

 intro

deeplearning.ai

Recurrent Neural Networks

Why sequence models?

Examples of sequence data

Speech recognition		"The quick brown fox jumped over the lazy dog."	
Music generation			
Sentiment classification	"There is nothing to like in this movie."		
DNA sequence analysis		AGCCCCCTGTGAGGAACCTAG	
Machine translation	Voulez-vous chanter avec moi?	Do you want to sing with me?	
Video activity recognition		Running	
Name entity recognition	Yesterday, Harry Potter met Hermione Granger.	Yesterday, Harry Potter met Hermione Granger. Andrew Ng	

 notation



deeplearning.ai

Recurrent Neural Networks

Notation

Motivating example

NLP

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

$\xrightarrow{q \text{ words}}$ $x^{(1)}$ $x^{(2)}$ $x^{(3)}$... $x^{(t)}$... $x^{(q)}$

$O = \text{Yes, part of name.}$

$\rightarrow y:$ $y^{(1)}$ $y^{(2)}$ $y^{(3)}$ $y^{(4)}$... $y^{(15)}$

label $y^{(1)}$ $y^{(2)}$ $y^{(3)}$ $y^{(4)}$... $y^{(15)}$

$T_x = 9$

$T_y = 9$ length of sequence

$T_x^{(i)} = 9$

$T_y^{(i)}$

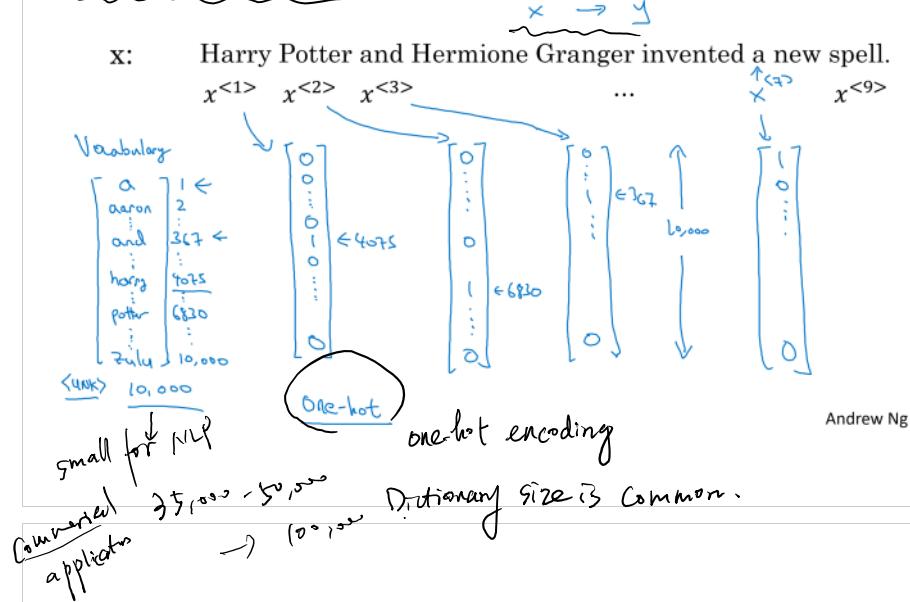
where are the people's name : task NER
application: search engines

could be other output representation.

(i) ;ⁱth training example

Andrew Ng

Representing words



Representing words

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

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

And = 367
 Invented = 4700
 A = 1
 New = 5976
 Spell = 8376
 Harry = 4075
 Potter = 6830
 Hermione = 4200
 Gran... = 4000

Andrew Ng

How to represent individual words?

① Come up with a vocab / a dictionary.

(later: What if encountering words not in your vocab? Create a category "unknown")

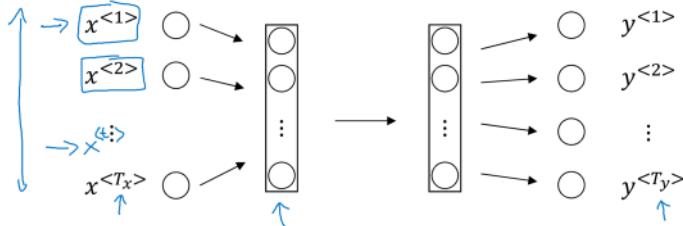


Recurrent Neural Networks

deeplearning.ai

Recurrent Neural Network Model

Why not a standard network?



