word embedding.

Sometimes, the feature representation might be hard for human understanding.

Sentiment Classification.

e.g Problem.

$$x \longrightarrow y.$$

The dessert is excellent.



Service was quite slow.



...

Question: You might won't have a large enough training set.

Simple sentiment classification e.g.

The dessert is excellent



8928 2468 4694 3180.

The $O_{8928} \rightarrow E \rightarrow e_{8928}$

dessert

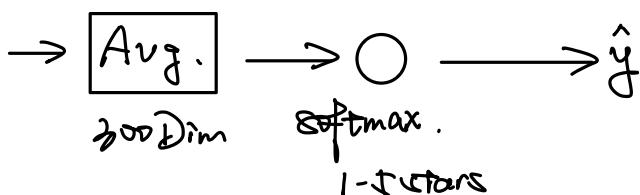
e_{2468}

is

e_{4694}

Excellent.

e_{3180}



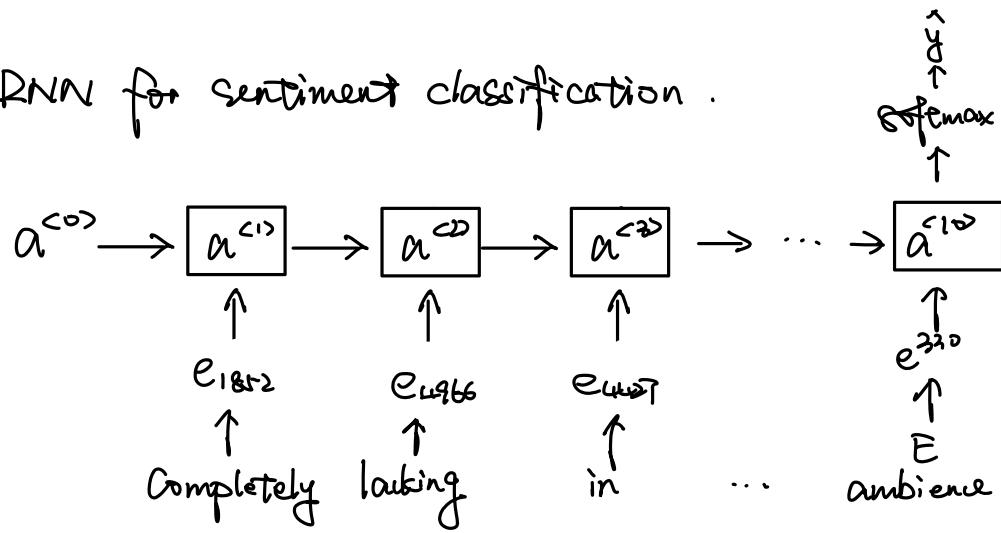
If E was trained on a large data set. E.G. 100B words.

Notice: word sequence doesn't matter.

Consider: "Complete lacking in good taste, good service, good ambience".

Word "good" appear multiple time. Model might think this sentence is a good review. (if we take "Avg" or "Sum").

RNN for sentiment classification.



So this might consider the completely negative effect of word "lacking" performed towards the sentiment of the sentence.

Debiasing word embeddings.

The problem of bias in word embedding.

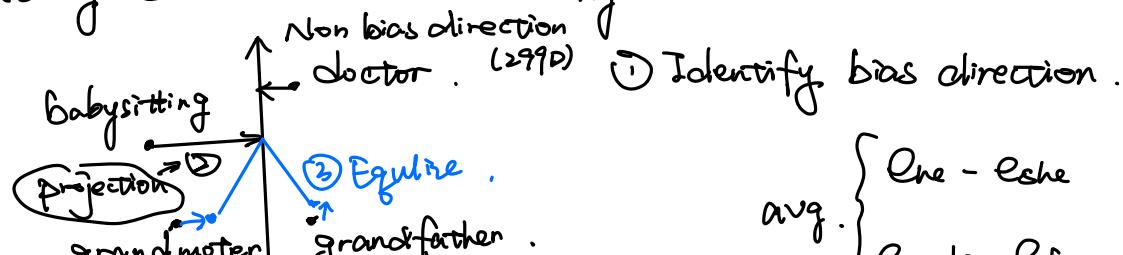
Man : Woman as King : Queen.

