

Course 5

Sequence Model

Sequence Model

Sequence Model.

X: "Harry Potter" and "Hermione Granger" invented a spell

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

$x^{<9>}$

$T_x = 9$

y: 1 1 0 1 1 0 0 0

$y^{<1>} \quad y^{<2>}$

$y^{<9>}$

i-th training
example

$x^{(i)<t>} \quad T_x^{(i)}$

$T_y = 9$

$y^{(i)<t>} \quad T_y^{(i)}$

Representation

a		$x^{<1>}$	$x^{<2>}$
Aaron	1	0	0
:	2	:	:
and	:	0	0
harry	4075	1	0
:		0	1
potter	6830	1	0
:		0	0
Zulm	20000	0	0

10,000

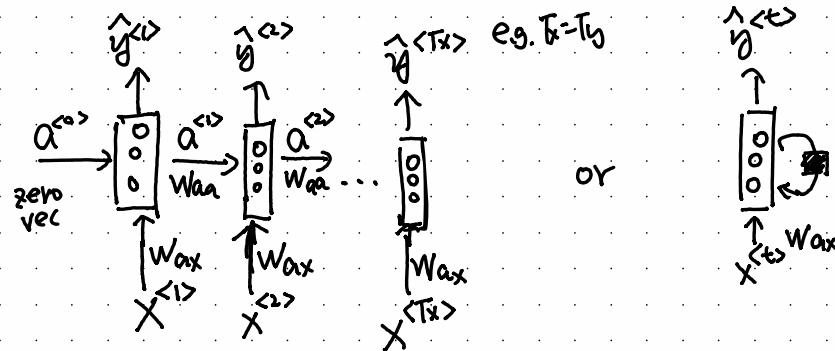
one-hot.

~30,000 for
commercial software

Recurrent NN

why not standard network

- input / output could be of different length
- doesn't share features across different pos.



/imitation: only use info "before" t

Δ B-RNN

$$a^{(0)} = 0 \quad a^{(1)} = g(W_{aa}a^{(0)} + W_{ax}x^{(1)} + b_a) \leftarrow \text{tanh, ReLU}$$
$$\hat{y}^{(1)} = g(W_{ya}a^{(1)} + b_y) \leftarrow \text{sigmoid}$$

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

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

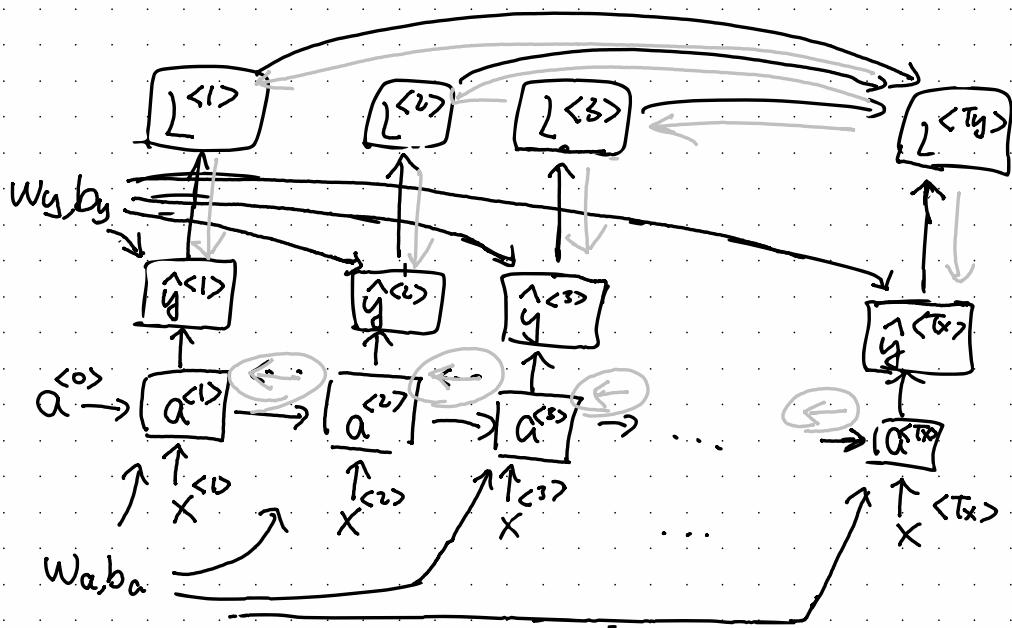
$$\hat{y}^{(t)} = g(W_{ya}\hat{a}^{(t)} + b_y)$$

↓ simplify

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

$$\begin{matrix} 100 \\ 100 \end{matrix} \left[\begin{matrix} W_{aa}; W_{ax} \\ 100; 10000 \end{matrix} \right] = \underline{W_a} \quad \left[\begin{matrix} \hat{a}^{(t-1)} \\ \hat{x}^{(t)} \end{matrix} \right] \begin{matrix} 100 \\ 10,000 \end{matrix} \quad \begin{matrix} 10,000 \\ 10,100 \end{matrix}$$

Backprop. through time

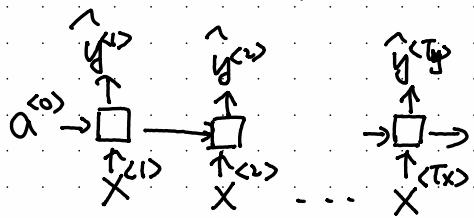


$$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -\hat{y}^{<t>} \log \hat{y}^{<t>} - (1 - \hat{y}^{<t>}) \log (1 - \hat{y}^{<t>})$$

