

Deeplearning course VII week 1

Recurrent neural networks

Supervised learning

Why sequence models

Speech recognition

x → "The quick brown fox jumped over the lazy dog"

Music generation

x cause y → y

Sentiment classification

e.g. single intefr → there is nothing to like in this movie → $\star\star\star\star\star$

DNA sequence analysis

ACGGCCATTCG → ...

Machine translation

French → English

Video activity recognition

movieclip → running

Name entity recognition

sentence → identify names

RNNs have a kind of memory = C

BRNN can work with info from past and future.

Notation example

NLP - Natural Language Processing

$x: (\text{Harry Potter})$

and $(\text{Hermione Granger})$ invented a new spell

$\rightarrow y:$

$x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad \dots \quad x^{(q)}$

$T_x^{(i)}$

$T_y^{(i)}$

input sequence length

$x^{(q+1)} \quad \dots \quad x^{(T_x)}$

$T_x = q+1$
cause different
 $T_y = q$

Representing words

Vocabulary:
(dictionary)

a	1
aaron	2
:	:
and	367
:	:
harry	4075
:	:
potter	6830
:	:
zalm	10000 words
Lunk	unknown word
quit small	one-hot
common 50k-100k	(One 1 and zeros everywhere)
or millions	

$n_x \in \mathbb{R}^{10000}$
shape $(10000,)$

input tensor
shape (n_x, m, T_x)

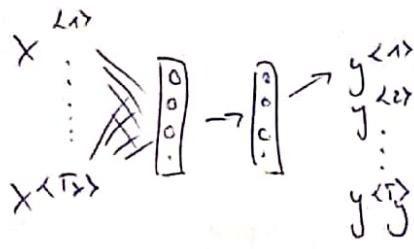
(i) mini-batch size (n_a, m)

$q^{(t)}$ is called

"hidden state"



Why not a standard network?

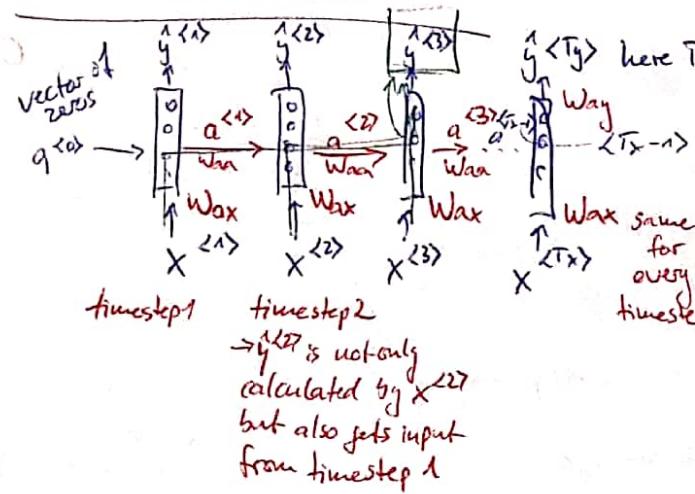


Problems: • Input, output can be different in lengths in different examples

- Doesn't share features learned across different positions of text
 - ↳ as in CNN: things learned in part of the image to generalize quickly to other parts of the image
- weight matrix will be enormous

(unidirectional)

Recurrent Neural Network



Weakness: only uses information that is earlier in the sequence to make a prediction

"He said, Teddy Roosevelt was not president"
"he said, Teddy bears are on sale"
↳ can't tell whether "Teddy" is a part of a person's name, in 1st sentence y<2>, in 2nd no.

Forward Propagation

activation (tanh/ReLU)

$$\begin{aligned} a^{(0)} &= 0 \\ a^{(1)} &= g(W_{aa} \cdot a^{(0)} + W_{ax} \cdot x^{(1)} + b_a) \\ y^{(1)} &= g(W_{ya} \cdot a^{(1)} + b_y) \\ &\text{sigmoid for binary classification} \\ a^{(2)} &= g(W_{aa} \cdot a^{(1)} + W_{ax} \cdot x^{(2)} + b_a) \\ y^{(2)} &= g(W_{ya} \cdot a^{(2)} + b_y) \end{aligned}$$

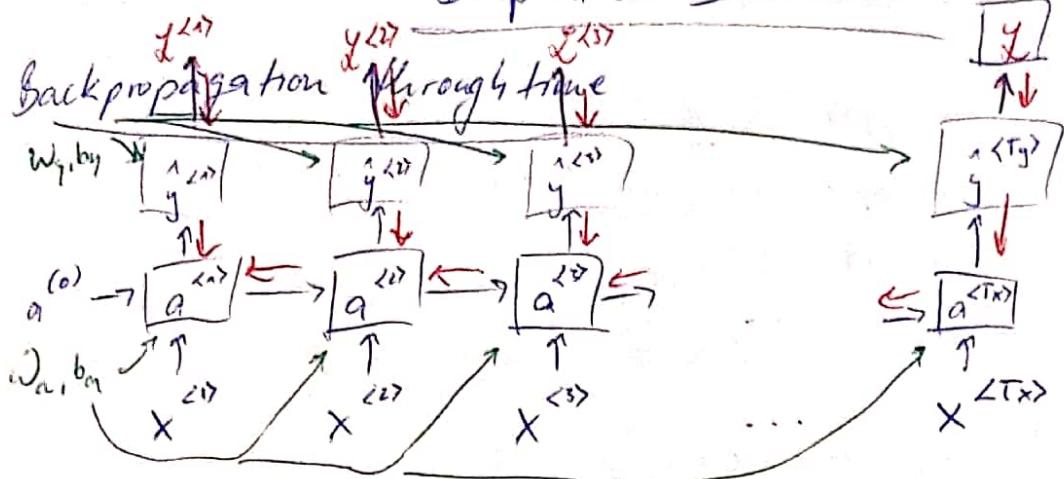
→ later: Bidirectional RNN

$a \leftarrow W_{ax} x^{(1)}$
↳ multiplied by some x-like quantity to compute some a-like quantity

Simplified RNN notation:

$$\begin{aligned} a^{(t)} &= g(W_{aa} \cdot a^{(t-1)} + W_{ax} \cdot x^{(t)} + b_a) \\ y^{(t)} &= g(W_{ya} \cdot a^{(t)} + b_y) \end{aligned}$$

$$\begin{aligned} a^{(t)} &= g(W_a [a^{(t-1)}, x^{(t)}] + b_a) \\ &\quad \text{[} W_a [W_a | W_{ax}] = W_a \text{ (stacked horizontally)]} \\ &\quad \text{[} (W_a; W_{ax}) \cdot \begin{pmatrix} a^{(t-1)} \\ x^{(t)} \end{pmatrix} = W_a \cdot a^{(t-1)} + W_{ax} \cdot x^{(t)} \end{aligned}$$



Backprop
through
time

Loss function

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

(loss of single prediction in
single time step)

overall loss for
entire sequence :
$$L(\hat{y}, y) = \sum_{t=1}^{Ty} L^{(t)}(y^{(t)}, \hat{y}^{(t)})$$

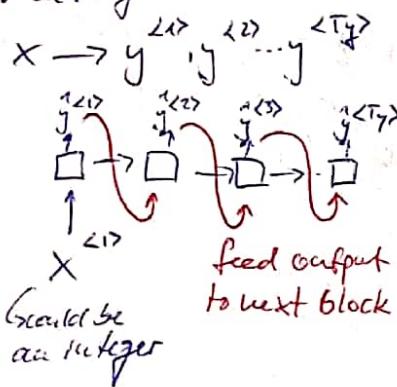
Different types of RNNs

Often $T_x \neq T_y$! or even different classes
(in & output lengths)

so far "many-to-many"

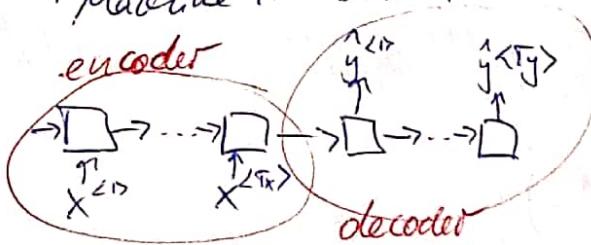
"One-to-many"

Music generation



other type of
"many-to-many"
with different T_x, T_y -lengths

→ Machine translation



Sequence classification

$X = \text{text}$ $y = 0/1 \dots 1 \dots 5$ suitable for
binary class.