Man : Computer-Programmer as Woman : Homemaker X.

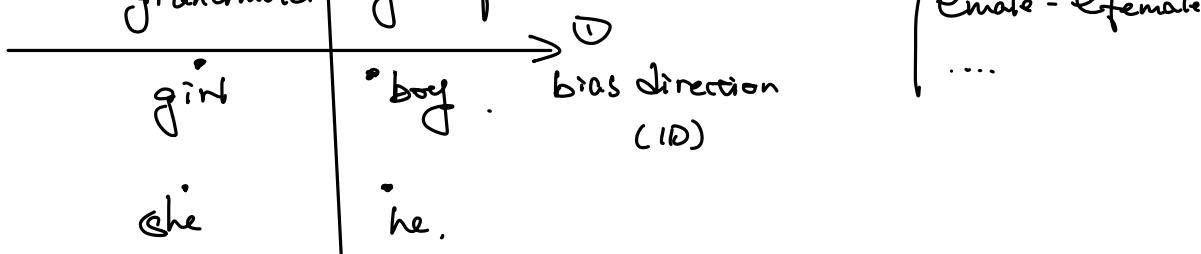
Father : doctor as Mother : Nurse . X .

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

Addressing bias in word embeddings.



avg. { One - Eshe
P1 - P2 .



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

Most words are not \leftarrow Not like words "grandmother" "father" which are definitional. And shall have specific gender info. Like word "doctor" or "nurse" which should not include gender info.

\downarrow
Train a linear classifier to decide.

③ Equivale pairs: e.g. $\text{Dist}(\text{grandmother}, \text{babysit}) \approx \text{Dist}(\text{grandfather}, \text{babysit})$

This is relative small. Can just hand pick them.

Basic Models.

Seq-to-Seg model.

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$

Jane visite l'Afrique en septembre

\rightarrow Jane is visiting Africa in September.

$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>}$

How to train a neural net can do this?

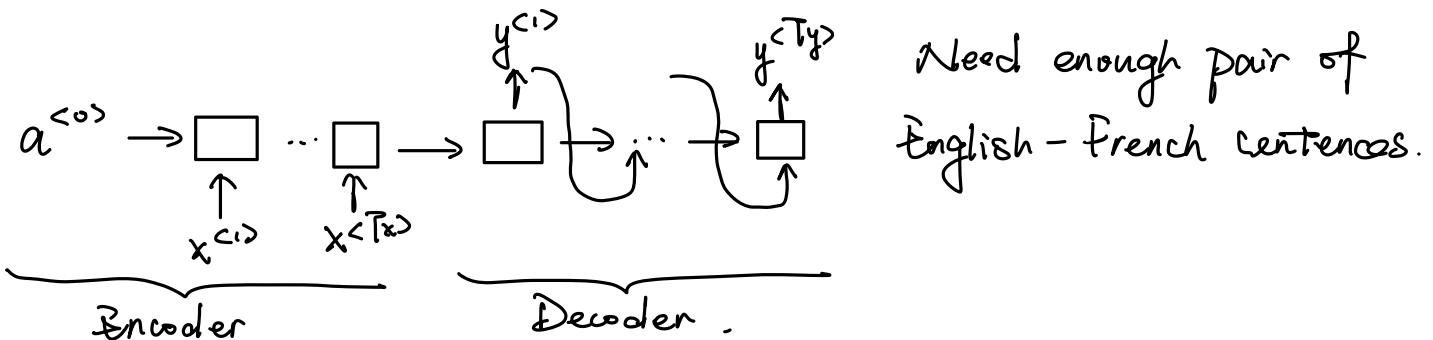
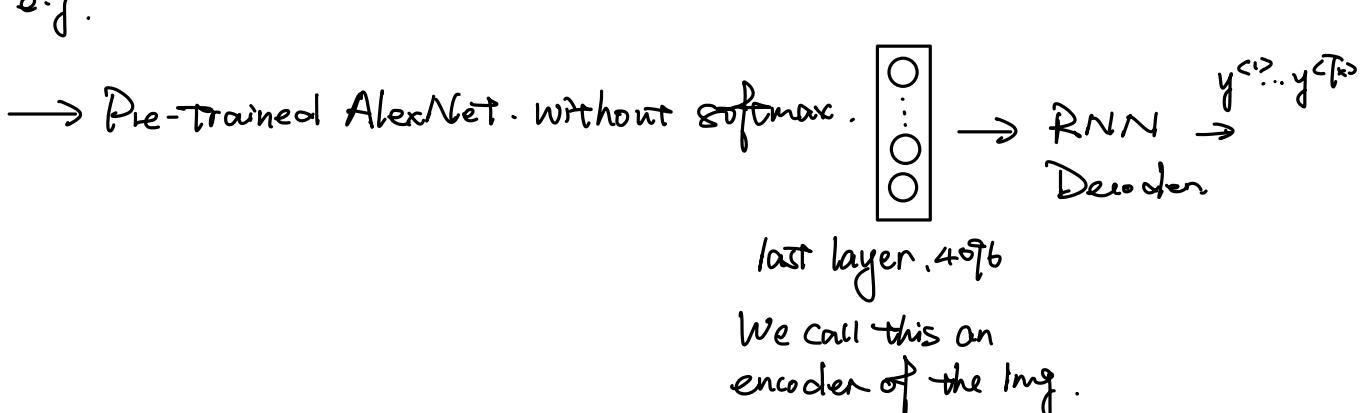