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

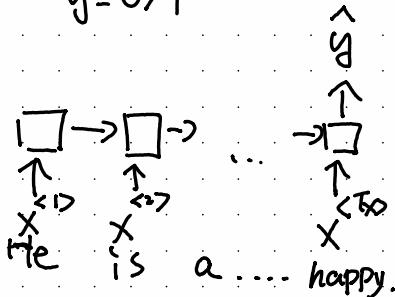
Different RNN types

$$T_x = T_y$$



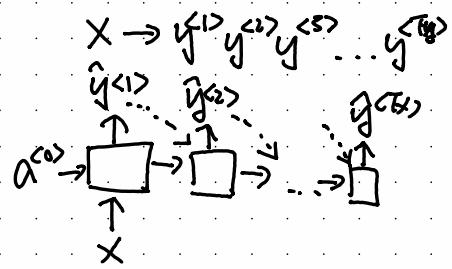
many-to-many

e.g. sentiment detection
 $X = \text{text}$
 $Y = 0/1$



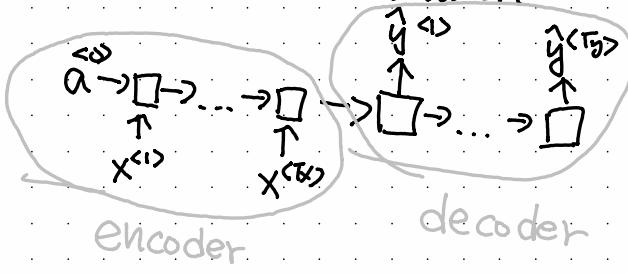
many-to-one

music generation



One-to-many

machine translation



many-to-many

Language Model

Speech Recognition

the apple and pair salad

the apple and pear salad \leftarrow higher p, although both sound the same

language model: $P(y^{<1>} \dots y^{<t>})$?

Training Set: large corpus of text.

- An apple a day ... doctor away. <EOS>

$y^{<1>} y^{<2>} y^{<3>} \dots y^{<t>} y^{<1>}$ Tokenize

- The Egyptian Mau is a type of cat

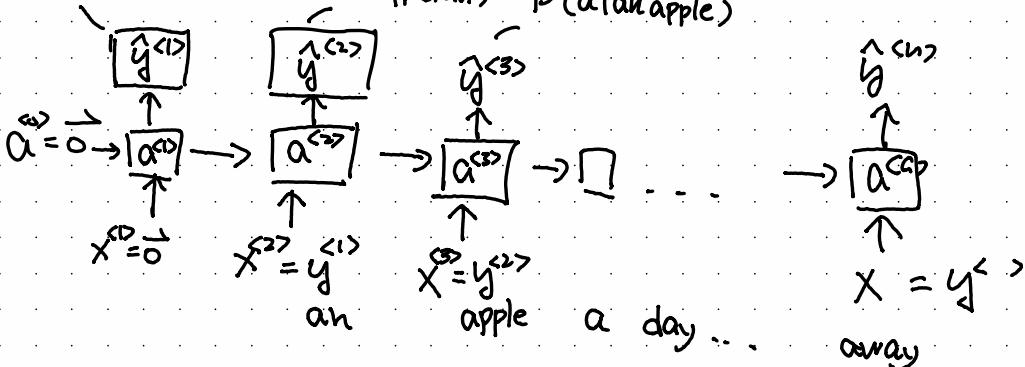
corpus

0	a
:	ab
0	acat
:	zebra
0	

\downarrow replace
<UNK>

unique token for unknown word.

$$P(\text{apple}|\text{an}) \quad P(\text{an}|\text{apple})$$



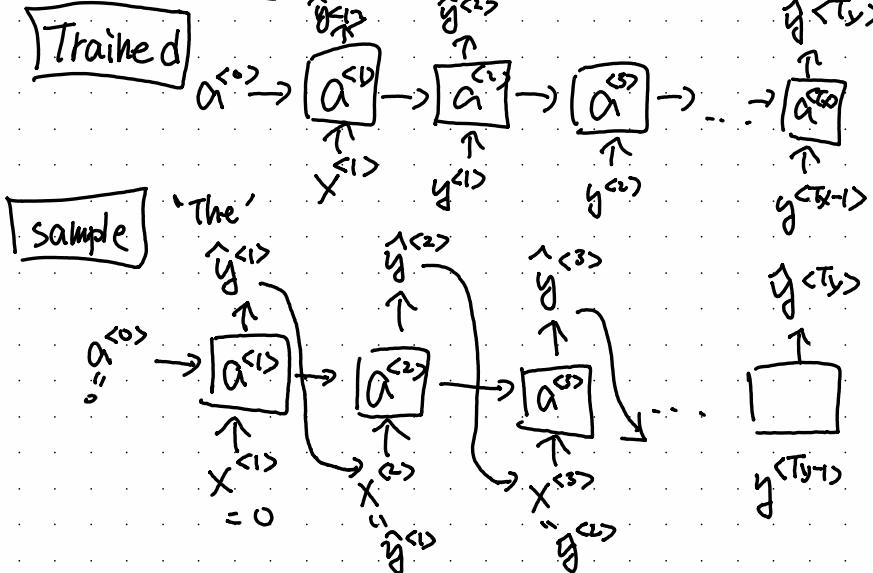
$$P(y^{<1>} | y^{<2>} | y^{<3>})$$

$$\mathcal{L}(y^{<t>} | y^{<t>}) = -\sum_i y_i^{<t>} \log p_i^{<t>}$$

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

- generate random sentence

Sampling a seq from trained RNN



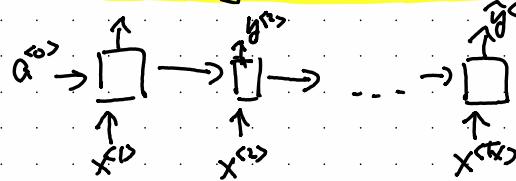
< np.random.choice >

- can also do "character-level" language model

- very long
- slower to train

- word level model is more popular.

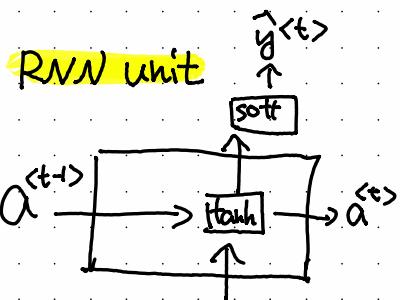
Vanishing gradients with RNN



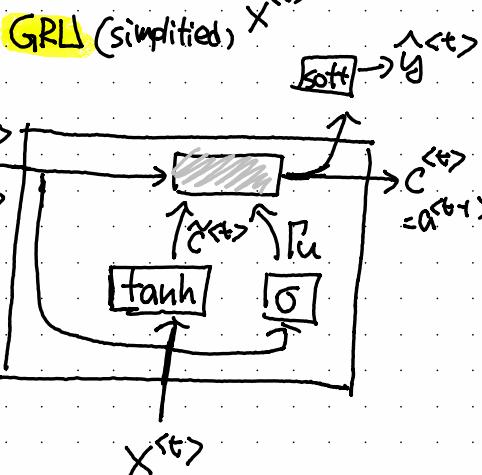
e.g. cats, which ~~are~~ are cute
is cute
- RNN cannot memorize long-term influence

exploding gradient - solve with clipping

GRU - Gated Recurrent Unit, 2014



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



C = memory cell

$$\underline{C}^{(t)} = a^{(t)}$$

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

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

update

$$\underline{C}^{(t)} = \Gamma_u \times \tilde{C}^{(t)} + (1 - \Gamma_u) \times \underline{C}^{(t-1)}$$