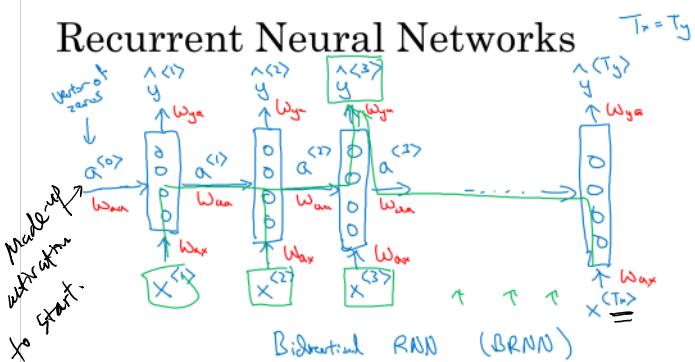
Problems: Why this does not work well?

- Inputs, outputs can be different lengths in different examples.
- Doesn't share features learned across different positions of text.

Andrew Ng

→ maybe padding, if there's a max length.
"Harry" appearing at a different position.

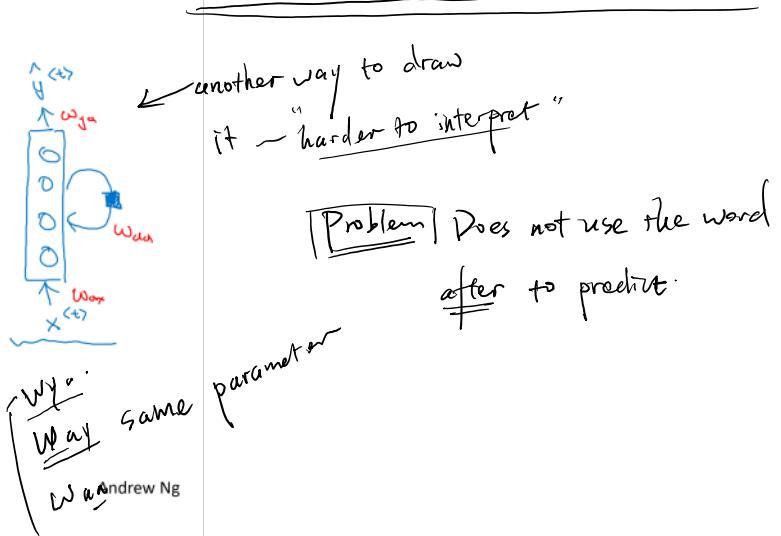
Recurrent Neural Networks



He said, "Teddy Roosevelt was a great President."

He said, "Teddy bears are on sale!"

What is a recurrent neural network.



[Problem] Does not use the word after to predict.

param
var

Forward Propagation $a \leftarrow W_a x^{(t)}$

$$\hat{y}^{(1)} \quad \hat{y}^{(2)} \quad \hat{y}^{(3)} \quad \hat{y}^{(T_y)}$$

从前文进阶部分，可以知道。

to simplify this notation.

W_a 和 W_a , W_a stack \rightarrow .

即

$$[W_a \mid W_a] \begin{bmatrix} a^{(t-1)} \\ \vdots \\ a^{(0)} \end{bmatrix} \begin{bmatrix} x^{(t)} \\ \vdots \\ x^{(0)} \end{bmatrix}$$

Next: backprop

$$\delta_{(t)} = \partial(\lambda^2 \alpha_{(t)} + p^2)$$

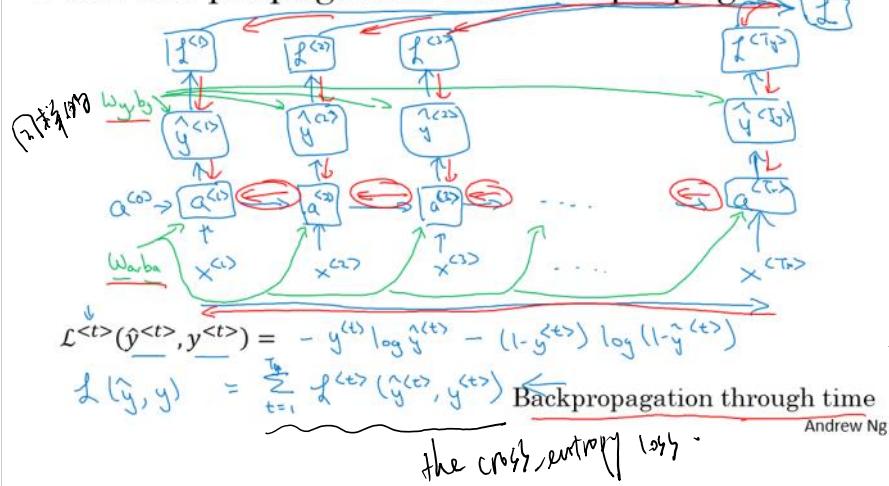
$$\delta_{(t-1)} = \partial(\lambda^2 \alpha_{(t-1)} + p^2)$$

