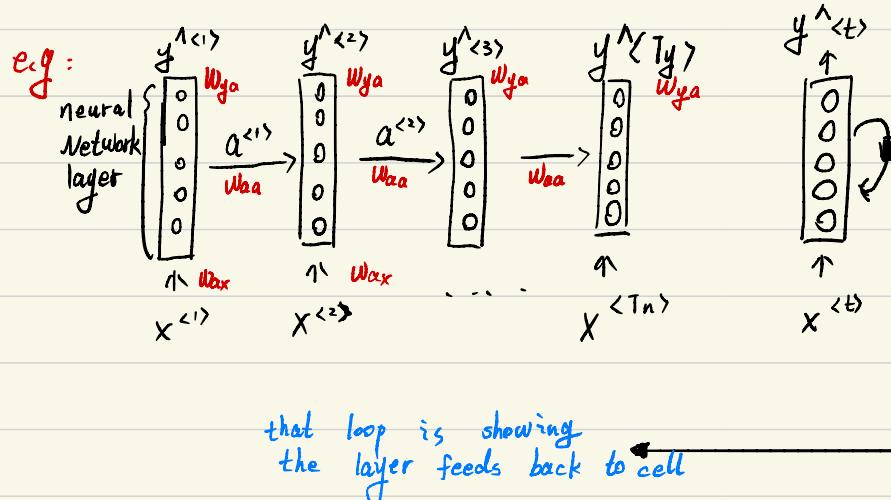


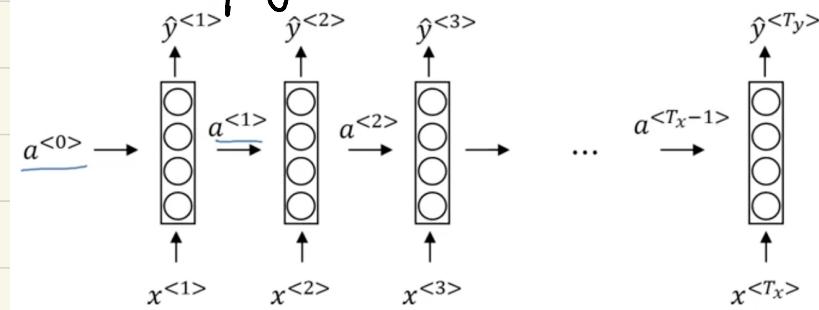
Recurrent Neural Network Model



Limitation for RNN:

the prediction for inputs uses from
input earlier not input later in sequence

• Forward Propagation



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

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

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

\leftarrow constant in linear regression

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

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

\downarrow simplified RNN Notation

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

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

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

$$\begin{bmatrix} W_{ba} & | & W_{ax} \end{bmatrix} = W_a$$

$$\begin{smallmatrix} \uparrow 100 \\ \downarrow 100 \rightarrow \\ \leftarrow 100 \rightarrow \end{smallmatrix} \quad \begin{smallmatrix} \uparrow 100 \\ \downarrow 100 \rightarrow \\ \leftarrow 100 \rightarrow \end{smallmatrix} \quad (100, 1100)$$