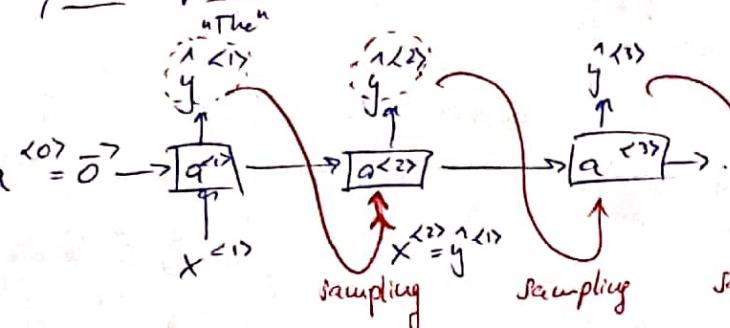
\boxed{y} read the entire
sentence and
 \dots generate just
one output

There is ... more

"many-to-one"

Sampling a sequence from a trained RNN

Sampling:



→ np.random.choice → randomly chosen sentence

Vocabulary = {a, arrow, ..., z, }

or = {a, b, c, ..., z, , 0, ..., 9, A, ..., Z}

Characters

cat
T T T
y z o y z

- pro: don't have to worry about unknown tokens
- con: character-based language model is able to give "man" a non-zero probability!
- con: much more and longer sequences
- con: difficult to capture long-range dependencies

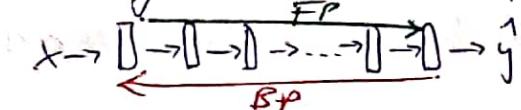
Vanishing gradients with RNNs

(bigger problem than exploding gradients)

long-range dependencies "The cat, which ate a lot..., was full"
"The cats, ... , were full"

↳ distorts predictability

Vanishing gradient problem arises for very deep NNs as y^1 has a hard time to backpropagate to layer $\alpha^{[1]}$ (activation layer 1)



updates may become so large that they overshoot the optimal value during back propagation

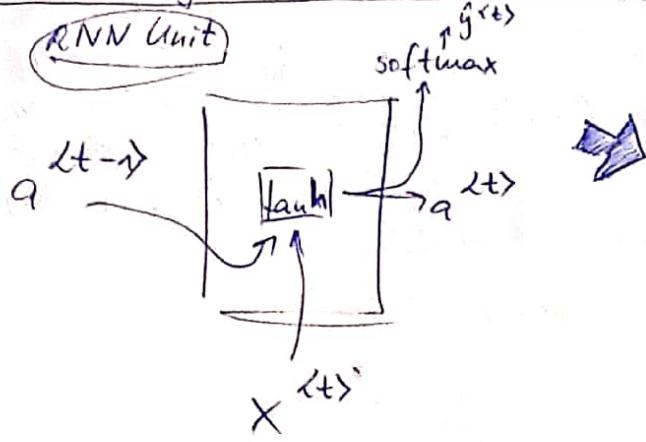
gradient clipping avoids exploding gradients

Gated Recurrent Unit (GRU) → helps with vanishing gradients

[Cho, 2014 On the properties of neural translation: Encoder-decoder...]

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

RNN Unit



GRU (simplified)

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

$$\text{overwriting every timestep} \quad c^{(t)} = \tanh(W_c [c^{(t-1)}, x^{(t)}] + b_c)$$

will decay if updated or not!
is it 0 or 1 update

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

The cat, which already ate ... was full

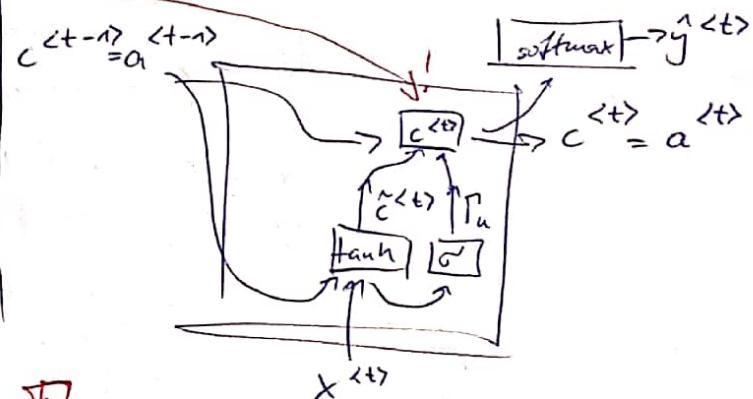
Full GRU Γ_u relevant

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

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

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

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



Γ_u can be very close to zero (0.000001) that it is sufficiently from vanishing problem!

Enables very long-range dependencies!

↳ For signal to BP without vanishing we need $c^{(t)}$ to be highly dependent on $c^{(t-1)}$

(VI)

Long Short Term Memory (LSTM) \Rightarrow for long distance learning, more powerful than GRU

GRU

$$\begin{aligned} \text{candidate for update } & \tilde{c}^{(t)} = \tanh(W_c [f_r * c^{(t-1)}, x^{(t)}] + b_c) \\ \text{update gate } & f_u = \sigma(W_u [c^{(t-1)}, x^{(t)}] + b_u) \\ \text{relevance gate } & f_r = \sigma(W_r [c^{(t-1)}, x^{(t)}] + b_r) \\ \cancel{\text{candidate}} \quad & c^{(t)} = \underbrace{f_u * \tilde{c}^{(t)}}_{\tilde{c}^{(t)}} + (1 - f_u) * c^{(t-1)} \\ \cancel{\text{candidate}} \quad & a^{(t)} = c^{(t)} \end{aligned}$$

[Hochreiter & Schmidhuber 1997, LSTM]

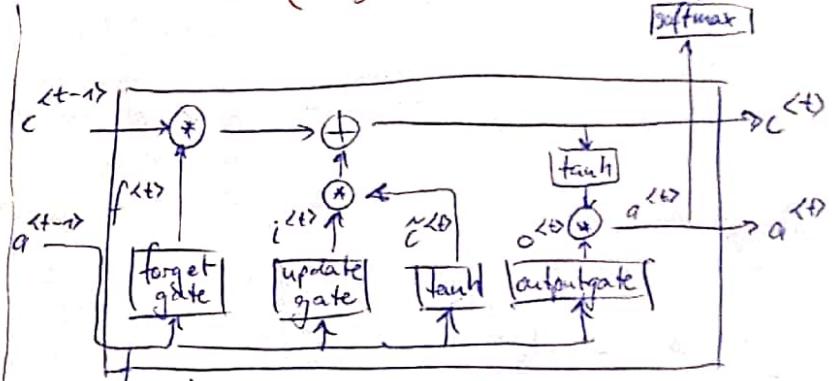
GRU came up later actually and is a simplification

LSTM

$$\begin{aligned} \tilde{c}^{(t)} &= \tanh(W_c [a^{(t-1)}, x^{(t)}] + b_c) \\ f_u &= \sigma(W_u [a^{(t-1)}, x^{(t)}] + b_u) \text{ instead of } f_u = \sigma(W_u [c^{(t-1)}, x^{(t)}] + b_u) \text{ having one update gate} \\ f_f &= \sigma(W_f [a^{(t-1)}, x^{(t)}] + b_f) \text{ controlling both terms } \\ f_o &= \sigma(W_o [a^{(t-1)}, x^{(t)}] + b_o) \text{ 2 separate terms} \\ c^{(t)} &= f_u * \tilde{c}^{(t)} + f_f * c^{(t-1)} \\ \rightarrow & \text{gives } c \text{ the option to keep the old value } c^{(t-1)} \text{ and then add it to new } \tilde{c}^{(t)} \\ \boxed{\text{↳ separate update and forget gates!}} \end{aligned}$$

$$a^{(t)} = f_o \tanh(c^{(t)})$$

softmax



key take aways $\rightarrow [a^{(t-1)}, x^{(t)}]$ to compute all gate values
 \rightarrow get $c^{(t)}$ from $c^{(t-1)}$