Image Captioning.

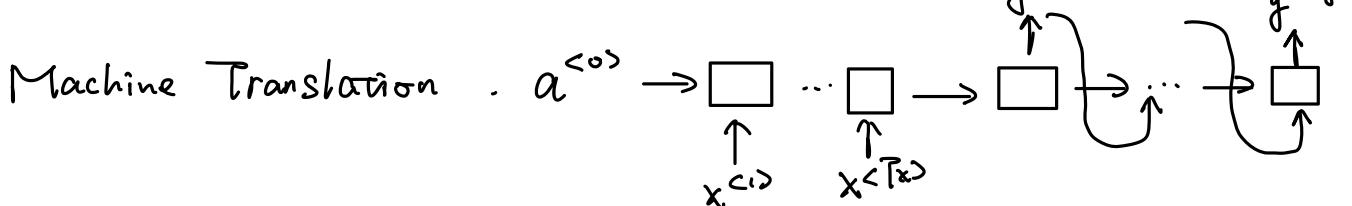
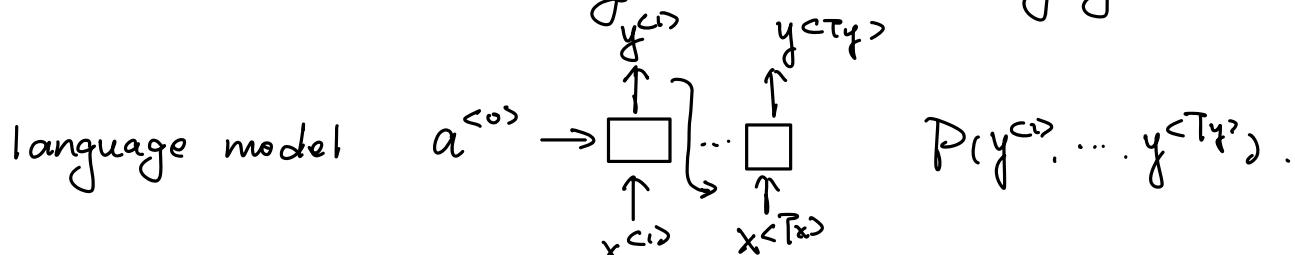
Cat

Some img.



Picking The Most Likely Seq.

Machine translation as building a conditional language model.



The decoder part really looks like the language model.

But "conditional" on the encoder part. $P(y^{<1>} \dots y^{<Ty>} | x^{<1>} \dots x^{<Tx>})$.