(gamma)

the cat, which ate .. . is full

$$\begin{aligned} C^{(t)} &= 1 \dots = \\ \Gamma_u &= 0, 0, \dots, 1 \end{aligned}$$

GRL (Full)

simpler, and can build bigger model

$$\tilde{h} \quad \tilde{C}^{(t)} = \tanh (W_c [\tilde{r}_r \times C^{(t-1)}, x^{(t)}] + b_c)$$

$$u \quad \tilde{r}_u = \sigma (W_u [C^{(t-1)}, x^{(t)}] + b_u)$$

$$r \quad \tilde{r}_r = \sigma (W_r [C^{(t-1)}, x^{(t)}] + b_r) \quad - \text{relevance}$$

$$h \quad C^{(t)} = \tilde{r}_u \times \tilde{C}^{(t-1)} + (1 - \tilde{r}_u) \times C^{(t-1)}$$

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

LSTM, 1997 more complex, but more powerful.

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

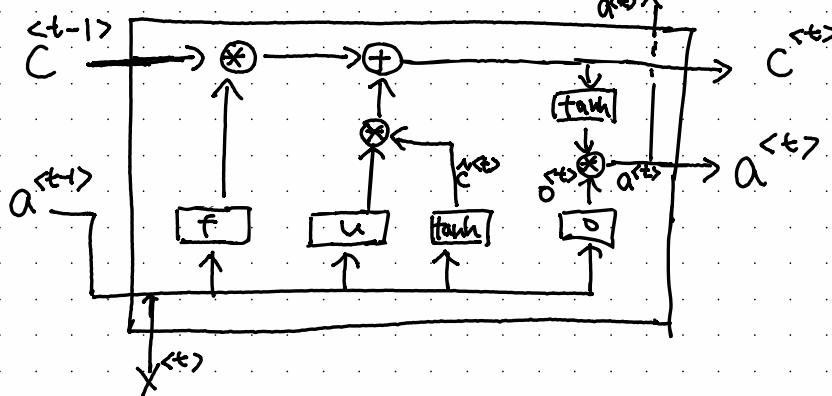
$$\tilde{r}_u = \sigma (W_u [a^{(t-1)}, x^{(t)}] + b_u) \quad - \text{update} \quad \begin{matrix} \text{peephole connection} \\ C^{(t-1)} \text{ affects} \end{matrix}$$

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

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

$$C^{(t)} = \tilde{r}_u \times x^{(t)} + \tilde{r}_f \times C^{(t-1)}$$

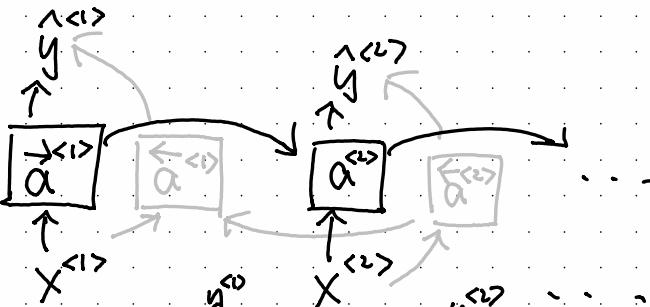
$$a^{(t)} = \tilde{r}_o \times \tanh (C^{(t)})$$



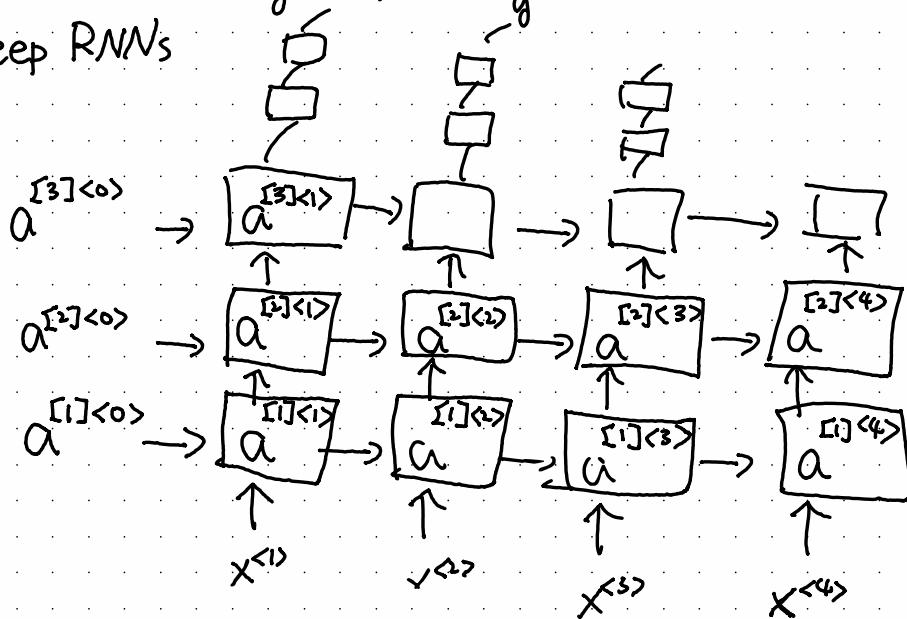
$\tilde{r}_u, \tilde{r}_f, \tilde{r}_o$

Bidirectional RNN (BRNN)

He said Teddy bears are cute
getting into from future
He said Teddy Roosevelt was a good man