prechore connection \rightarrow adding earlier cell memory

$$[a^{(t-1)}, x^{(t)}, c^{(t-1)}]$$

\Rightarrow more flexible and powerful

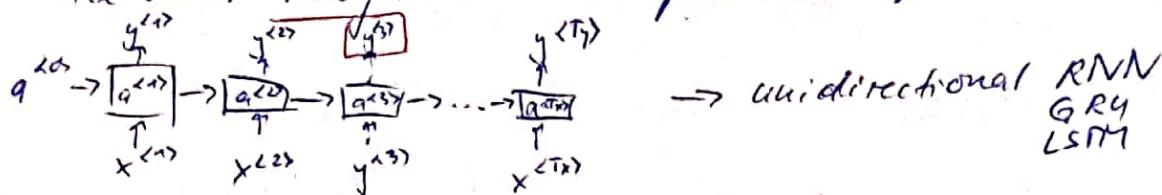
Bidirectional RNN

→ commonly BRNN with LSTM for NLP

Getting information from the future

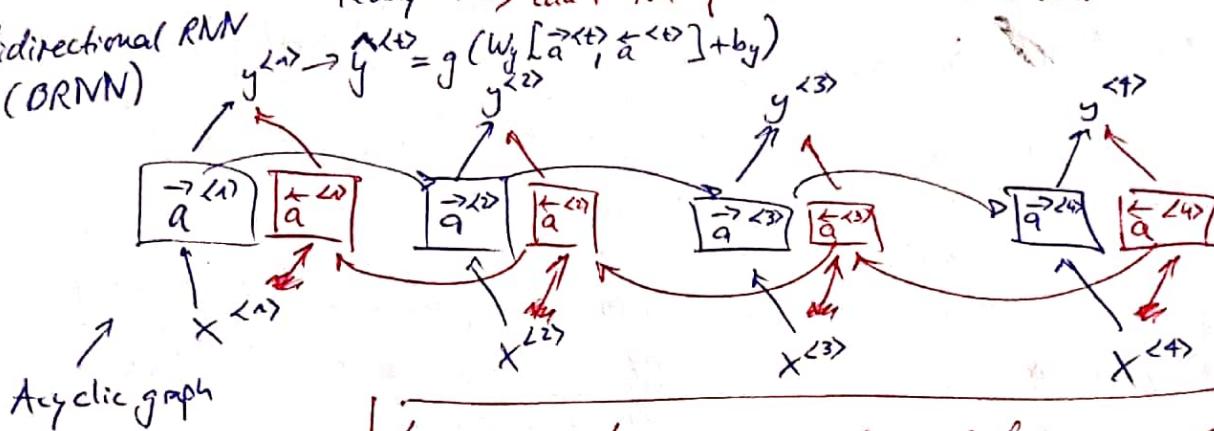
He said, "Teddy bears are on sale!"

He said, "Teddy Roosevelt was a great President"



"Teddy" → can't tell if it's the name or a bear

Bidirectional RNN (BRNN)

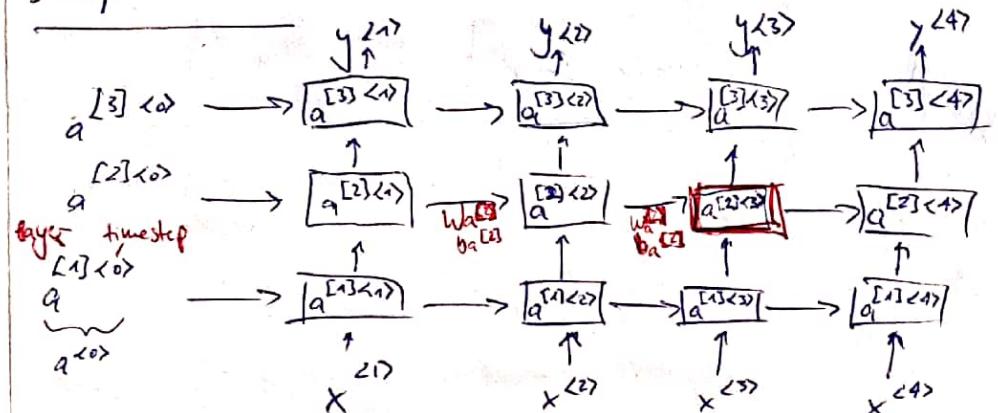


↳ allows to include past and future parts into predictions in the middle of a sequence

→ contrast: need entire sequence before a prediction can be made

↳ not for real-time NLP!

Deep RNNs



$$\text{i.e. } a^{[2]}_{<3>} = g(W_a^{[2]} [a^{[2]}_{<2>}, a^{[1]}_{<3>}] + b_a^{[2]})$$

for RNNs 3 layers
can be already a lot
But outputs like $y^{[2]}$
can be added with a
from layer 2 deep NN with
not horizontal (or
temporal connection
and then compute
 $y^{[2]}$ $\rightarrow y^{[2]}$
↳ $\rightarrow [a^{[2]}_{<2>}] \rightarrow \text{deep NN}$

NLP and word EmbeddingsWord representation

$V = \{v_1, v_2, \dots, v_{|V|}, \langle \text{UNK} \rangle\}$ $|V|=10000 \rightarrow$ One-hot representation

Man (5391)	Woman (9853)	Orange (6259)
$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

Weakness of this representation:

Treats each word as a thing itself
and enables no cross-relation

"I want a glass of orange juice."

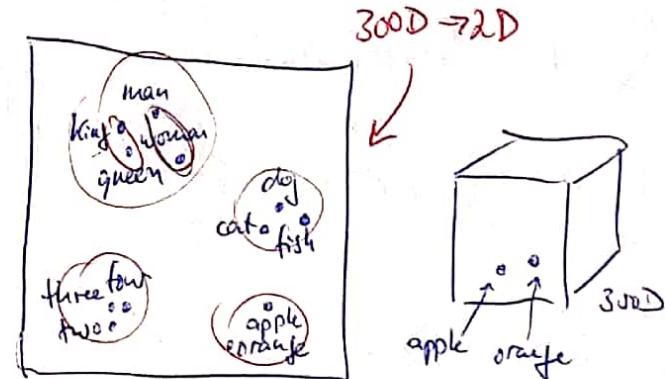
"I want a glass of apple ____?"

algorithm does not

know that apple
juice and orange juice
are both popular things

Factorized representation: word embedding
because the inner product of the one-hot
vectors of "apple" and "orange" is zero!

	Man (5391)	Woman (9853)	King	Queen	Apple	Orange
Gender	-1	1	-0.45	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	0.75	0.97
size						
cost						
action	XER	300				
:						
features						
300	(300, 1)					



t-SNE

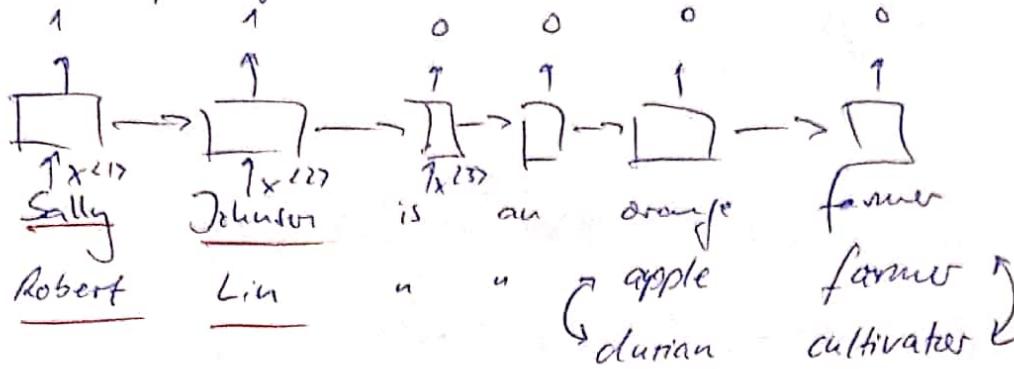
[van der Maaten and Hinton, 2008, Visualizing data with t-SNE]

Embeddings group reasonable words together

One-hot vectors do not do a good job at capturing similarities between words (every 1-hot-vec has the same Euclidean distance)

Using word embeddings

Named entity recognition example



examining lots of unlabeled text (From internet download)

Gallows to carry out transfer learning

→ 1 billion words - 100B words

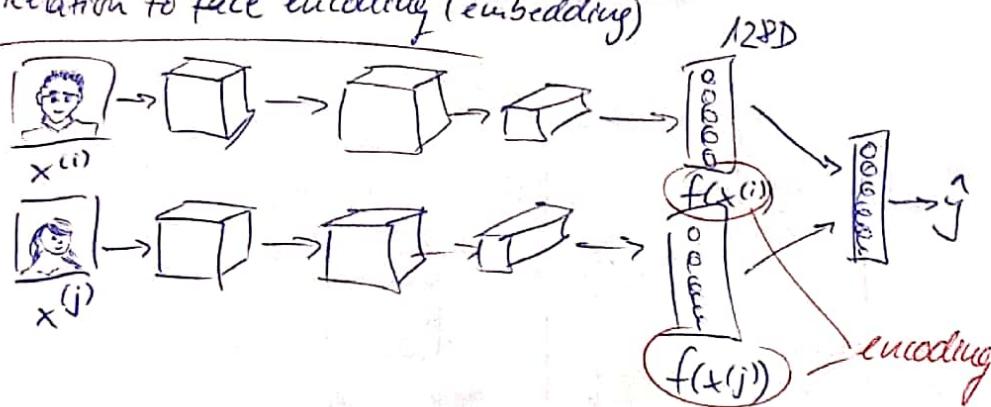
BRNN

→ 100k words → pass to smaller labeled data-set

Transfer learning and word embeddings

- A) 1. Learn word embeddings from large text corpus (1-100B words)
(or download pre-trained embedding online)
- B) 2. Transfer embedding to new task with smaller training set (100k words)
→ pos: can use lower-dim feature vectors 10'000 one-hot → 300-dim vector dense vector
3. Optional: Continue to finetune the word embeddings with new data