Backpropagation through time

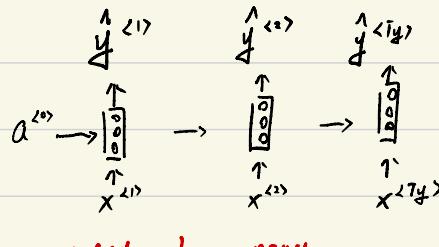
Loss function

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

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

reverse of Forward propagation

Different type of RNN

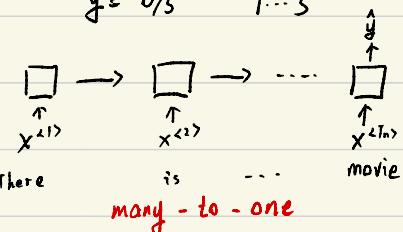


many - to - many

e.g.: Sentiment classification

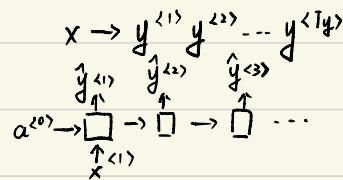
$x = \text{text}$

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



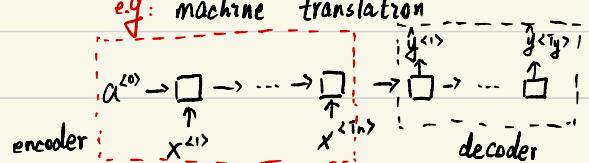
There is ... many - to - one

e.g.: music generation

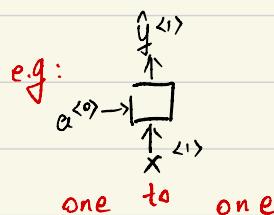


one - to - many

e.g.: machine translation



many - to - many



one to one

Language modeling with RNN

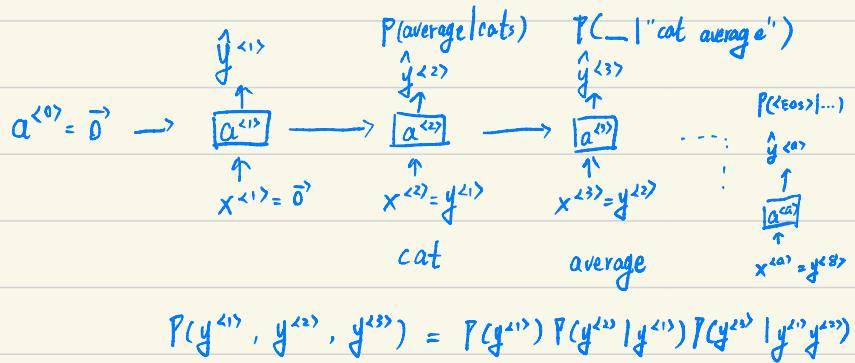
- training set: large corpus of english text

Tokenize (标记化)

e.g.: - Cat average 15 hours of sleep a day <eos>
 $y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad \dots \quad y^{<20>} \quad y^{<21>}$

The Egyptian Mau is a breed of cat <eos>
<unk>

→ RNN model



Sampling novel sequences

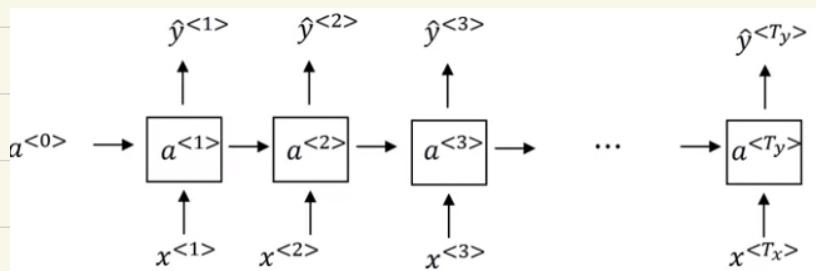
character-level language model
no worry with <UNK>; too long

Sequence generation

Vanishing gradients with RNN (梯度消失)