Deep RNNs



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

Word Embeddings

one-hot representation

$$32 \rightarrow \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 0 \\ ? \\ ? \end{bmatrix} \leftarrow \text{dot prodct is } 0, \text{ doesn't reflect semantic distance}$$

apple orange

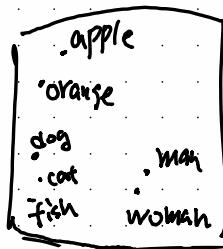
featurized representation: word embedding

	man	woman	apple	orange	...
gender	-1	1	0	0	
age	0.02	0.03	0.02	0.03	
size	0.4	0.4	0.01	0.02	
food	0.01	0.02	0.95	0.94	
e.g.					
word					

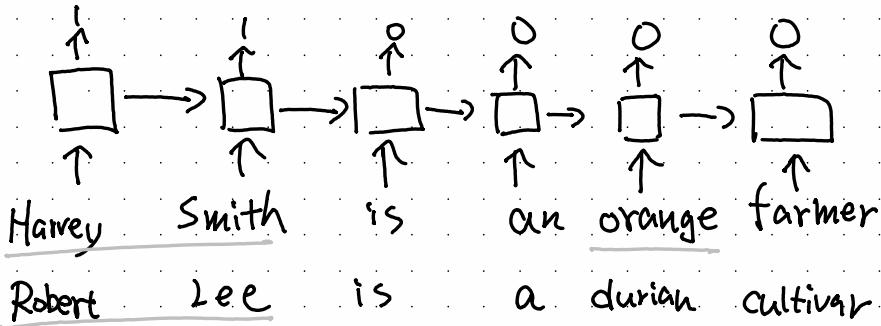
Vishalize word embedding
embed 3D to 2D for viz

H-SNE

"embed" words into feature space



ex: Name entity recognition



↑
name not corporate name

↑
what if we have fruit
not trained before?

Transfer Learning

1. learn from large text corpus
(or download pre-trained embedding online)
2. transfer embedding to new tasks w. smaller set
3. (optional) continuously finetune word embedding

properties of word embedding, 2013 Mikolov

e.g. man → woman

king → ?

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

find word w

$$\arg \max_w \underline{\text{sim}}(E_w, E_{\text{king}} - E_{\text{man}} + E_{\text{woman}})$$

30~75% accuracy in papers.

	man	woman	king	queen
gender	:	:	:	:
age	:	:	:	:
royal	:	:	:	:
:	:	:	:	:

Cosine Similarity

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

$$\text{man:woman} = \text{boy:girl}$$

$$\text{RMB:China} = \text{Rupee:India}$$

Embedding Matrix

300
(features)

E

1,000,000
(words)
corpus

$E \cdot O = []_{(300, 1)}$

One-hot vec

1,000,000

The diagram illustrates the multiplication of an Embedding Matrix E and a One-hot vector O . The Embedding Matrix E has dimensions 300 (features) by 1,000,000 (words and corpus). The One-hot vector O has dimensions 1,000,000 by 1. The resulting product $E \cdot O$ is a single column vector with dimensions 300 by 1, representing the embedding for word j .

$E \cdot O_j = e_j$ - embedding for word j
(in practise, AP2 extract the column)

Learning Word Embeddings

history: difficult \rightarrow simple

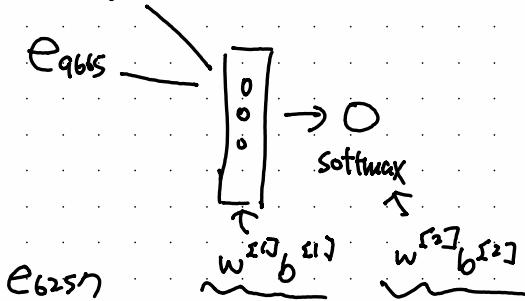
Neural Language Model, 2003, Bengio

$$I \quad O_{4343} \rightarrow E \rightarrow e_{4343}$$

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

want
a
glass
of
Orange

Juice



$$\text{learn } E, W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]} \quad 4 \times 300 = 1200 \text{ dim}$$

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

Context:

last 4 words

last / next 4 words

last 1 word

"nearby" 1 word (skip gram)

Word2Vec, 2d3, Mikolov

skip-gram

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

<u>context</u>	<u>target</u>	<u>randomly chose</u>
orange	juice	
Orange	glass	
Orange	my	e.g. vocab size = 10000.

Model

$$x \rightarrow y$$

Context c → Target t
"orange" "juice"

$$o_c \rightarrow E \rightarrow e_c \rightarrow o \rightarrow \hat{y}$$

$$\text{softmax } e^{o_t^T e_c} \\ \text{softmax: } P(t|c) = \frac{e^{o_t^T e_c}}{\sum_{j=1}^{10000} e^{o_j^T e_c}}$$

o_t : parameter associates with output t
i.e. target t being the label

slow. so,

$$L(\hat{y}, y) = - \sum_i y_i \log \hat{y}_i \quad y = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}$$

how to sample context c ?

too ~~common~~, the, a, of, ...

→ orange, apple...

$P(c)$ based on frequency



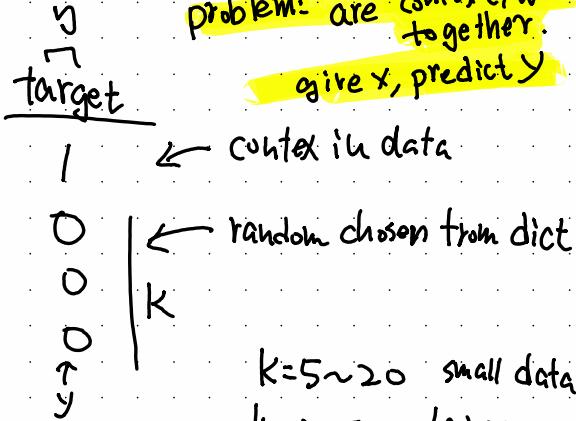
Negative Sampling

<u>Context</u>	<u>word</u>
Orange	juice
Orange	king
Orange	book
Orange	of
c ↑ t	t

2013, Mikolov.

problem: are context/word together.

give x, predict y



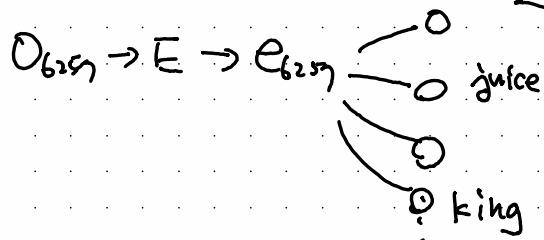
Model

softmax:

$$P(c|t) = \frac{e^{\theta_c^T e_t}}{\sum_j e^{\theta_c^T e_j}} \quad \leftarrow 10,000 \text{ softmax, slow!}$$