Relation to face encoding (embedding)



difference:

In face recognition we want to train a NN to classify an image that wasn't seen before and in NLP we work with a fixed vocabulary

(2)

Properties of word embeddings

[featureized matrix]

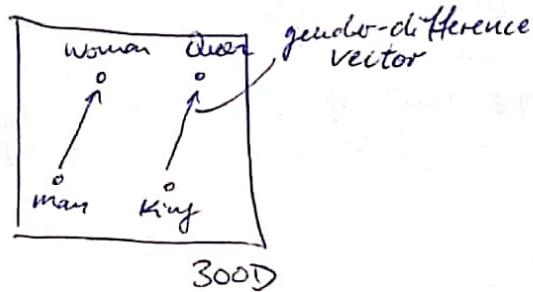
	Man	Woman	King	Queen
gender	$\begin{bmatrix} -1 \\ 0.01 \\ 0.03 \\ 0.09 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 0.02 \\ 0.02 \\ 0.01 \end{bmatrix}$	$\begin{bmatrix} -0.95 \\ 0.93 \\ 0.70 \\ -0.02 \end{bmatrix}$	$\begin{bmatrix} 0.97 \\ 0.95 \\ 0.67 \\ 0.01 \end{bmatrix}$
Royal				
Age				
Food				
lman	lwoman	lman - lwoman	$\approx \begin{bmatrix} -2 \\ 0 \\ 0 \end{bmatrix}$	

[Mikolov, 2013, Linguistic regularities in continuous analogy reasoning space word representation]

$$\text{Man} \rightarrow \text{Woman} \text{ as King} \rightarrow ? \quad \text{Queen} \quad \text{lking} - \text{lqueen} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \end{bmatrix} \quad \left. \begin{array}{l} \text{common} \\ \text{difference} \\ \text{in the} \\ \text{gender} \end{array} \right\}$$

$$\hookrightarrow \text{lman} - \text{lwoman} \approx \text{lking} - \text{lqueen}$$

high similarity

Analogies using word vectors

$$\rightarrow \text{Find word } w_i: \underset{w}{\operatorname{argmax}} \sim(\text{lw}, \text{lking} - \text{lman} + \text{lw})$$

↑ similarity
↓ accuracy reported 30% - 75%

Cosine similarity \rightarrow Measure similarity between 2 words

$$\sim(\text{lw}, \text{lking} - \text{lman} + \text{lw})$$

word embedding algorithm
(can learn)

man:woman as boy:girl
big:biggish as tall:taller

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

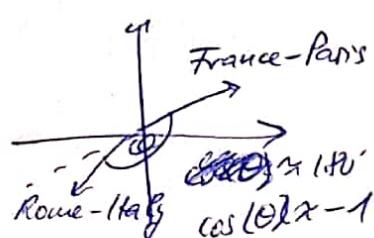
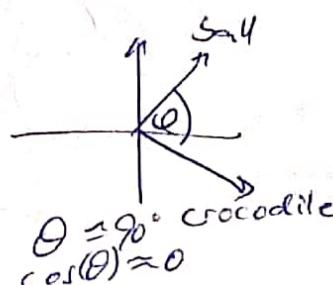
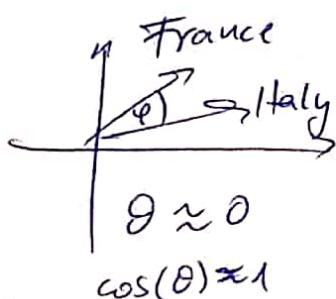
inner product tends to be large when similarity is high

Euclidean lengths

when vectors pointing in same direction

similarity is 1

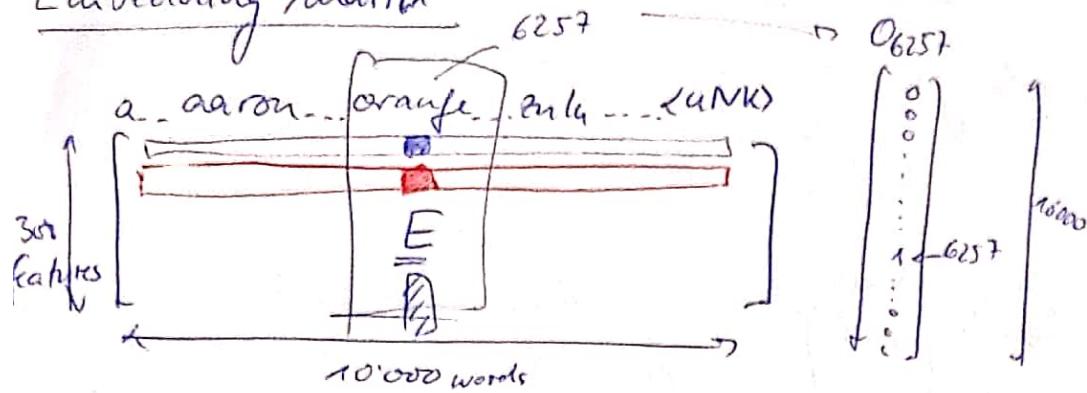
$$\|u - v\|^2 \rightarrow \text{measure of dis-similarity}$$



vectors similar but encode country and city the other way round

(3)

Embedding matrix



$$E \cdot O_{6257} = \begin{bmatrix} \text{blue box} \\ \text{red box} \\ \vdots \\ \text{black box} \end{bmatrix} = e_{6257} \quad e_w = \text{embedding for word } w$$

$(300, 10000) \quad (10000, 1)$

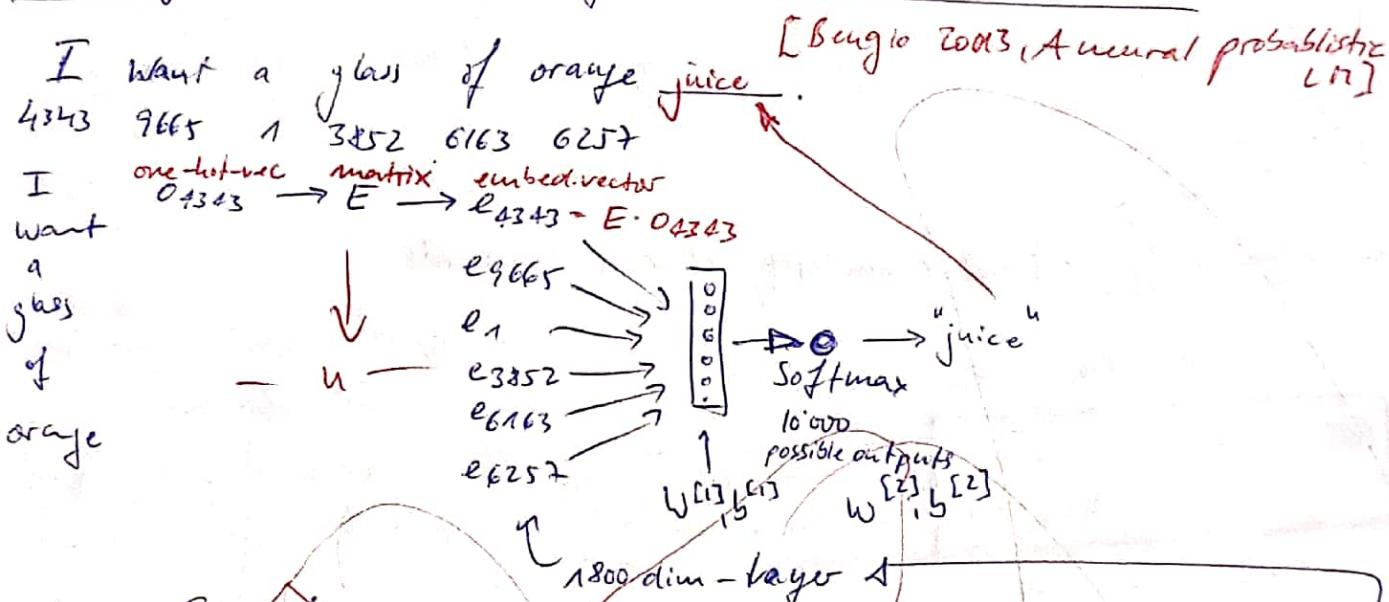
$(300, 1)$

$E \cdot O_j = e_j \quad (\text{Initialize } E \text{ randomly})$

In practice use specialized function to look up an embedding

Keras → Embedding layer calls column of interest

Learning Word Embeddings: Word2vec & GloVe



using 300-dimensional word embeddings \rightarrow 6 words $\rightarrow 6 \times 300 = 1800$ dimension vector

→ in practice look at four previous words to make a prediction: $4 \times 300 \rightarrow 1200$ dimensional vector

↳ can deal with sentences of any size

DL spec // course II // week 2

Other context/target pairs Skipgram model

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

context: last four words.