The cat, which ate ..., was full

The cats, which ate ..., were full



exploding gradients \rightarrow sort address by gradient clipping

Gated Recurrent Unit (GRU) solve vanishing gradient with RNN

$C = \text{memory cell}$

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

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

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

\swarrow update

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

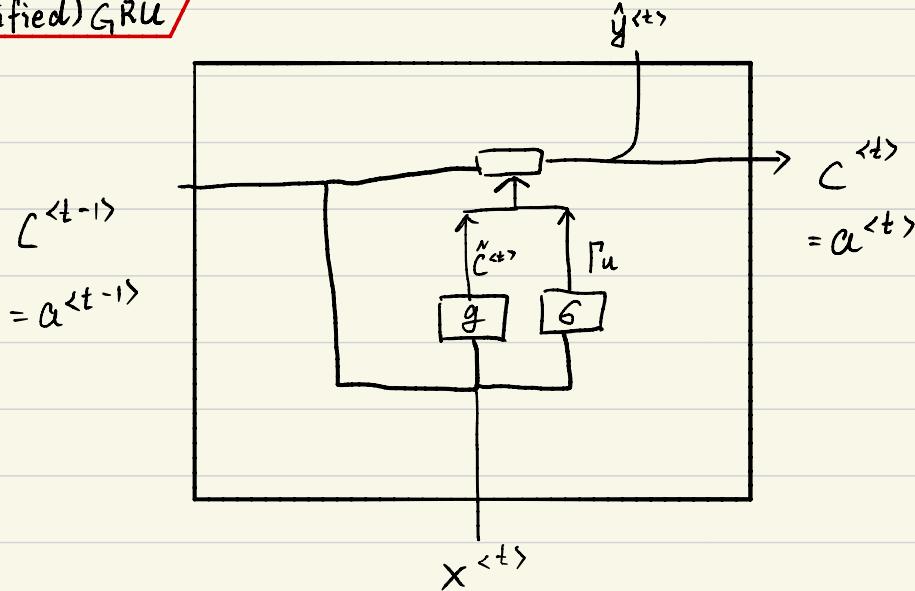
e.g.:

The cat, which ... , was full

$$\rightarrow C^{<t>} = 1 \quad \dots \quad = 1$$

$$\Gamma_u = 1 \quad (\text{determine if update to } 1)$$

(simplified) GRU



full GRU

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

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

LSTM long short Term Memory

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

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

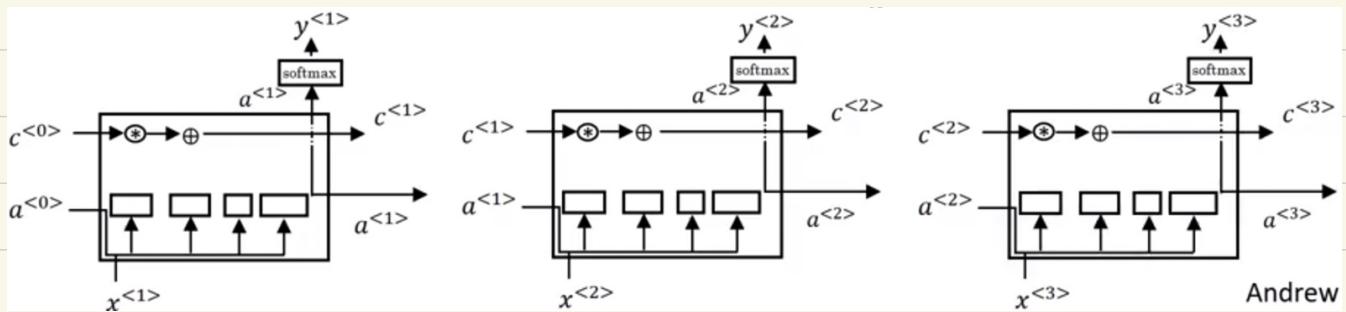
forget $\rightarrow \Gamma_f = G(W_f[a^{<t-1>}, x^{<t>}] + b_f)$