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

orange: O_{6257}



10,000 binary classification

each iteration, train K+1 classifier

Sampling:

freq. $P(w_i)$

$$\Rightarrow P(w_i) = \frac{f(w_i)^{\frac{3}{4}}}{\sum f(w_i)^{\frac{3}{4}}} \quad \text{heuristic}$$

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

Glove Word, 2014. Pennington.
(global vector)

C, t

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

$$x_{ij}^* = X_{ij}$$

$$\min \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) (\Theta_i^T e_j + b_i + b_j' - \log x_{ij})^2$$

t c
↓ ↓
weighting

$$f(x_{ij}) = 0 \text{ if } X_{ij}=0 \quad \text{if } \log 0 = 0$$

Θ_i, e_j are symmetric

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

this is...
during
zulu

note:

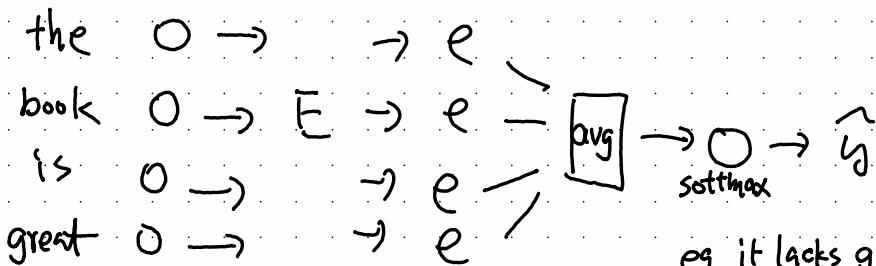
~~✓~~ e_i might not align with defined interpretable feature

axis of the embedding matrix

Sentiment classification

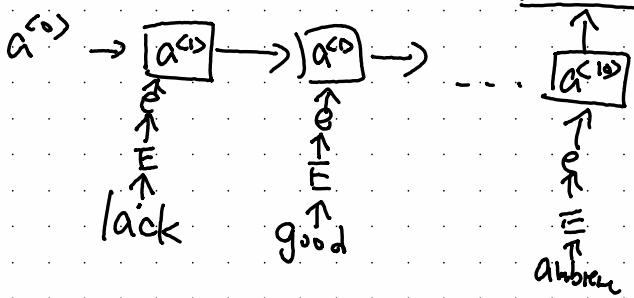
X		Y	
the book is great		****	data set is smaller.
the book lacks depth		*	
reads Ok		**	
:			

Model 1.



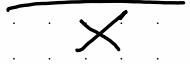
eg. it lacks good insight,
good depth, good taste
 problem!

Model 2 - RNN



many-to-one

Debiasing Word Embeddings, 2016.

Man = Computer Programmer as Woman = Homemaker 

Word embedding reflect gender, ethnicity used to train the model.

1. identify

identify bias axis

2. neutralize

project onto axis except definitive terms

3. equalize

Move definitive pairs to same distance around axis

Seq to Seq Model 2014

e.g. translation

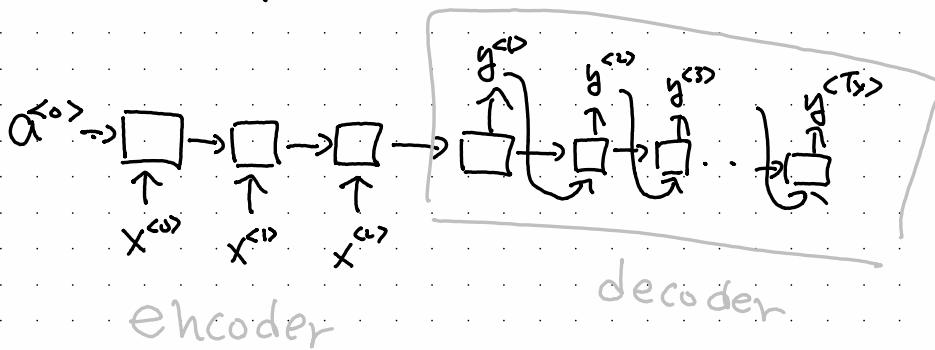
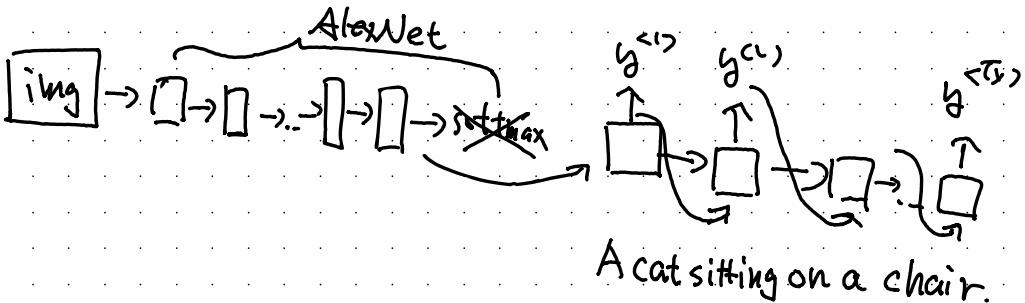


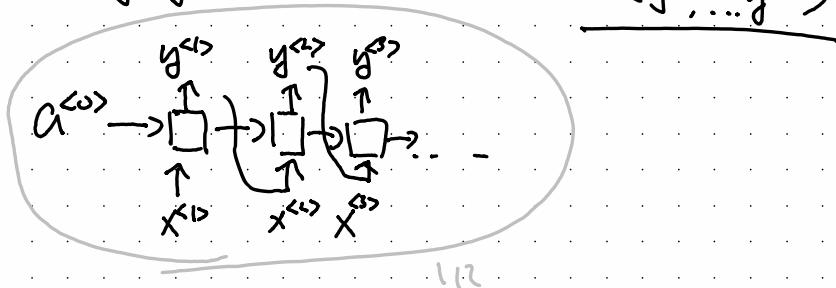
Image Captioning, 2014, 2015



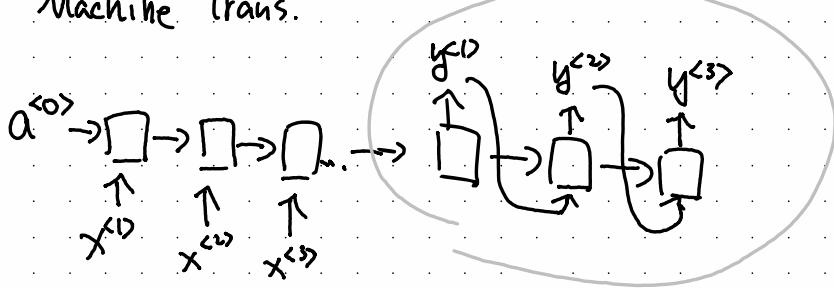
seq2seq vs language model

machine trans as building conditional language model.