Woodson 8/11/81 *for lack of escape* —.

- 4 words on left & right

- " — topology with "

* Last 1 word

* Nearly 1 word

Word2vec-algorithm (Mikolov 2013, efficient estimation of word representation in vector space)
Skip-grams

Skip-grams

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

Context

(randomly
chosen)

drauf

Tarjet

judge (1 w. letter)

glass (2 before)

my (...)

Model details

Vocab size = 10'000k

Context c ("orange") \rightarrow Target t ("juice")
6257 4,871

one-hot
vector

$$O_c \rightarrow E \rightarrow O_c = E \cdot O_c \xrightarrow{\text{Softmax}} O \xrightarrow{y^1} y$$

E has a lot of parameters

Softmax:
probability of
t parameterized
by c

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

θ_E = parameter associated with output

$$\mathcal{L}(\hat{y}_i | y_i) = - \sum_{j=1}^{1000} y_i \log \hat{y}_i$$

$$C^y = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad \leftarrow 4834$$

ER^{10,000} \rightarrow
probabilities for all target words

Problems with softmax classification

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

→ has to sum overall words in vocab in each iter.
 ↳ solution: Hierarchical softmax classifier



→ scales in $\log N$

→ common words rather on top and vice versa

How to sample context c ?

→ the, of, a, and, to, ...

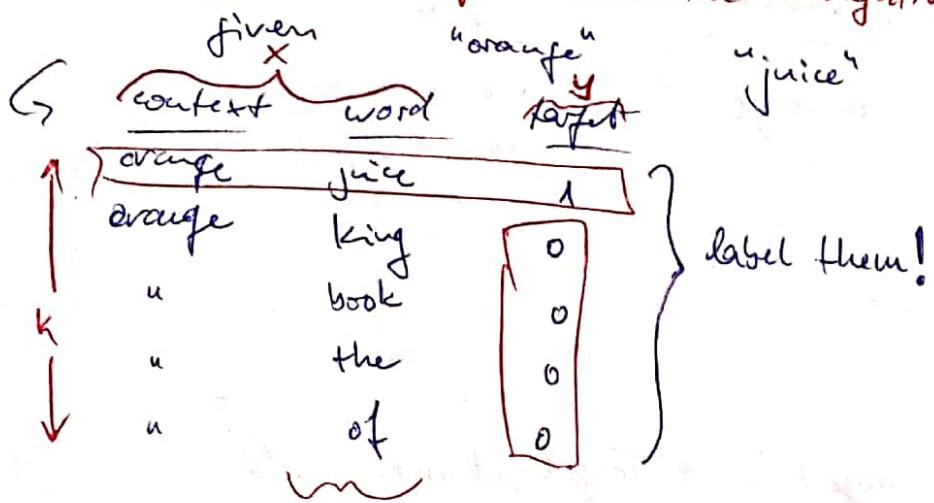
→ separate common from uncommon words

Negative sampling

[Nikolov 2013, Distributed representation of words → plurals and their compositionality]

Defining a new learning problem

→ have one positive example and choose k -negative words to train the model more efficiently



How to sample negative context words?
 → empirical frequency $f(w_i)$ in language

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

Model

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

logistic regression-like

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

$$\theta_{est} \rightarrow E \rightarrow e^{\theta_{est}^T c}$$

$k = 5-20$ for smaller datasets

$k = 2-5$ for large datasets

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

for every positive example we have k -negative examples to train the model
 reduces computation cost

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

c, t

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

$x_{ij} = x_{ji}$ - symmetrical relationship

Model

$$\boxed{\text{minimize}_{\Omega, e} \sum_{i=1}^{10000} \sum_{j=1}^{10000} f(x_{ij}) \quad (\Omega_i^T e_j + b_i + b_j' - \log x_{ij})^2}$$

now symmetric

weighting term

$$f(x_{ij}) = 0 \text{ if } x_{ij} > 0$$

$\Omega_i^T e_j$ - inner product as indicator

$0 \log 0 = 0$

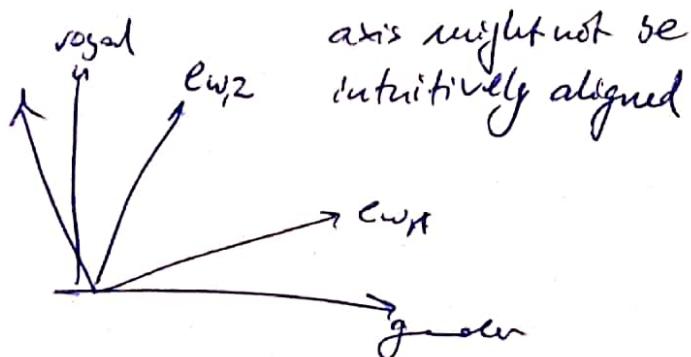
$\Omega_i^T e_j$ are symmetric

now related are c and t ?

	Man	Woman	King	Queen
Man	1			
Woman		1		
King			1	
Queen				1

- one way to train the algorithm is to initialize both randomly and run GD and then for every word take the average

$$\boxed{(final) e_w = \frac{e_w + \delta_w}{2}}$$



Sentiment classification

restaurant ratings x map $\rightarrow y$

The dessert is excellent

4 stars

Service was slow

2 stars

nothing special

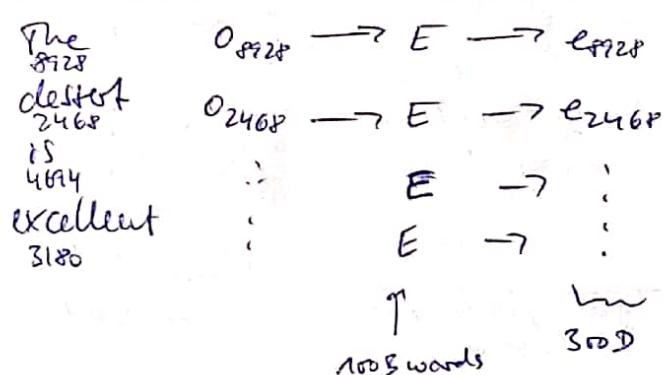
3 stars

bad

1 stars

"completely lacking in good taste,
good service and good
ambience"

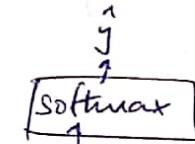
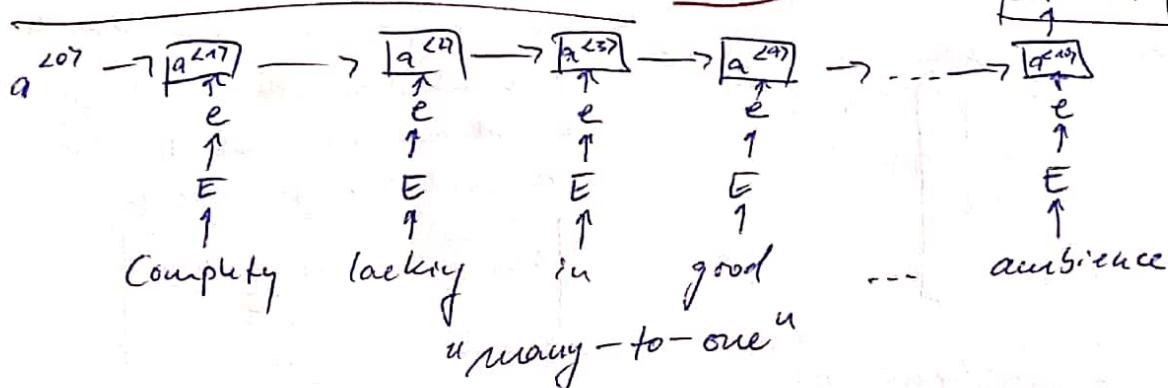
Simple sentiment classification model



→ cons: ignores words' order
G bad sentence can have the word "good" often

RNN for sentiment classification

better



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

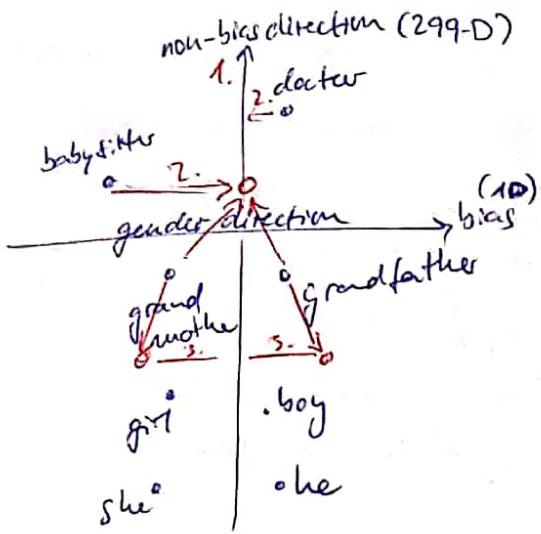
bias: gender, ethnicity...