Finding the most likely translation.

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

English. French

Jane visite l'Afrique en septembre

- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in Sept.
- In Sept. Jane will visit Africa.
- } argmax $P(y^{<1>} \dots y^{<Ty>} | x)$.

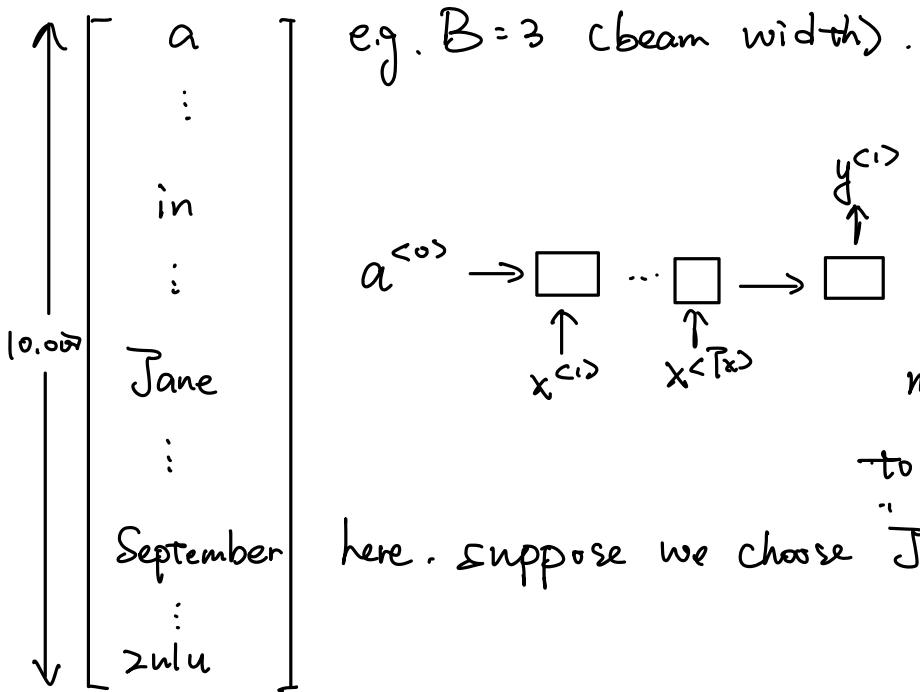
Why not greedy search. $P(y^{(1)}|x)$

Pick the most likely first word. Then most likely second, and ...

We need a search algorithm instead of listing all the sentences.

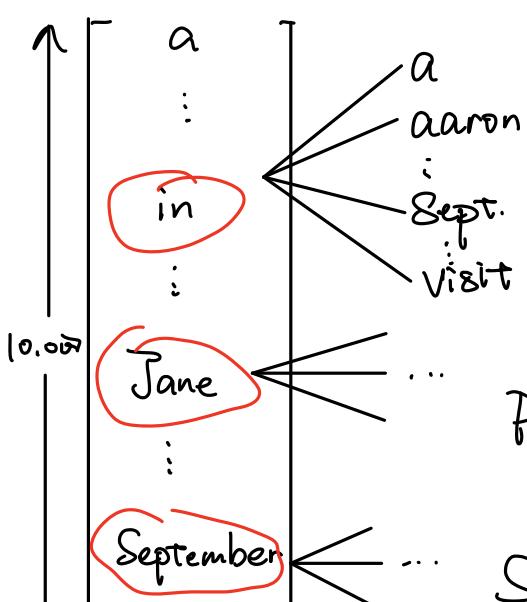
Beam Search.

Step 1.



First thing generate the most likely first words.
Instead of greedy pick the most likely, but have a param B to get the number B of most likely here. Suppose we choose "Jane", "in", "September"

Step 1. Step 2.



Suppose here we use first word "in".
 $y^{(1)} = \text{in}$ $y^{(2)}$. $P(y^{(2)}|x, \text{"in"}).$

Actually we want the most likely pair.