$$\alpha_{(t)} = \partial(\lambda^2 \alpha_{(t-1)} + N^2 x_{(t)} + p^2)$$

notation NN building

backprop

Forward propagation and backpropagation



Back propagation through time
the cross entropy loss

rnn different type

T_x need not $= T_y$ the length of input and output need not be equal.

Recurrent Neural Networks

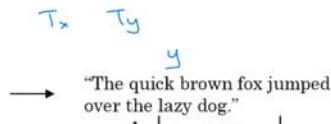


deeplearning.ai

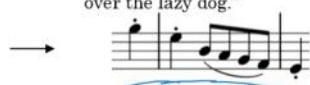
Different types of RNNs

Examples of sequence data

Speech recognition



Music generation



Sentiment classification

"There is nothing to like in this movie."



DNA sequence analysis

AGCCCCCTGTGAGGAACCTAG

→ AGCCCCTGTGAGGAAC $\color{red}{T}$ AG

Machine translation

Voulez-vous chanter avec moi?

→ Do you want to sing with me?

Video activity recognition



→ Running

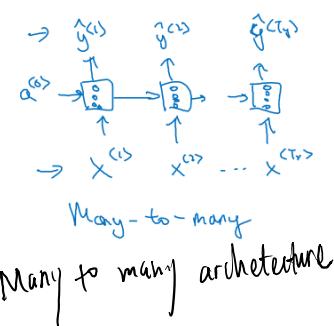
Name entity recognition

Yesterday, Harry Potter met Hermione Granger.

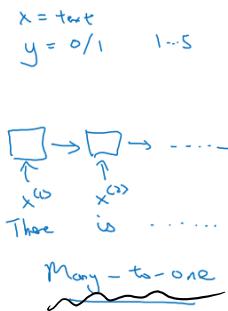
→ Yesterday, Harry Potter met Hermione Granger.
Andrew Ng