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 if the text used to train the model.



1. identify bias direction

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

→ average them → find out gender direction

2. Neutralize: For every word that is not definitional, project to get rid of bias
project "doctor", "babysitter" to axes to reduce/eliminate gender component

3. Equalize parts (hand-pick)

Grandmother — Grandfather

girl — boy

→ move grandmother/father to that distance to non-bias axis is equal!

DL Spec / course V / week 3

Basic Models

Sequence models & Attention mechanisms

[Sutskever 2014 seq2seq learning w NN]

Sequence to sequence model

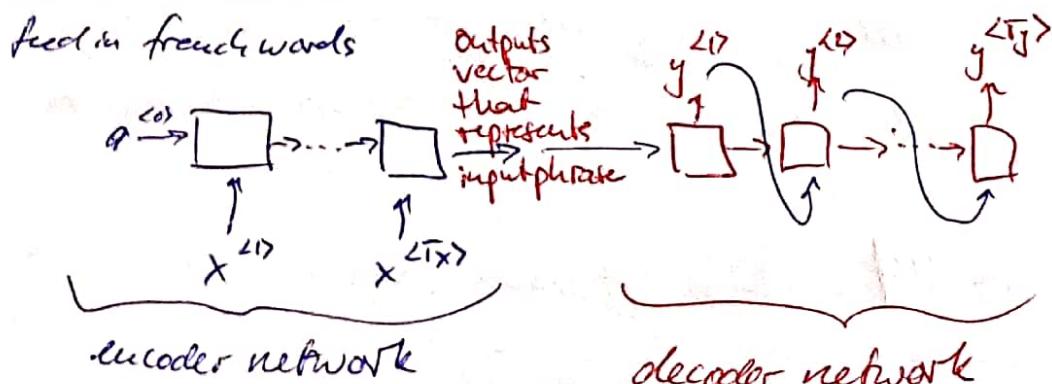
$x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad x^{(4)} \quad x^{(5)}$
 Jane visite l'Afrique en septembre

[Cho 2014, learning phrase representations using RNN encoder-decoder for statistical machine trans.]

→ Jane is visiting Africa in September.

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

first encoder Network



decoder network
outputting translation
word after word

Image captioning

$y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad y^{(4)} \quad y^{(5)} \quad y^{(6)} \quad y^{(7)}$
 A cat is sitting on a chair



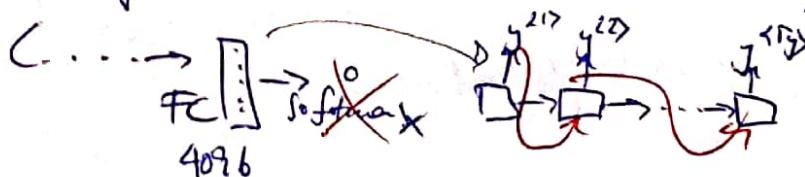
→ use a pre-trained AlexNet CNN

[Mao, 2014, Deep captioning with multi-modal RNN]

[Vinyals 2014, Neural image caption generator]

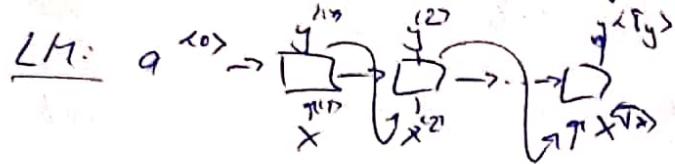
[Karpathy, Feifei, 2015, Deep visual-semantic alignments for generating image descriptions]

→ get rid of Softmax and keep the 4096 FC

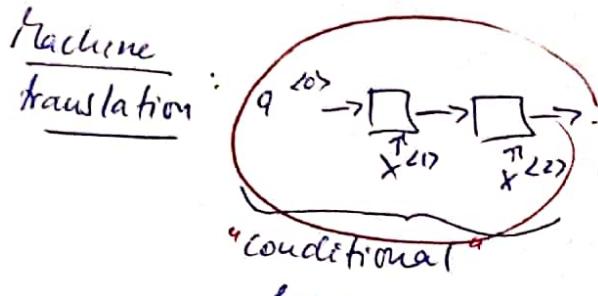


Picking the most likely sentence

Machine translation or building a conditional language model (LM)



LM gives a probability of a sentence
(generate novel sentences as well)



Instead of modelling the probability of a sentence

trying to estimate the probability of some english translation

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

is now modelling the probability output english translation

$$P(y^{(1)}, \dots, y^{(T_y)} | x^{(1)}, \dots, x^{(T_x)}) \text{ conditioned on some input french sentence}$$

English

French

→ do not use random sampling!

→ rather: find English sentence y that maximizes that conditional probability

$$\underset{y^{(1)}, \dots, y^{(T_y)}}{\operatorname{argmax}} P(y^{(1)}, \dots, y^{(T_y)} | x)$$

Need an algorithm that can find the value y that maximizes

→ Beam Search

Why not a greedy search? DL spec course II week 3

→ generate the first most likely word according to
conditional language model } $P(y^{1 \dots T} | x)$

→ after picking 1st word pick the 2nd which is most likely
and so on

What want is the sequence joint probability $P(y^{1 \dots T} | x)$

(A) → "Jane is visiting Africa in September"

(B) → "Jane is going to be visiting Africa in September"

$P(A|x) < P(B|x)$ } ↗ is going is more likely than "is visiting" so the 2nd phrase will be chosen which is not optimal
and only because of the word-by-word-prediction.

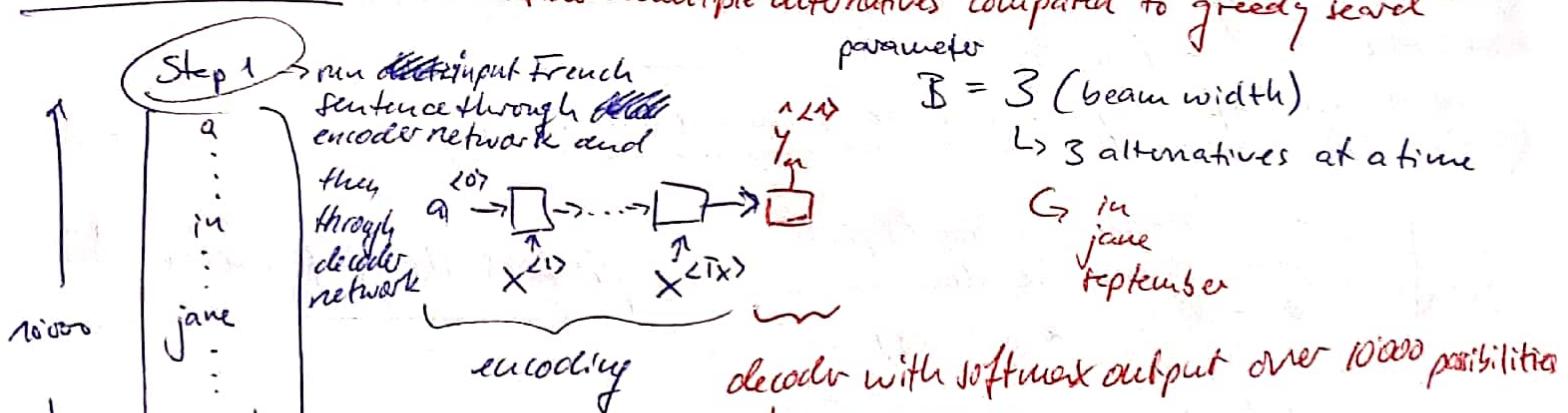
total # of words in an English sentence is exponentially large
if 10'000 w in dict. and a sentence of 10 words long

then there are $10^{'000}^{10}$ possibilities

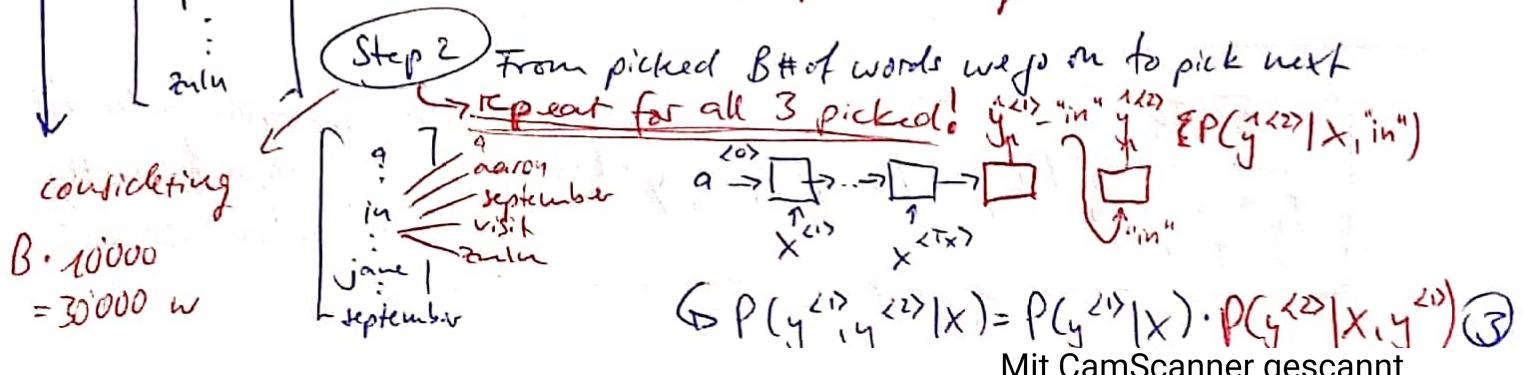
so the most common thing to do is use an approximate search algorithm.
→ will try to pick the sentence y that maximizes the conditional probability