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

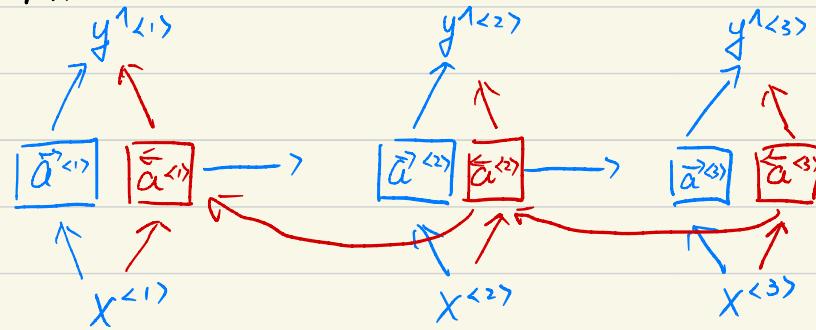
↑ output

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

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

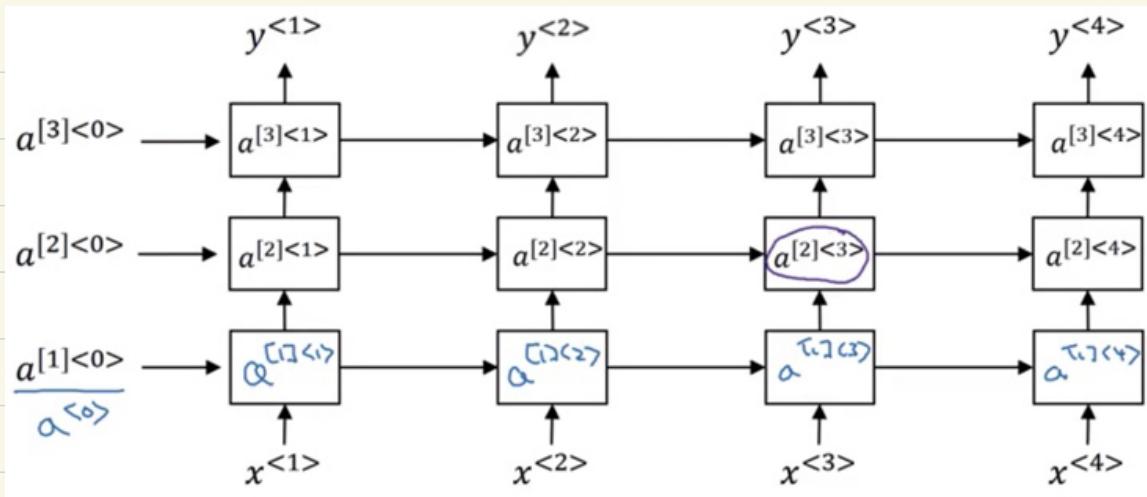


Bidirectional RNN



allow $a^{<t>}$ to take front and back information
 need all data (when speech, need complete and then can translate)

Deep RNN



$$a^{[2]}<3> = g(W_a^{[2]} [a^{[2]}<2>, a^{[1]}<3>] + b_a^{[2]})$$

Word embedding 词向量

1-hot representation $|V| = 10\,000$

Man	Woman	king	queen	apple	Orange
(5391)	(9853)	(4914)	(7151)	(456)	(6257)
$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}$
5391 \rightarrow	9853 \rightarrow				

Featurized representation

gender	$\begin{bmatrix} -1 \\ 0.01 \\ 0.03 \\ 0.04 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 0.02 \\ 0.02 \\ 0.01 \end{bmatrix}$	-0.95	0.97	0.00	0.01
Royal	$\begin{bmatrix} 0.01 \end{bmatrix}$	$\begin{bmatrix} 0.02 \end{bmatrix}$	0.93	0.95	-0.01	0.00
Age	$\begin{bmatrix} 0.03 \end{bmatrix}$	$\begin{bmatrix} 0.02 \end{bmatrix}$	0.7	0.69	0.03	-0.02
Food	$\begin{bmatrix} 0.04 \end{bmatrix}$	$\begin{bmatrix} 0.01 \end{bmatrix}$	0.02	0.01	0.95	0.97

e_{5391} e_{9853}

e.g.: I want a glass of orange — apple —

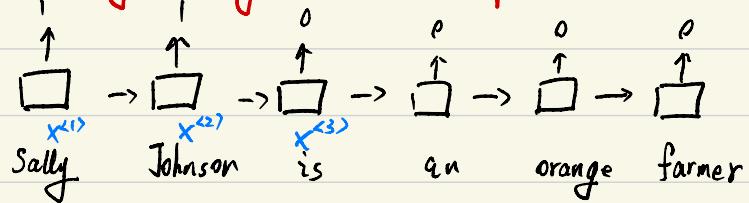
most features are similar

Visualizing word embeddings

t-SNE

Using word embeddings

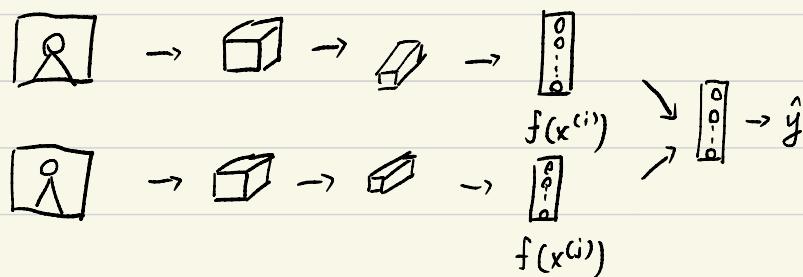
e.g: Named entity recognition example



transfer learning and word embedding

1. Learn word embedding from large text corpus (1-100B words)
2. transfer embedding to new task with smaller training set
3. Optional: Continue to finetune the word embeddings with new data.