Examples of RNN architectures

$$T_x = T_y$$



$$T_x = T_y$$



$$\text{Sentiment classification}$$

$$x = \text{text}$$

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

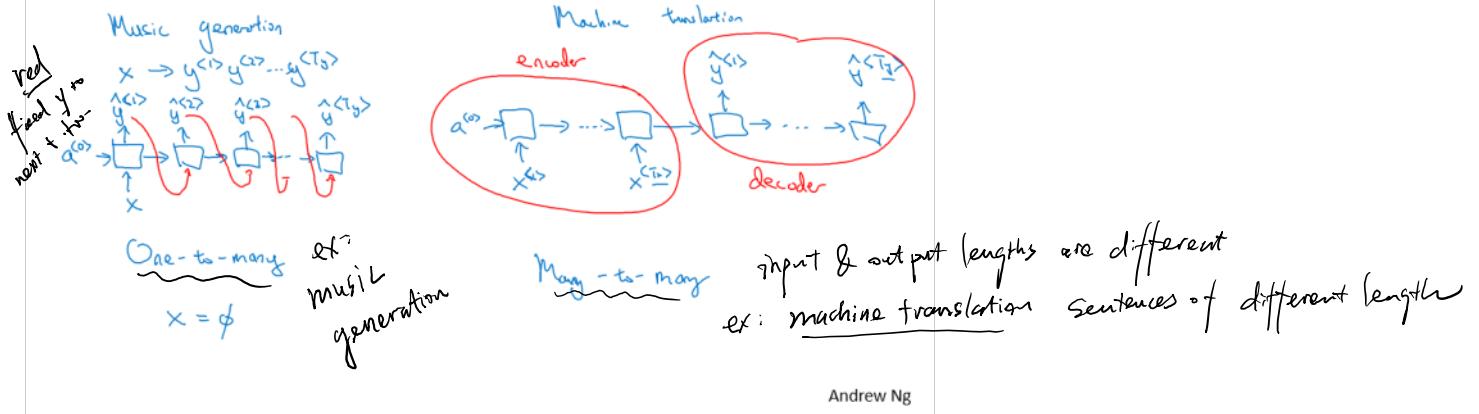


Many to one

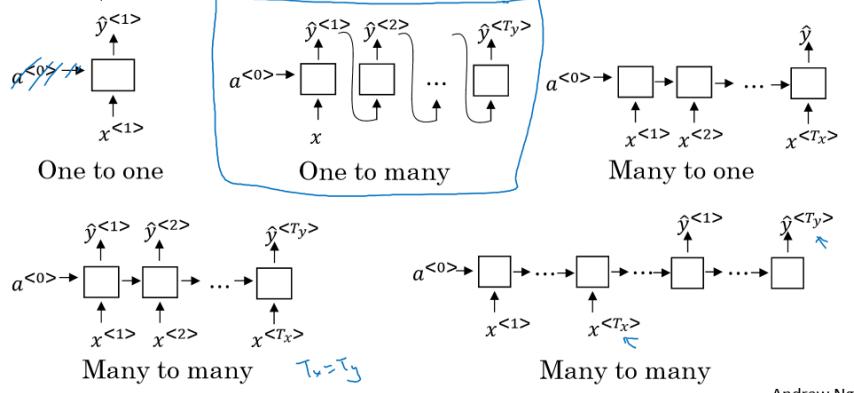
example : movie review

Andrew Ng

Examples of RNN architectures



Summary of RNN types



★ ★ good summary and graph.

later "attention-based architecture."
not in this list.

up next: sequence generation



deeplearning.ai

Recurrent Neural Networks

Language model and sequence generation

What is language modelling?

Speech recognition

The apple and pair salad.

→ The apple and pear salad.

$$P(\text{The apple and pair salad}) = 3.2 \times 10^{-13}$$

$$P(\text{The apple and pear salad}) = 5.7 \times 10^{-10} \quad \text{more likely.}$$

$P(\text{sentence}) = ?$ $P(y^{(1)}, y^{(2)}, \dots, y^{(T)})$ a particular sequence of words

Andrew Ng

Language modelling with an RNN

Training set: large corpus of english text.

Tokenize — form a vocab

ignore punctuation here

Cats average 15 hours of sleep a day. $\langle \text{EOS} \rangle$

$y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad \dots \quad y^{(8)} \quad y^{(9)}$

 $x^{(t)} = y^{(t-1)}$

The Egyptian ~~Mau~~ is a bread of cat. $\langle \text{EOS} \rangle$

10,000

$\langle \text{UNK} \rangle$

Andrew Ng

what if some words are not in vocab?

recode as [unknown]

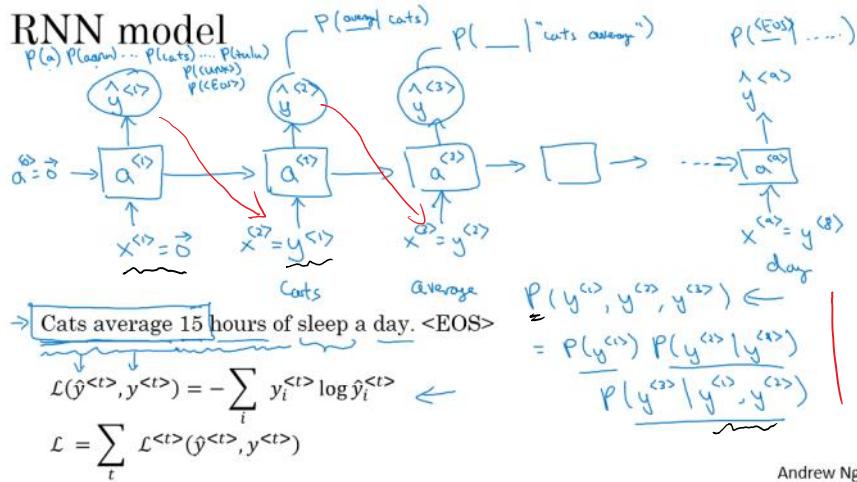
one-hot encoding $\begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$

tiny y conditional probability $p(y|z)$

Model the chance the sentences with RNN

这模型不是 model probability.
是? prob model sigmoid.
但不是 sigmoid ?