Language Model:



Machine Trans.



Conditional
language model

$$\frac{P(y^{<1>} \dots y^{<n>} | x^{<1>} \dots x^{<n>})}{e.g. P(English | French)}$$

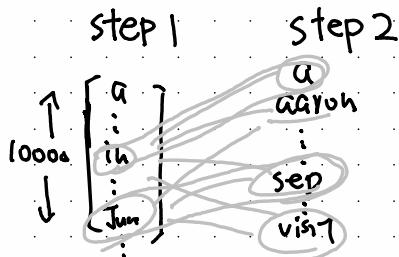
$\underset{y^{<1>} \dots y^{<n>}}{\operatorname{argmax}} P(y^{<1>} \dots y^{<n>} | x)$
French sentence

(
beam search

△ greedy search usually result in worse sentence

e.g. picking common word like "going", "the", ... etc

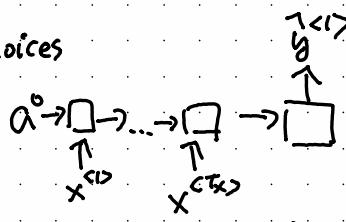
Beam Search



$B=3$ beam width \curvearrowright keep track of top B choices

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

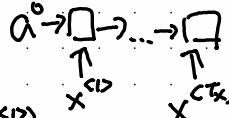
Memorize 3 good choices



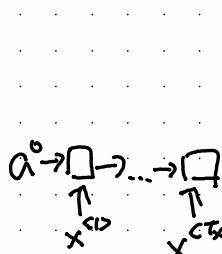
Step 2

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

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



$P(y^{<2} | x, 'in')$



jane

$P(y^{<2} | x, 'jane')$

beam width

eventually pick $B=3$ most likely output
then moves on to next step.

step 3

in september



keep B NNs.

jane is



jane visit



$B=1$ ~ beam search
"greedy" search.

Refinement: Length normalization

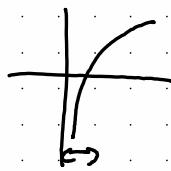
$$\underset{y}{\operatorname{argmax}} \prod_{t=1}^{T_y} P(y^{ct} | x, y^{<1>} \dots y^{<t-1>})$$

- multiplication of P mitigate value
→ tend to pick short sentence

$$\Rightarrow \underset{y}{\operatorname{argmax}} \sum_{y=1}^{T_y} \log P(y^{ct} | x, y^{<1>} \dots y^{<t-1>})$$

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

$$\alpha = 0.7$$



/ large $B \rightarrow$ slow, better
small $B \rightarrow$ fast, faster.

Production $B \sim 10$

research $= B \sim 100 \sim 1000$

to squeeze performance

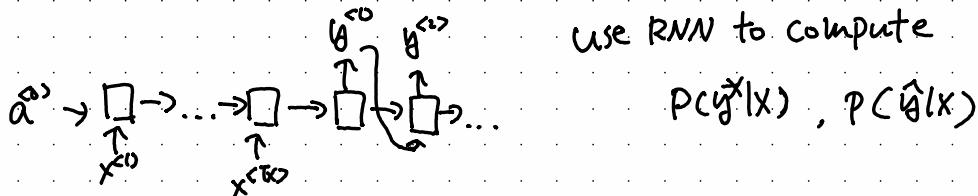
△ Compared to BFS, DFS, beam search doesn't guarantee best solution
- heuristic search

Error analysis in Beam Search.

Jane visite l'Afrique en septembre

Human: Jane visit Africa in September (y^*)

Machine: Jane visit Africa last September (\hat{y})



use RNN to compute

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

Case 1.

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

beam search can improve

Case 2.

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

RNN model is at fault.

→ y^* is better, but RNN says \hat{y} is better, thus RNN is at fault.



do the above analysis for all dev set,

then decide what to do next. - increase B

or

update RNN

Bleu score: bilingual evaluation underway, 2002

- deal with equally good translations

Bleu: Unigram

Sentence: Le chat est sur le tapis

ref 1 : the cat is on the mat

ref 2 : There is a cat on the mat.

:

clipping # occurrence

MT output: the the the ...

Count_{clip}('the')

precision: $\frac{7}{7}$

modified: $\frac{2}{7}$

Count('the')

Bleu score on bigrams

MT: The cat the cat on the mat.

	Count	Count _{clip}
the cat	2	
cat the	1	
the cat	1	
cat on	1	
on the	1	
the mat	1	

$$P_i = \frac{\sum_{\text{eg}} \text{Count}_{\text{clip}}(\text{n-gram})}{\sum_{\text{eg}} \text{Count}(\text{n-gram})}$$

n-gram

$$P_n = \frac{\sum_{\substack{\text{n-gram} \\ \text{eg}}} \text{Count}_{\text{clip}}(\text{n-gram})}{\sum_{\text{n-gram} \in \text{eg}} \text{Count}(\text{n-gram})}$$