relation to face encoding (embedding)



Properties of word embeddings

analogies

	Man	Women	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
Food	0.04	0.01	0.02	0.01	0.95	0.97

C₅₃₉₁

E₉₈₅₃

How to get

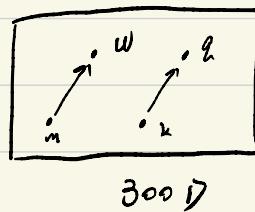
$\text{man} \rightarrow \text{woman} \Rightarrow \text{king} \rightarrow ?$

$$\text{eman} - \text{ewoman} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{eking} - \text{equeen} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

So $\text{eman} - \text{ewoman} \approx \text{eking} - \text{equeen}$

analogies using word vectors



$$\text{eman} - \text{ewoman} \approx \text{eking} - \text{eq}$$

$$\arg \max_w \text{sim}(\text{ew}, \text{eking} - \text{eman} + \text{ewoman})$$

Cosine similarity

$$\text{sim}(\text{ew}, \text{eking} - \text{man} + \text{ewoman})$$

$$\rightarrow \text{sim}(U, V) = \frac{U^T V}{\|U\|_2 \|V\|_2}$$



e.g.: Ottawa : Canada as Nairobi : Kenya

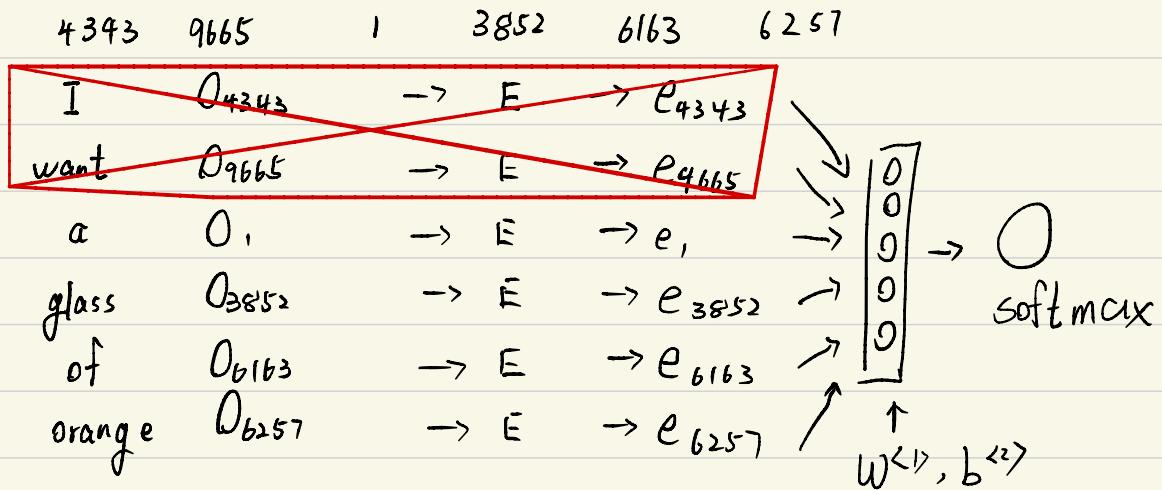
Big : Bigger as Tall , Taller

Embedding matrix

a	aaron	---	orange	6757	---	zulu	<UNK>
300					10000		

Learning word embeddings

e.g.: I want a glass of orange —.



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

context: Last 4 words

4 words or left 8 right

a glass of orange ? to go along with

Last 1 word orange ?

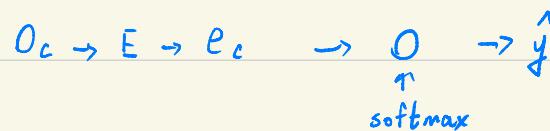
Nearby 1 word (skip gram) glass ?