Hum exercise 用看是怎样的。



Sample novel sequences.

sample sequence

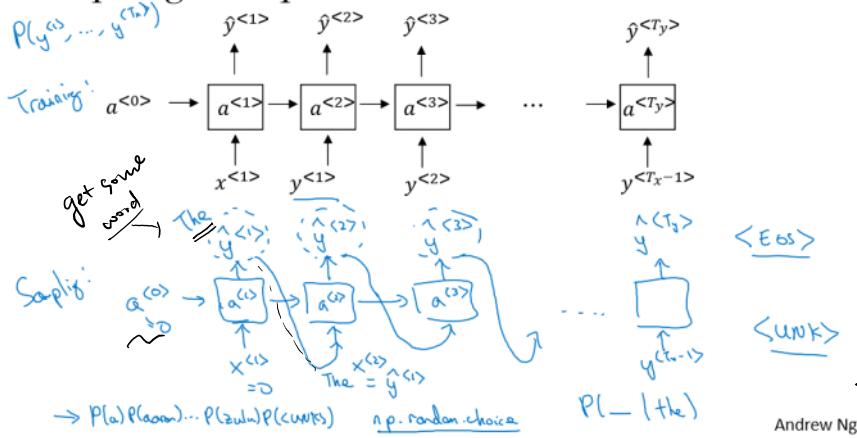


Recurrent Neural Networks

deeplearning.ai

Sampling novel sequences

Sampling a sequence from a trained RNN

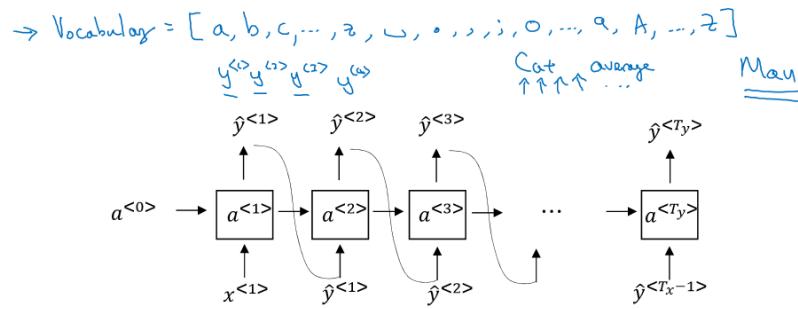


randomly sample ~~visually~~ the softmax distribution

if $[EOS]$ in your vocab, keep sampling until you hit it.
If not, can set the target # of words
can force out Euler)

Character-level language model

→ Vocabulary = [a, aaron, ..., zulu, <UNK>] ←



can do character-level

Same rationale for character level

adv: No need to worry about unknown words.

disad: much longer sequences!

• ⇒ more computationally expensive

• usually word level is used
- character level

$x^{<1>} \quad \hat{y}^{<1>} \quad \hat{y}^{<2>} \quad \dots \quad \hat{y}^{<T_x-1>}$

Andrew Ng

Usually word level is used
But in some cases character level
Rarely used for cases you
need to deal with lots of
unknown words.

Sequence generation

News

President enrique peña nieto, announced
sench's sulk former coming football langston
paring.

"I was not at all surprised," said hich langston.

"Concussion epidemic", to be examined. ↪

The gray football the told some and this has on
the uefa icon, should money as.

Shakespeare

The mortal moon hath her eclipse in love.

And subject of this thou art another this fold.

When lesser be my love to me see sabl's.

For whose are ruse of mine eyes heaves.