Beam Search → can consider multiple alternatives compared to greedy search

Step 1



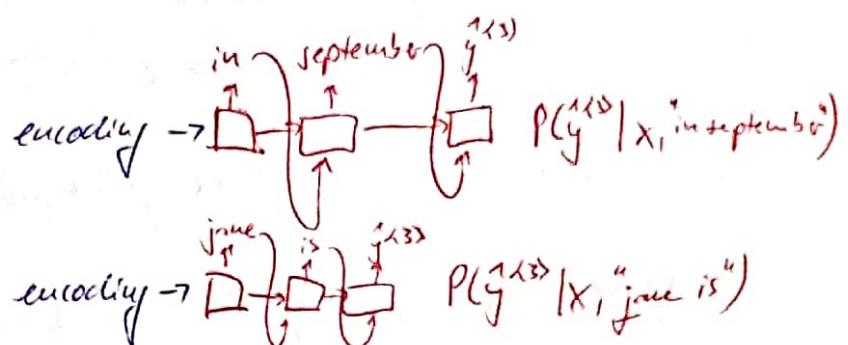
Step 2



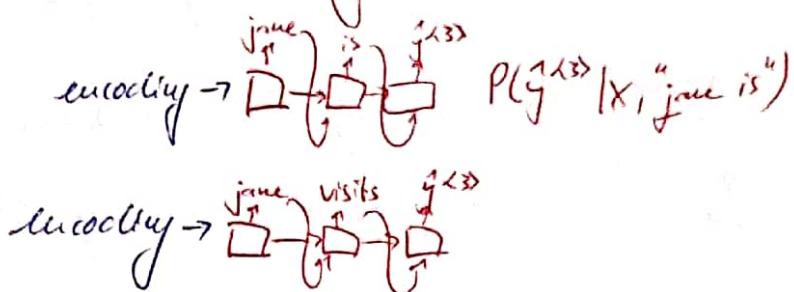
If Beam search decides that 3 most likely are not including "september" as starting word it is also rejecting "September" as first word for translation.

again for Beam search ($B=3$)

in september aaron
jane
zulu



jane is aaron
zulu



jane visits aaron
africa
zulu

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

\Rightarrow hopeful outcome is "jane visits africa in september" <EOS>

refinements to Beam search

Length normalization:

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

(product of probabilities)

In practice:

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

inserting a log here

tums the product of a log
into a sum of a log and increases numerical stability

~~maximizing $\log P(y | x)$
should give the same
as $\max P(y | x)$~~

prefers short sentences!

next modification

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

$x = 0.7$

normalization reduced
penalty significantly
for outputting longer
sentences

(4)

Beam search discussion

Beam width B ? \rightarrow size of B steers computational time

- | | |
|-------------------------------------|---|
| 1
large B : | better results bc more probs but slower |
| ~ 10
~ 10
up to 100 | small B : worse results, faster |
| increased > 1000 | |

Error Analysis in Beam Search (BS)

Jane visits l'Afrique au septembre

Human: Jane visits Africa in September (y^*) good translation!

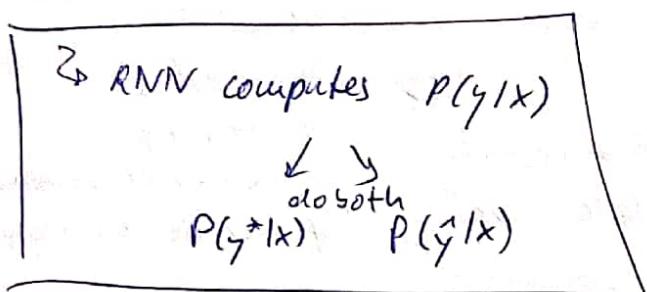
Algorithm: Jane visited Africa in Sept. (\hat{y}) bad, changed meaning!

encoder \rightarrow decoder

\rightarrow RNN

\rightarrow Beam search

\hookrightarrow B we know from ANN error analysis that getting more training data never hurts so BS, BUT might not help always.



Case 1: $P(y^*|x) > P(\hat{y}|x)$: BS chose \hat{y} . But y^* attains higher $P(y|x)$

\hookrightarrow BS failed

Case 2: $P(\hat{y}|x) \geq P(y^*|x)$: y^* is a better translation than \hat{y} but RNN predicted $P(y^*|x) < P(\hat{y}|x)$

\hookrightarrow RNN failed

Error Analysis Process → go through devset and find mistakes

Human	Algorithm	$P(y^* x)$ RNN	$P(\hat{y} x)$ BS	At fault?
Jane...	Jane...	2×10^{-10}	1×10^{-10}	BS
...	RNN
		:	:	BS
				BS
				RNN
				:

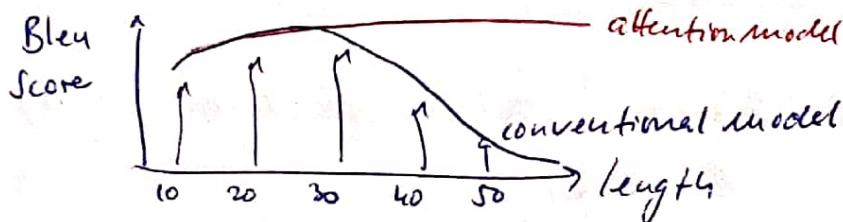
Figures out what fraction of errors are "due to" beam search
vs. RNN model
(RNN generates objective function!)

Bleu score (optional) [Papineni 2002, A method for automatic evaluation of machine translation]
Bilingual evaluation underway → how good is machine translation

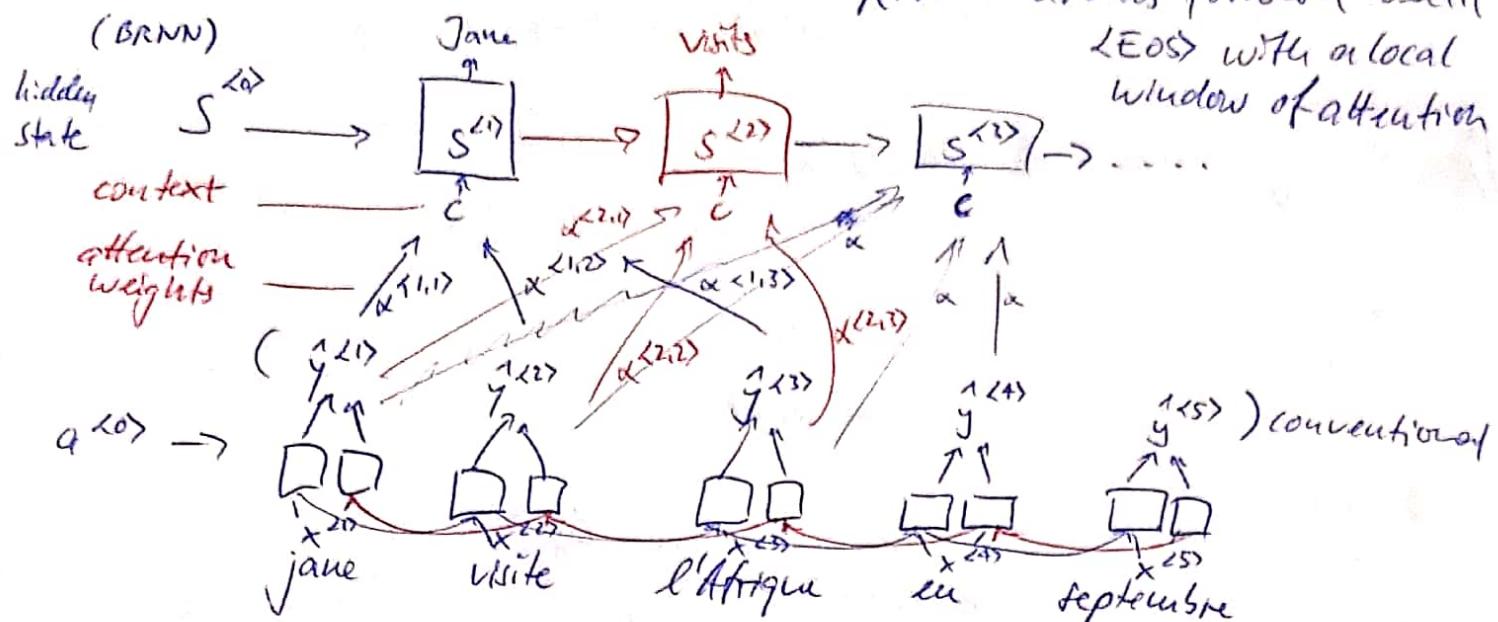
Attention Model intuition [Bahdanau 2014, Neural MT by jointly learning to align + translate]

input very long French sentence encoding → decoding (memorize sentence)