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

Since we use $B=3$. Step 2 we need to consider all

↓ [2nd] 30,000 possibilities. And we pick the most likely 3.

e.g. Say: "in sept.". "Jane is". "Jane visits".

Also. we reject "Sept." as the first word.

Everytime we have number of B of copies of the network. And we do thorough search over all vocab. And extend network. Cont.

e.g. Step three: "in Sept. Jane". "jane is visiting". "Jane visits Africa" finally we may end up with "Jane visits Africa in Sept.".

Refinement to Beam Search.

Length Normalization.

$$\underset{y}{\operatorname{argmax}} \frac{\prod_{t=1}^{Ty}}{P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})}.$$

This number always smaller than 1.

So we consider: (first optimization)

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

Notice: This equation tends to prefer short sentences!

So we add a normalize term.

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

Where α is a hyperparam you can tune. ① If $\alpha=1$. normalizing over length. ② $\alpha=0$. No normalization.

Beam Search discussion.

Beam width B:

① large B. Better result, slower

② Small B. worse result. faster.

in practice. we use 10. production 100. Not common for 1000↑

Remember. Beam search is not guaranteed to find the exact max!

Error Analysis on Beam Search.

Example:

Jane visite l'Afrique en septembre

Human: Jane is visiting Africa in September (y^*)

Alg: Jane visited Africa last Sept. (\hat{y}).

Increasing the Beam Width barely hurts.

RNN computes $P(y|x)$.

What you can do. is use the RNN model. to compute $P(y^*|x)$ and $P(\hat{y}|x)$. This help you to identify the error.

Case ① $P(y^*|x) > P(\hat{y}|x)$. (human perform better).

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

Conclusion: Beam Search is at fault.

Case ②. $P(y^*|x) < P(\hat{y}|x)$.

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

Conclusion: RNN model is at fault.

Error Analysis Process.

Human .	Alg (Give sth. wrong)	$P(y^* x)$	$P(\hat{y} x)$, At fault.
Jane visits Africa in Sept.	Jane visited Africa last Sept.	2×10^{-6}	1×10^{-10} B.

Continue go through all the bad example - the Alg. returned .

We can figure out what fraction of errors are "due to" beam search V.S. RNN model .

Bleu Score. How to evaluate when there are multiple good result?

French: Le chat est sur le Tapis .

Reference ①: The cat is on the mat. } Both
Reference ②: There is a cat on the mat. } Correct.

Bleu: Bilingual evaluation, underway .

Extreme e.g. MT output: the the the the the the . (7 words)

Precision: $\frac{T}{7}$. whether the words in the output appear in reference .

Which is not a good evaluation .

Modified precision: "the" appear in ① two time . in ② one time .

So the upper bound is 2.

$\frac{2}{7} \leftarrow$ Count clip ("the") in reference .
 $\frac{1}{7} \leftarrow$ Count ("the") in output

Bleu score on bigram .

Reference ①: The cat is on the mat .

Reference ②: There is a cat on the mat

MT output: The cat the cat on the mat . (\hat{y})

Possible bigram . Count : Count Clip .

"the cat" 2 1

"cat -the"	1	0	Modified precision: $\frac{\sum \text{CountClip}}{\sum \text{Count}} = \frac{4}{6}$.
"cat on"	1	1	
"on -the"	1	1	
"-the mat"	1	1	

Bleu score on unigram.

$$P_1 = \frac{\sum_{\text{unigram} \in \hat{y}} \text{CountClip}(\text{unigram})}{\sum_{\text{unigram} \in \hat{y}} \text{Count}(\text{unigram})}$$

Bleu score on n-gram.

$$P_n = \frac{\sum_{n\text{-gram} \in \hat{y}} \text{CountClip}(n\text{-gram})}{\sum_{n\text{-gram} \in \hat{y}} \text{Count}(n\text{-gram})}$$

This allows you to evaluate the "overlap" of references and the predicted output.

Bleu Details.

P_n = Bleu score on n-gram only. e.g. P_1, P_2, P_3, P_4 .