→ Assignment content

Andrew Ng

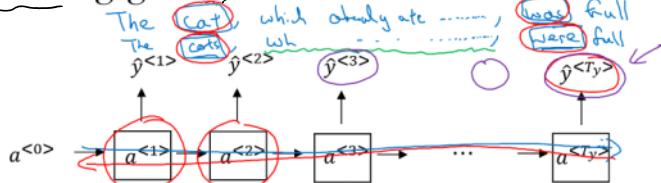
vanish gradient



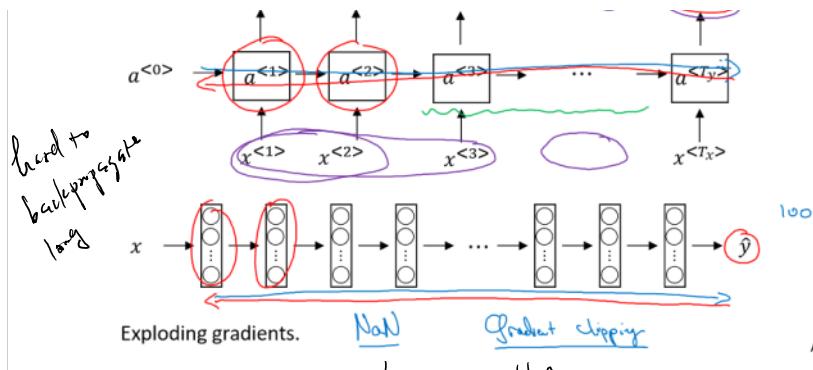
Recurrent Neural Networks

Vanishing gradients with RNNs

Vanishing gradients with RNNs



Very long term dependency
the RNN introduced above does not take care
of this long term dependency



Vanishing gradient is a bigger problem.
But exploding gradient

is also bad — easy to observe though. NaN of param
Solution: clip it gradient clipping

Next few videos: focus on vanishing gradient

of this 10^{-10}

Need to memorize for a very long time.
previous. ~~if~~ local influence
can't back propagate to many steps
before

Andrew Ng

gru

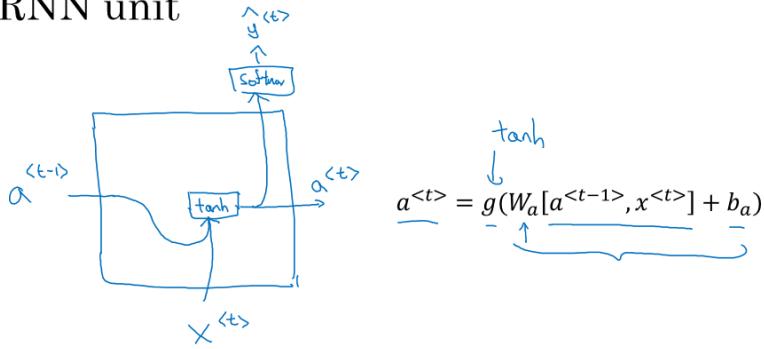
Recurrent Neural Networks



deeplearning.ai

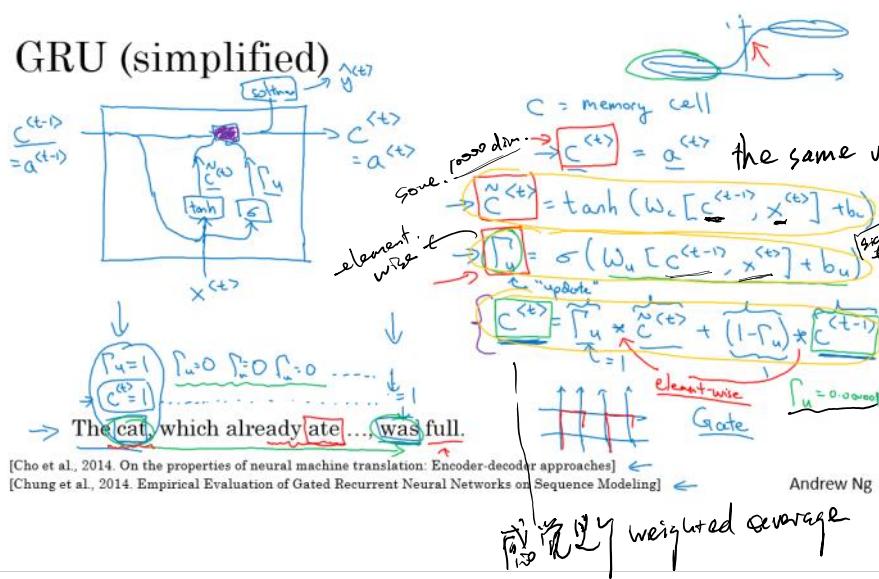
Gated Recurrent Unit (GRU)

RNN unit



Andrew Ng

GRU (simplified)



C: memory cell

Q 請問 Γ_u 是用來做甚麼？

the same value

gate, decide whether to update [binary]
the gate decide whether to update

if $\Gamma_u = 1$: update $\Gamma_u \in \{0, 1\}$ 否則?
if $\Gamma_u = 0$: not update. still the old c-

interesting!