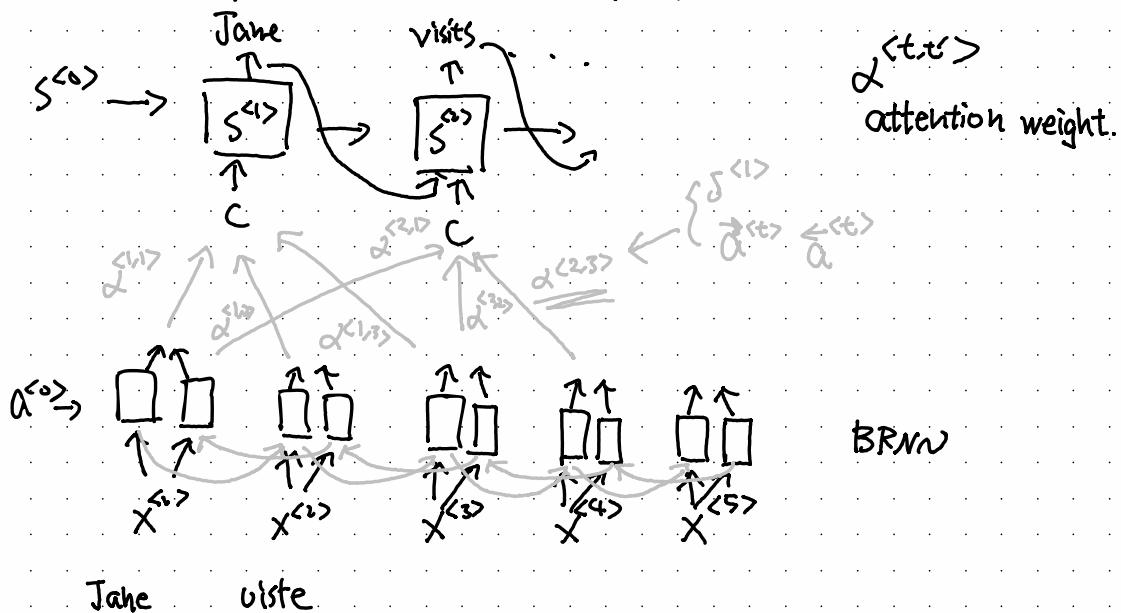
$$\text{Combined: } P = \exp\left(\frac{1}{n} \sum_{i=1}^n P_i\right)$$

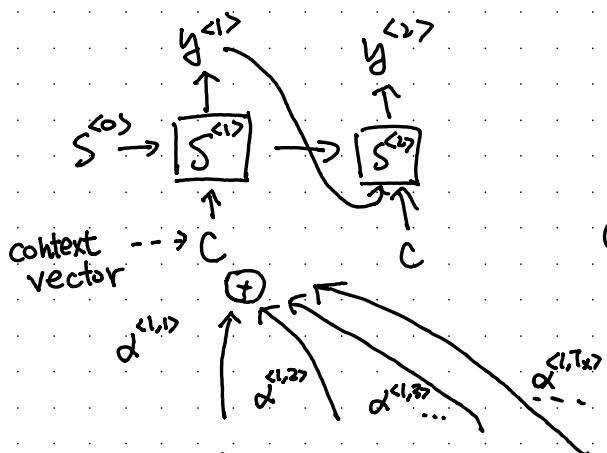
BP, brief penalty \rightarrow penalize short trans

$$B_P = \begin{cases} 1 & \text{if MT output longer than human} \\ e^{(1 - \frac{\text{human_len}}{\text{MT_len}})} & \end{cases}$$

Attention Model, 2014

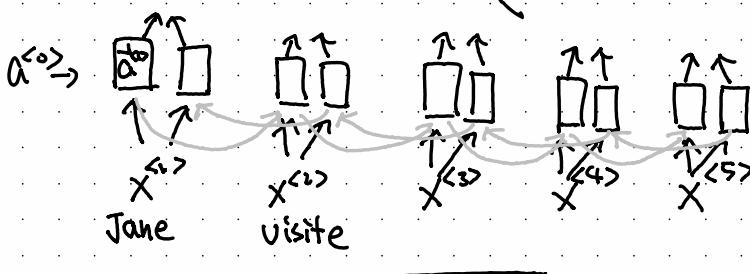
- cannot handle long sequences





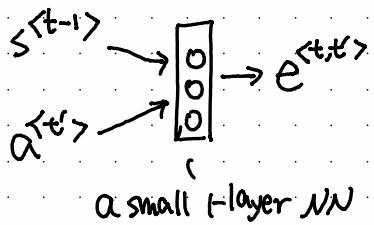
$$\begin{aligned} a^{<t>} &= (\vec{a}^{<t>}, \vec{a}^{<t>}) \\ \sum_t a^{<t>} &= 1 \\ C^{<1>} &= \sum_{t'} \alpha^{<1,t>} a^{<t>} \end{aligned}$$

$\alpha^{<t,t>}$ is amount of attention
 $y^{<t>}$ pay to $a^{<t>}$



$$\alpha^{<t,t>} = \frac{e^{a^{<t,t>}}}{\sum_{t'=1}^{T_x} e^{a^{<t,t'>}}} \quad \leftarrow \text{similar to softmax}$$

runtime



$T_x \cdot T_y$ parameters
 $a^{<t,t'>}$
usually T_x, T_y are not huge.

2015, image captioning, \rightarrow pay attention to a part of the picture

speech recognition

$X \rightarrow Y$
audio transcript

data: 300h ~ 3000h

Commercial: 100,000 h

- Attention model works well

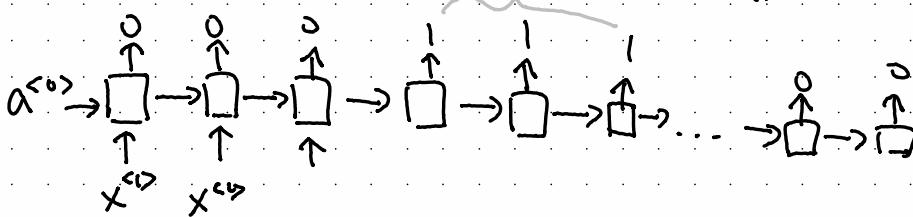
- CTC cost [2006]

Connectionist temporal classification tt...h_eee...w_...ggg...
 $\Delta T_x \gg T_y$, e.g. 1000 Hz in audio \rightarrow the quick brown fox

\rightarrow collapse repeated chars not separated by "blank"

trigger word detection

multiple / for balanced
data set.