Combine Bleu score: $BP \left(\exp \left(\frac{1}{4} \sum_{n=1}^4 P_n \right) \right)$

BP = brevity penalty. Very short output, easy get high score.

This penalize short output.

$$BP = \begin{cases} 1 & \text{if MT output length} > \text{reference length} \\ \exp(1 - \text{MT output length} / \text{reference length}) & \text{otherwise.} \end{cases}$$

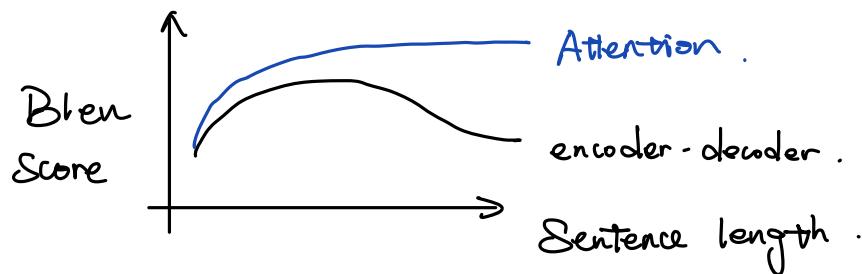
Attention model intuition.

The problem of long seq.

Suppose we have a long French sentence. and want to translate to English.

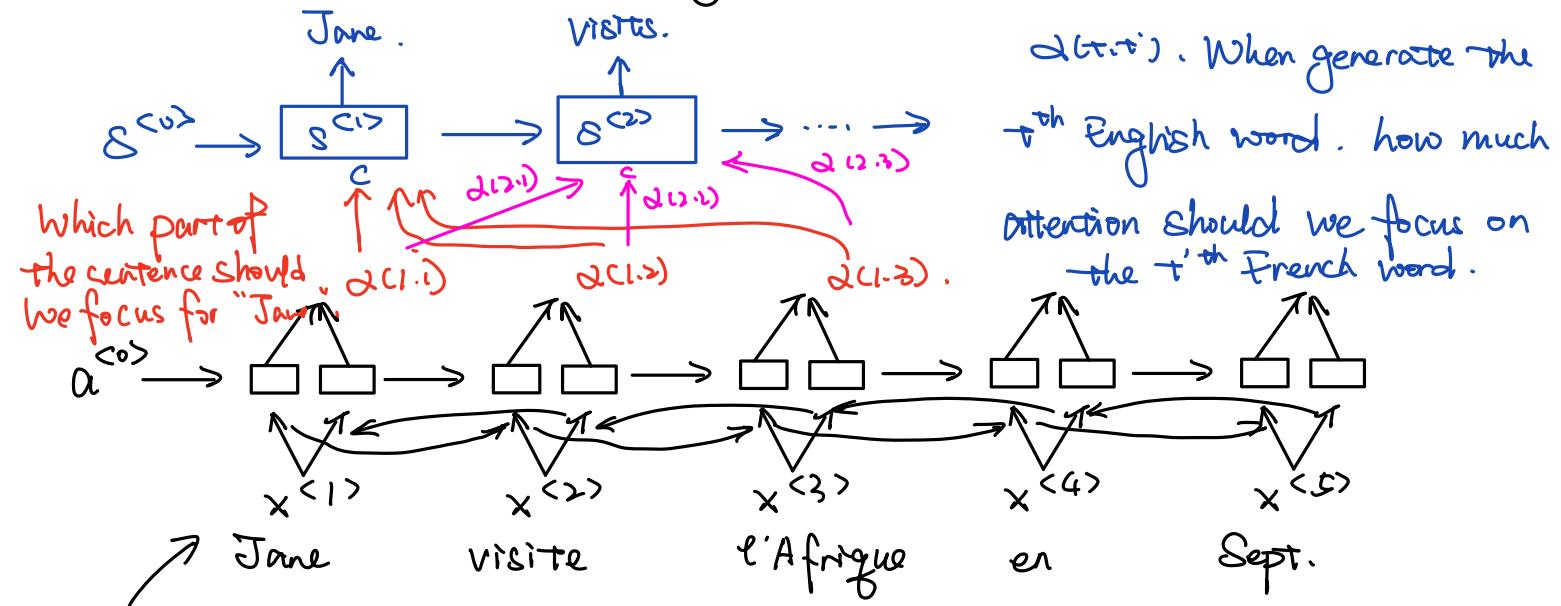
What human translation did is translate part by part. But current

encoder-decoder model need to memorize it all.