Γ_u close to zero

9

Full GRU

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

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

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

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

|
alternative notation

The cat, which ate already, was full.

LSTM

? relevant

?

得讀甚麼？GRU
my paper

GRU one of the most commonly used
method to tackle
vanishing gradient &
model Long Term memory

Andrew Ng



deeplearning.ai

Recurrent Neural Networks

LSTM (long short term memory) unit

GRU and LSTM

GRU

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

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

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

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

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

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

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

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

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

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

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

the paper is difficult to read . more about the
vanishing gradient problem

more powerful & general model than GRU.

3 Gates slightly more complicated.

Andrew Ng

[Hochreiter & Schmidhuber 1997. Long short-term memory]

LSTM units

GRU

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

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

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

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

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

[Hochreiter & Schmidhuber 1997. Long short-term memory]

LSTM

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

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

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

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

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

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

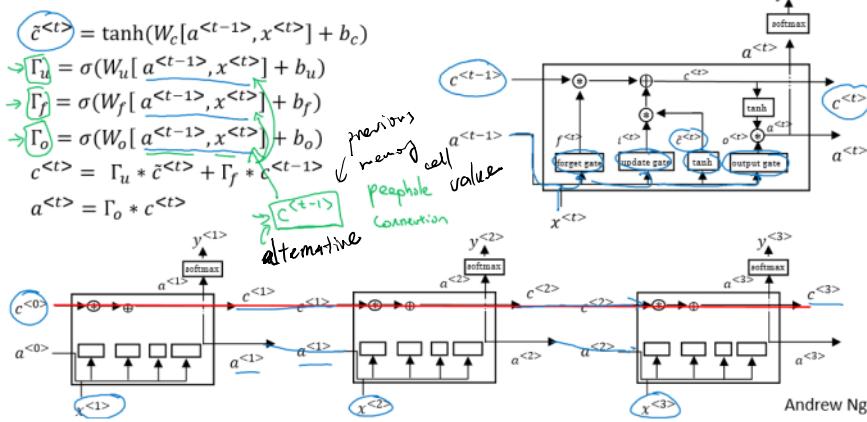
Andrew Ng

One variation of LSTM.

Peephole connection

1/3

LSTM in pictures



good at memorizing things for a long time

GRU or LSTM? No consensus. LSTM comes much earlier.

GRU is treated as a simplified model.

Andrew Ng: GRU simpler. easier to build a bigger model

LSTM is a more historically proven model.

(Many teams will use it as the default.)

But GRU is gaining momentum in recent yrs

bidirectional rnn



deeplearning.ai

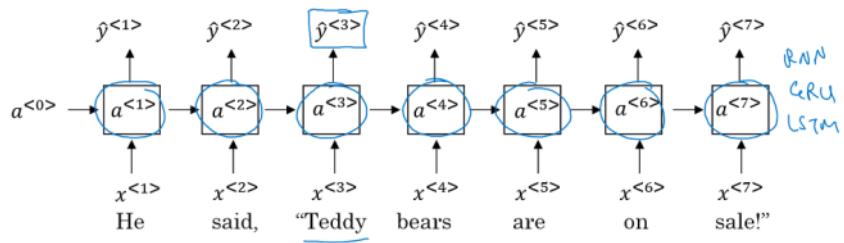
Recurrent Neural Networks

Bidirectional RNN

Getting information from the future *later words help determine earlier words*

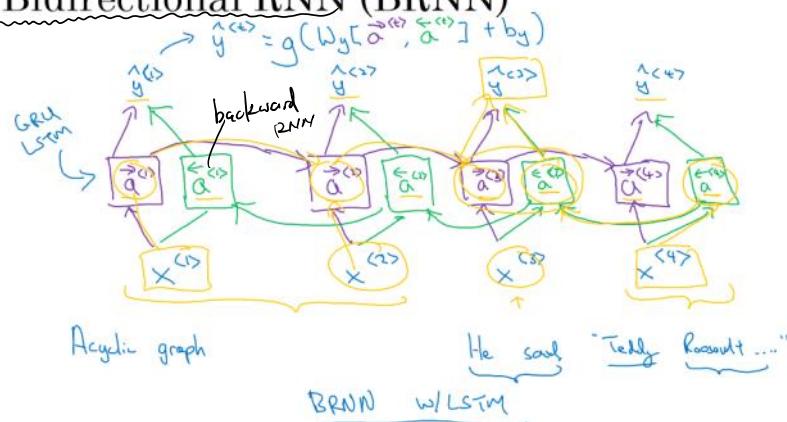
He said, "Teddy bears are on sale!"