Transformer. 2017

RNN → GRU → LSTM

coh:

increase complexity
sequential

transformer:

- attention + CNN

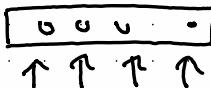
- self-attention. $A^{<1>} A^{<2>} \dots$

- multi-head attention. multiple copies

words

↑

parallel.



self-attention

transformer attention

$A(q, k, v)$ - attention-based vector representation

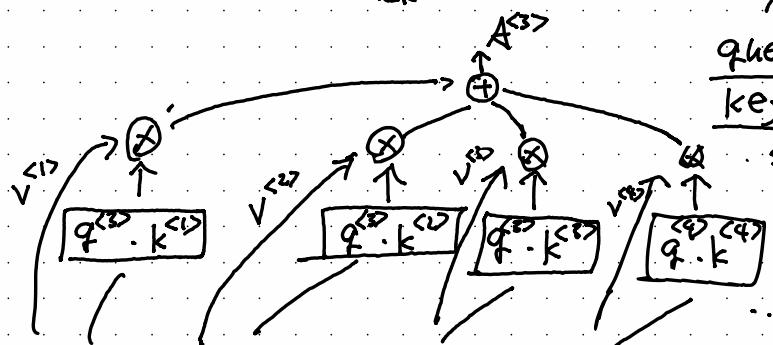
$$= \sum_i \frac{e^{q \cdot k^{<i>}}}{\sum_j e^{q \cdot k^{<j>}}} v^{<i>}$$

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

$$q^{<3>} = W^Q x^{<3>}$$

$$k^{<3>} = W^K x^{<3>}$$

$$v^{<3>} = W^V x^{<3>}$$



query - ask question

key - used to calculate

similarity between word and query

value - for plugging into the representation.

$q^{<1>} k^{<1>} v^{<1>} \quad q^{<2>} k^{<2>} v^{<2>} \quad q^{<3>} k^{<3>} v^{<3>} \quad q^{<4>} k^{<4>} v^{<4>}$
 $x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>}$
 Jane visite l'Afrique en Septembre

↓
more flexible than fixed word embedding

Multi-head Attention

for example

#heads = 8 and? head 1 w_i^Q, w_i^K, w_i^V - what's happening

head 2 w_2^Q, w_2^K, w_2^V - when?

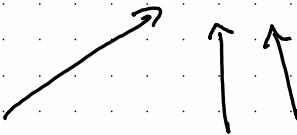
; ; - how?

MultiHead(Q, K, V)



= Concat(head₁, head₂, head₃, ...)

Attention($w_i^Q Q, w_i^K K, w_i^V V$)



w_i^Q, w_i^K, w_i^V

$q^{<1>}, k^{<1>}, v^{<1>}$

$q^{<2>}, k^{<2>}, v^{<2>}$

$x^{<1>}$

Jane

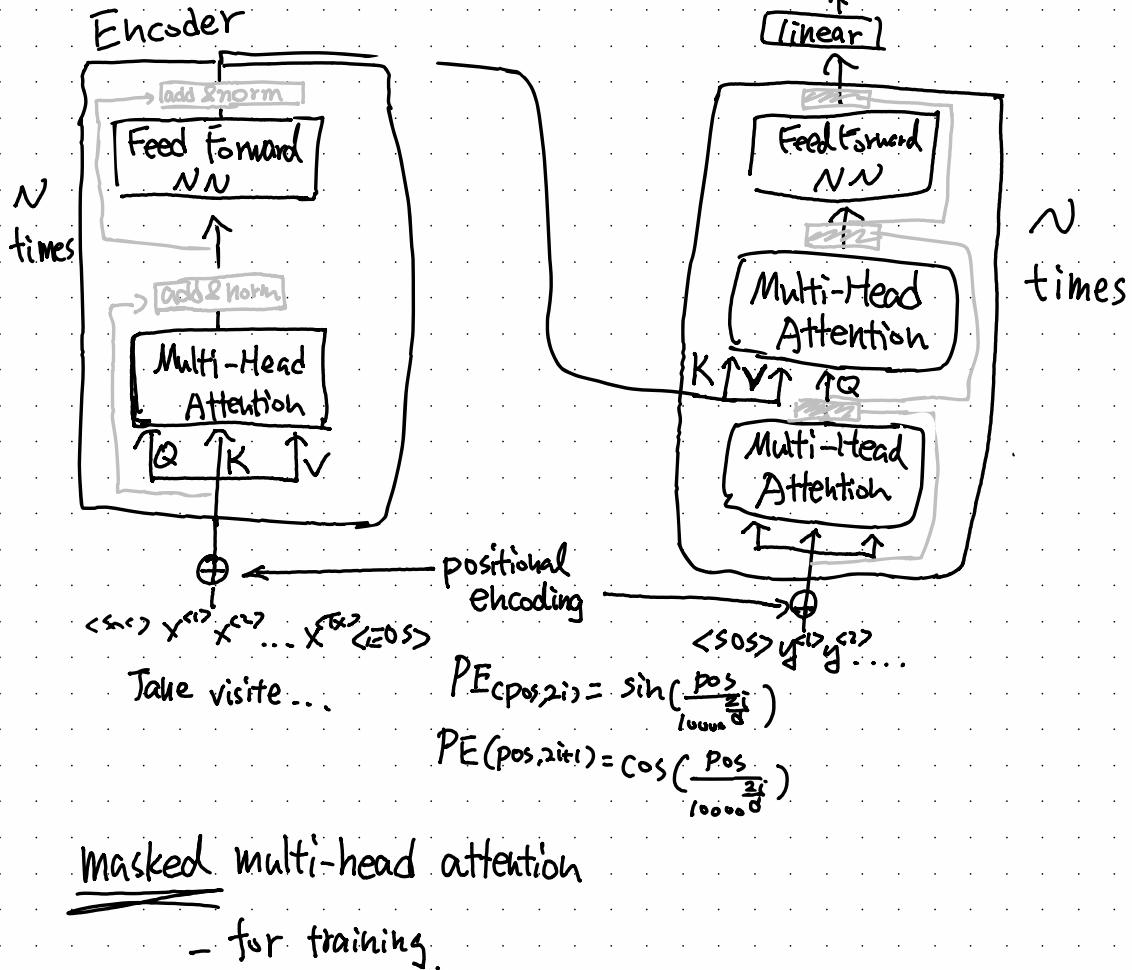
$x^{<2>}$

$q^{<3>}, k^{<3>}, v^{<3>}$

$x^{<3>}$

visite

Transformer Network



masked multi-head attention

- for training.