Attention model intuition. Illustrate with short sentence).

And we use another RNN to generate translation.



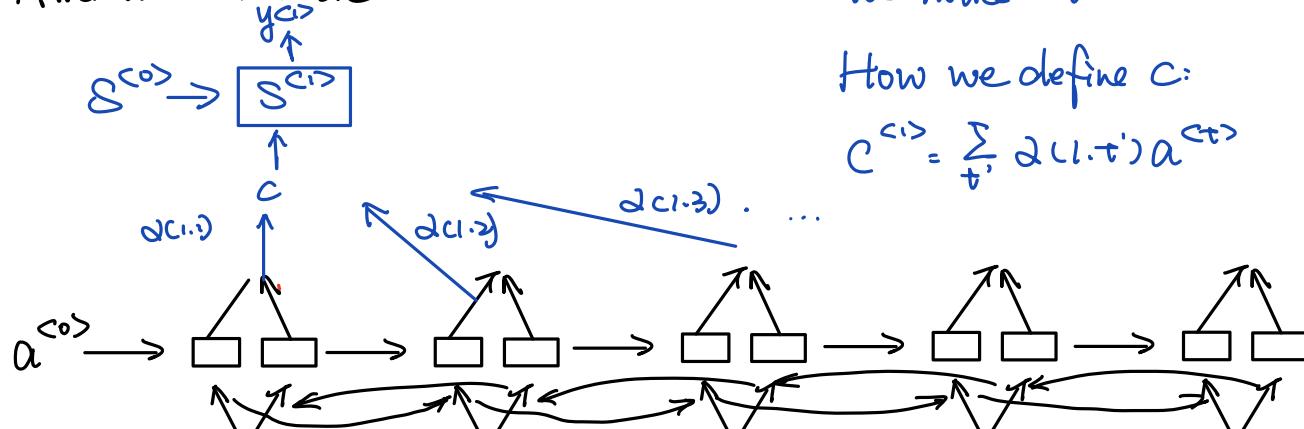
Suppose we use a BRNN (May be BiLSTM, BiGRU),

Attention model gives a set of weight. how much you should focus.

e.g. first word "Jane". Combine all the attention weights. $\alpha^{(1,1)} \sim \alpha^{(1,5)}$

Also there're weights for $\alpha^{(2,1)}, \alpha^{(3,3)}, \dots$

Attention Model .



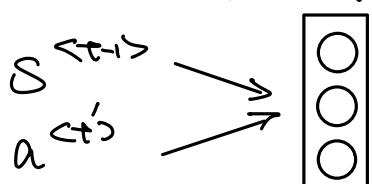
$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$
 Jane visite l'Afrique en Sept.

At every timestamp, use $\alpha^{<\leftrightarrow>} = (\overset{\rightarrow}{\alpha}^{<\leftrightarrow>} \cdot \overset{\leftarrow}{\alpha}^{<\leftrightarrow>})$
 forward backward.

Computing attention $\alpha^{<t,t'>}$. (Same as $\alpha^{(t,t')}$ above).

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

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



how to get e ? Use a small network.

$e^{<t,t'>}$ should depend on $S^{<t-1>} \cdot a^{<t'>}$. How? Don't know

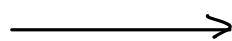
So we use GRU.

Downside of this Alg: Runtime: $(Tx \cdot Ty)$, Quadratic cost.

Speech recog.

Problem.

x



y

audio clip

transcript.



\longrightarrow "The quick brown fox".

Airpressure across time.

we used to use phonemes. de kwik brown.

basic sound units:

We can use Attention model for speech recog.

CTC cost for speech recog.

(Connectionist temporal classification)