Word2Vec

skip gram Model

Vocab size = 10 000 k

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



softmax: $P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$ (need long computation time)

↑ solve that

Hierarchical softmax



Negative Sampling

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

$k = 5 - 20$ smaller dataset

$k = 2 - 5$ larger dataset

model

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

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

selecting negative example

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10000} f(w_j)^{3/4}}$$

GloVe word vectors (global vectors for word representation)

e.g.: I want a glass of orange juice to go along

with my cereal

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

$$\text{minimize} \quad \sum_{i=1}^{10000} \sum_{j=1}^{10000} f(X_{ij}) \underbrace{(\theta_i^T e_j + b_i + b_j - \log X_{ij})^2}_{\text{weighting term}}$$

sentiment classification

$x \rightarrow y$
The dessert is excellent
Service was quite slow

challenge: x have a huge label data set

Simple sentiment classification model

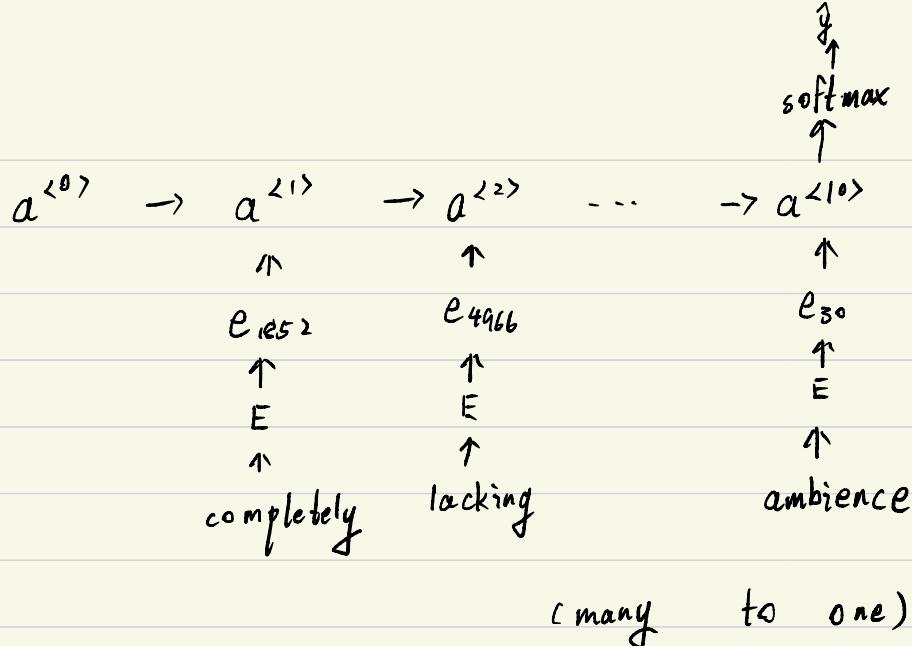
The dessert is excellent
8928 2468 4694 3148

The 08928 → E → e_{8928}
dessert 02468 → E → e_{2468} → Avg → O softmax → \hat{y}
is 04694 → E → e_{4694}
excellent 03148 → E → e_{3148} / 300D

But not good at this example

"Completely lacking in good taste, good service"

RNN for sentiment classification



Debiasing word embedding

the problem of bias in word embedding

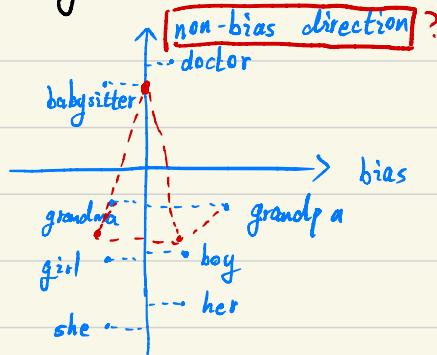
Man:Woman as King:Queen

Man:Computer_Programmer as Woman:Homemaker X

Father:Doctor as Mother:Nurse X

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

addressing bias in word embeddings



1. identify bias direction
 - {
 - Che - Eshe
 - Cmale - Cfemale
 2. neutralize : For every word that is not definitional, project to get rid of bias
 3. equalize pairs