→ human would not memorize whole sentence
but a few words at a time + translate "live"
+ store in activations
to generate translation



Bi-directional RNN here:



RNN marches forward until EOS with a local window of attention

Contexts change with advancing words to translate

Attention Model \rightarrow pay attention to a part of a sentence or a fine parts from BRNN \leftarrow forward backward

Grad step
generates a new context

$$\sum_{t'} \alpha^{<1,t'} = 1$$

$$c^{<1>} = \sum_{t'} \alpha^{<1,t'} a^{<t>}$$

$$a^{<t>} = (\overrightarrow{a^{<t>}}, \overleftarrow{a^{<t>}})$$

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

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

[Bahdanau...]

[Xu, 2015, Show attention and tell: neural image caption generation with visual attention]

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

use softmax to sum to 1

how compute vectors $e^{<\cdot>}$?

\rightarrow Use small NN

$$s^{<t-1>} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow e^{<t,t'>}$$

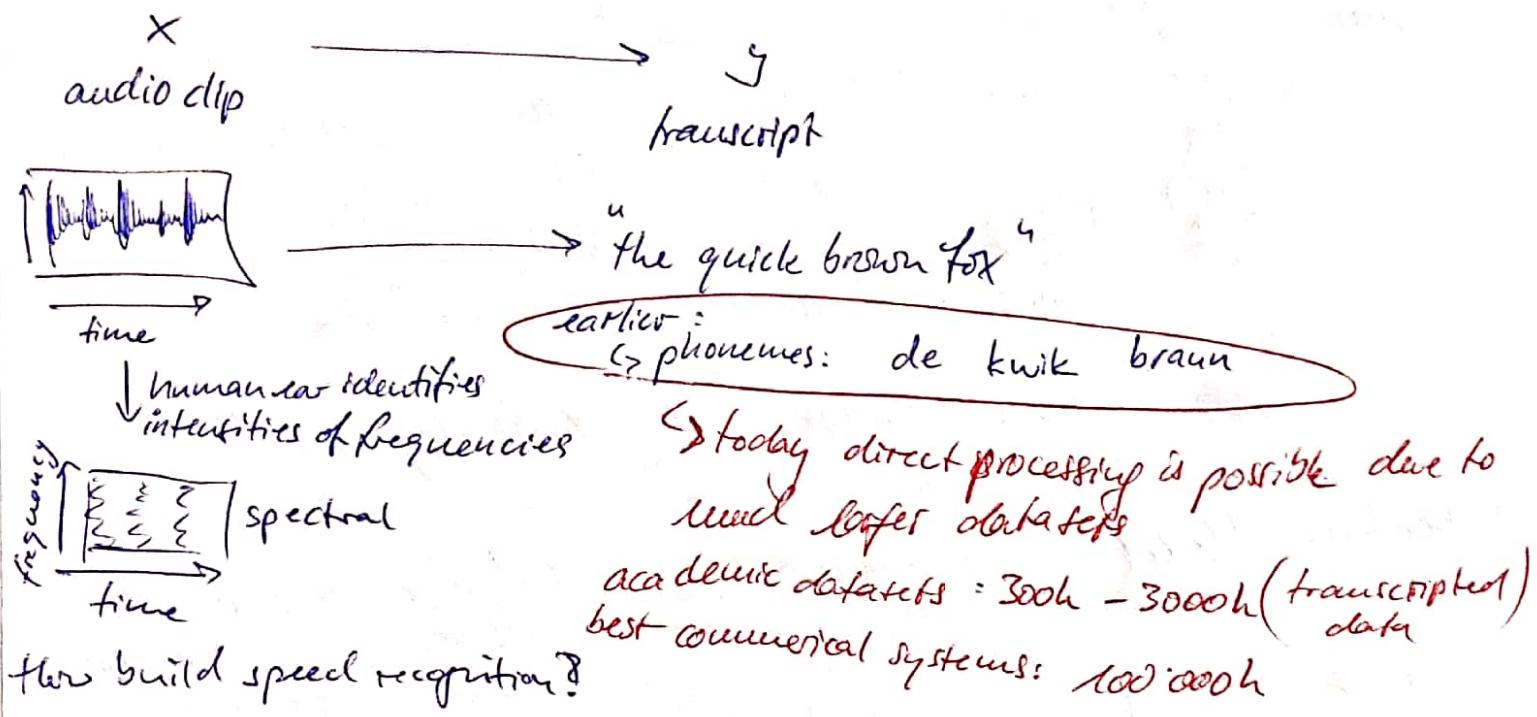
Con: takes quadratic time to run $\Theta(n^2)$

Attention examples

July 20th 1969 \rightarrow 1969-07-20

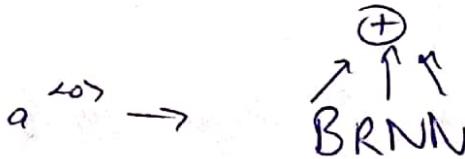
23 April, 1564 \rightarrow 1564-04-23

Speech recognition - Audio data



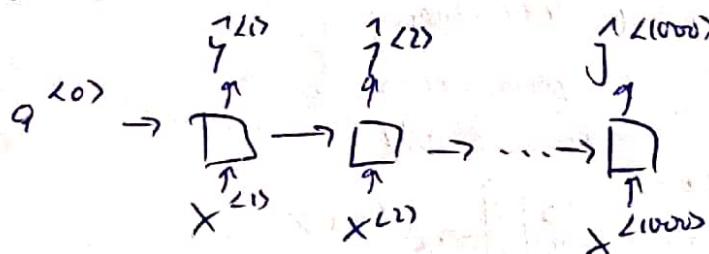
could do:

Attention model \rightarrow output transcript



what works well is: CTC cost for speech recognition data with RNN
 [Graves 2006, Labeling unsegmented sequence (connectionist temporal classification)]

"the quick brown fox"

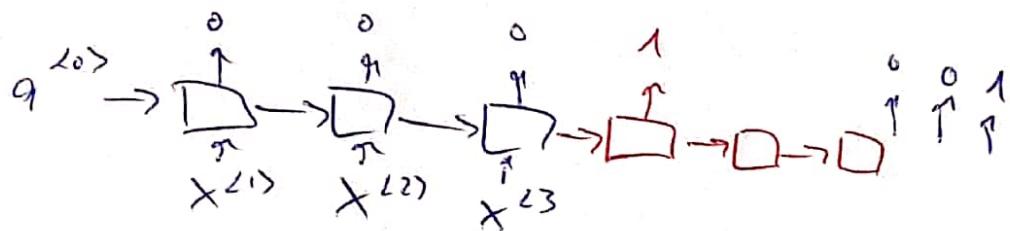


CTC cost allows
 $ttt_h_eee_--_u_--999-$
 \hookrightarrow correct output for "the
 quick brown fox"

Rule: collapse repeated characters not

Trigger Word Detection

Alexa, Baidu, Hey Siri, OK Google

Algorithm

Teacher forcing | Teacher forcing | Teacher forcing

disadvantage: creates

unbalanced

training set with

more 0s than 1s

Jet target

label

again

target label

→ fit: 0000 111 000 111

→ make it output more often

- | | | | |
|--------------------------|-----------|---|--|
| 1 True X | 2 False ✓ | 3 | |
| slowly moving
between | → u - ✓ | u | |
-
- | | | | |
|-------------|----------|---|--|
| 4 No, RNN ✓ | 5 True ✓ | 6 sum over X ^{$\sum t^i +$} \dots for t^i X $\sum t^i & \dots$ for t | |
| → u ✓ | ✓ | ✓ | |
-
- | | | | |
|----------|-----------|---------------------------|--|
| 7 True ✓ | 8 large ✓ | 9 cookbook X → cookbook ✓ | |
| → ✓ | → ✓ | u | |
-
- | | | | |
|----------------------------|--|--|--|
| 10 t-th input X → Features | | | |
| ✓ | | | |