He said, "Teddy Roosevelt was a great President!"



Andrew Ng

Bidirectional RNN (BRNN)



Acyclic graph

can be GRU, LSTM blocks

Bidirection + LSTM commonly used

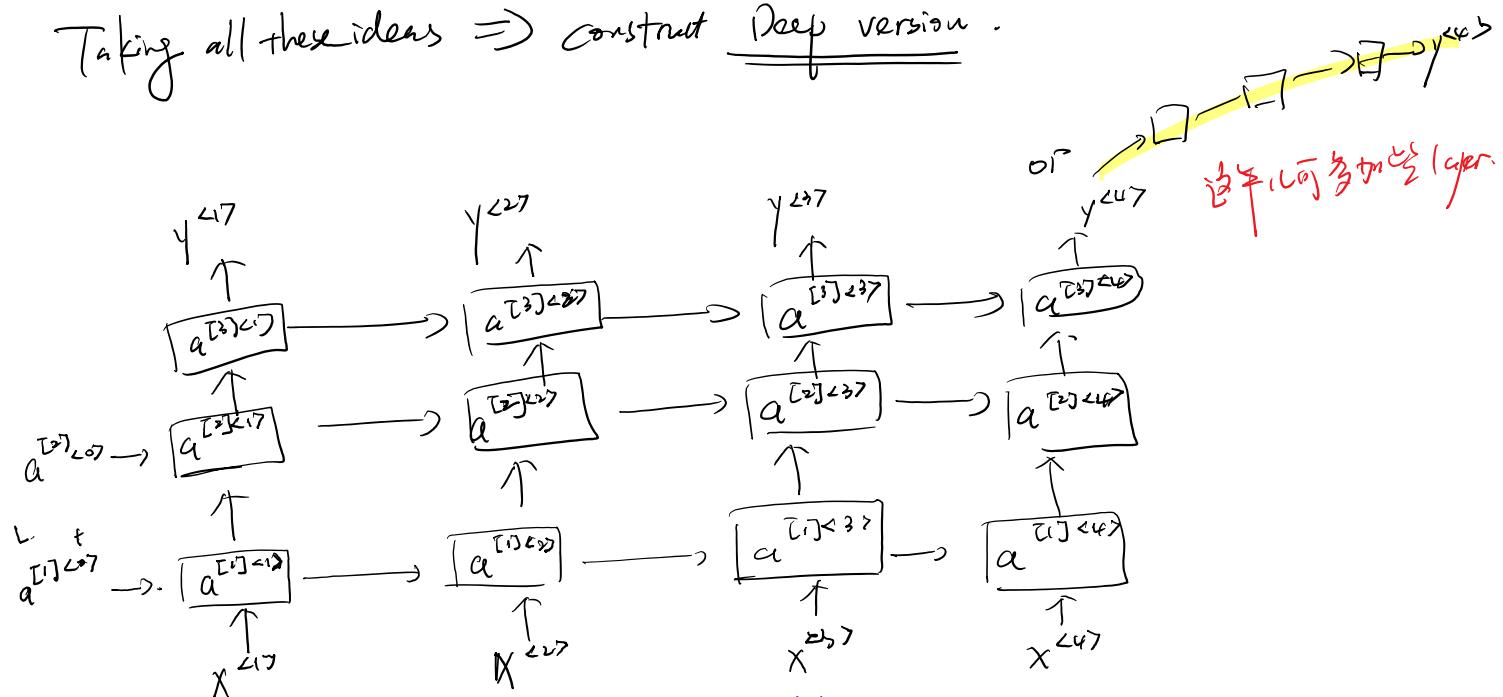
A reasonable first thing to try.

Andrew Ng

Disadvantage: Need the entire sequence before making prediction

e.g. Speech recognition. Start processing after the person stop talking

Taking all these ideas \Rightarrow construct Deep version.



(RNN, 4 layers (3 layers + 1 output))

Two RNN models are used.