# SIT330-770: Natural Language Processing

## Week 6 - Neural Networks for NLP

**Dr. Mohamed Reda Bouadjenek**

School of Information Technology, Faculty of
Sci Eng & Built Env

reda.bouadjenek@deakin.edu.au

Andrew Ng

Neural Networks and Deep Learning

(Optional)

# SIT330-770: Natural Language Processing

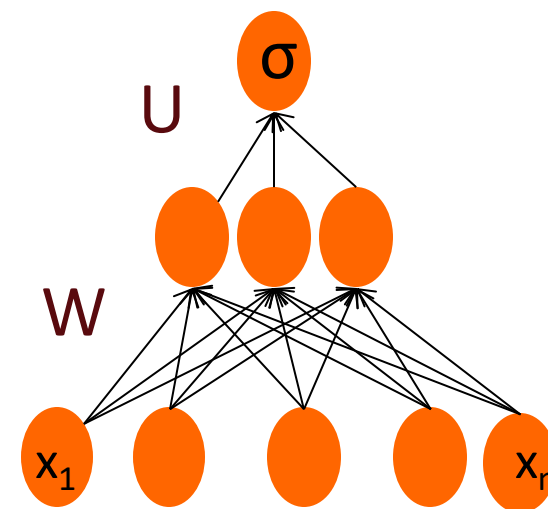## Week 6.11 - Applying feedforward networks to NLP tasks

**Dr. Mohamed Reda Bouadjenek**

School of Information Technology,
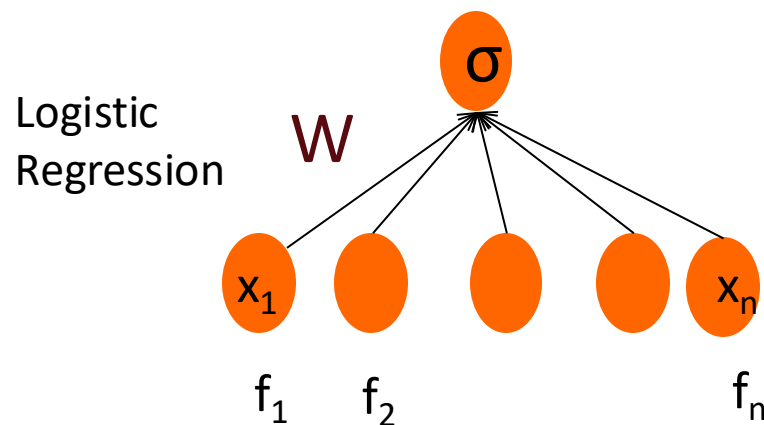Faculty of Sci Eng & Built Env

- Let's consider 2 (simplified) sample tasks:

  1. Text classification
  2. Language modeling

- State-of-the-art systems use more powerful neural architectures, but simple models are useful to consider!

# Classification: Sentiment Analysis

- We could do exactly what we did with logistic regression

- Input layer are binary features as before

- Output layer is 0 or 1

# Sentiment Features
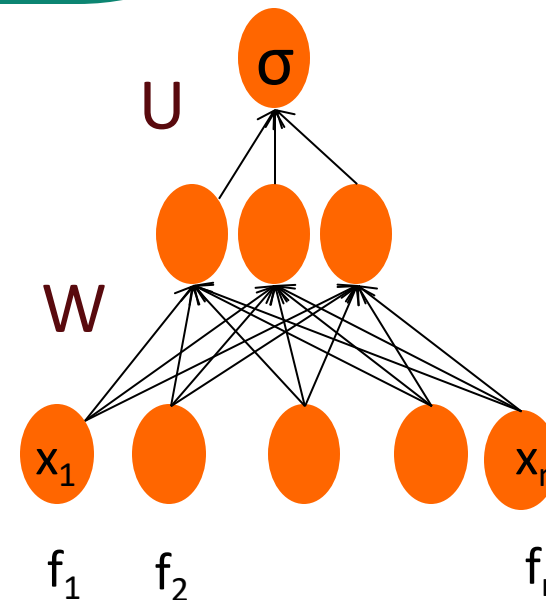
| Var | Definition |
| --- | --- |
| $x_1$ | count(positive lexicon) $\in$ doc) |
| $x_2$ | count(negative lexicon) $\in$ doc) |
| $x_3$ | $\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$ |
| $x_4$ | count(1st and 2nd pronouns $\in$ doc) |
| $x_5$ | $\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$ |
| $x_6$ | log(word count of doc) |

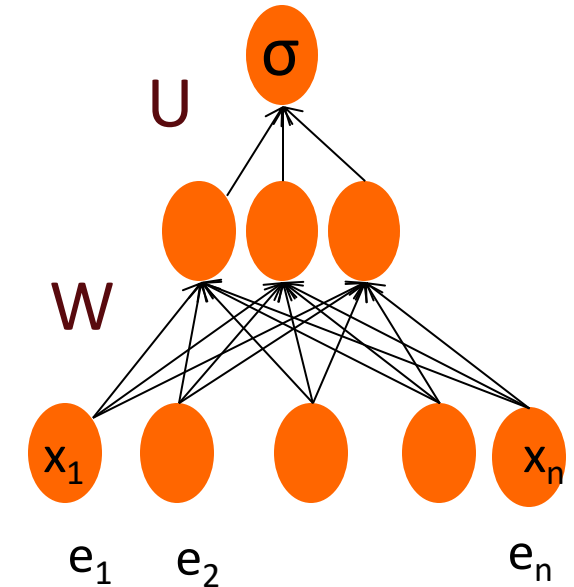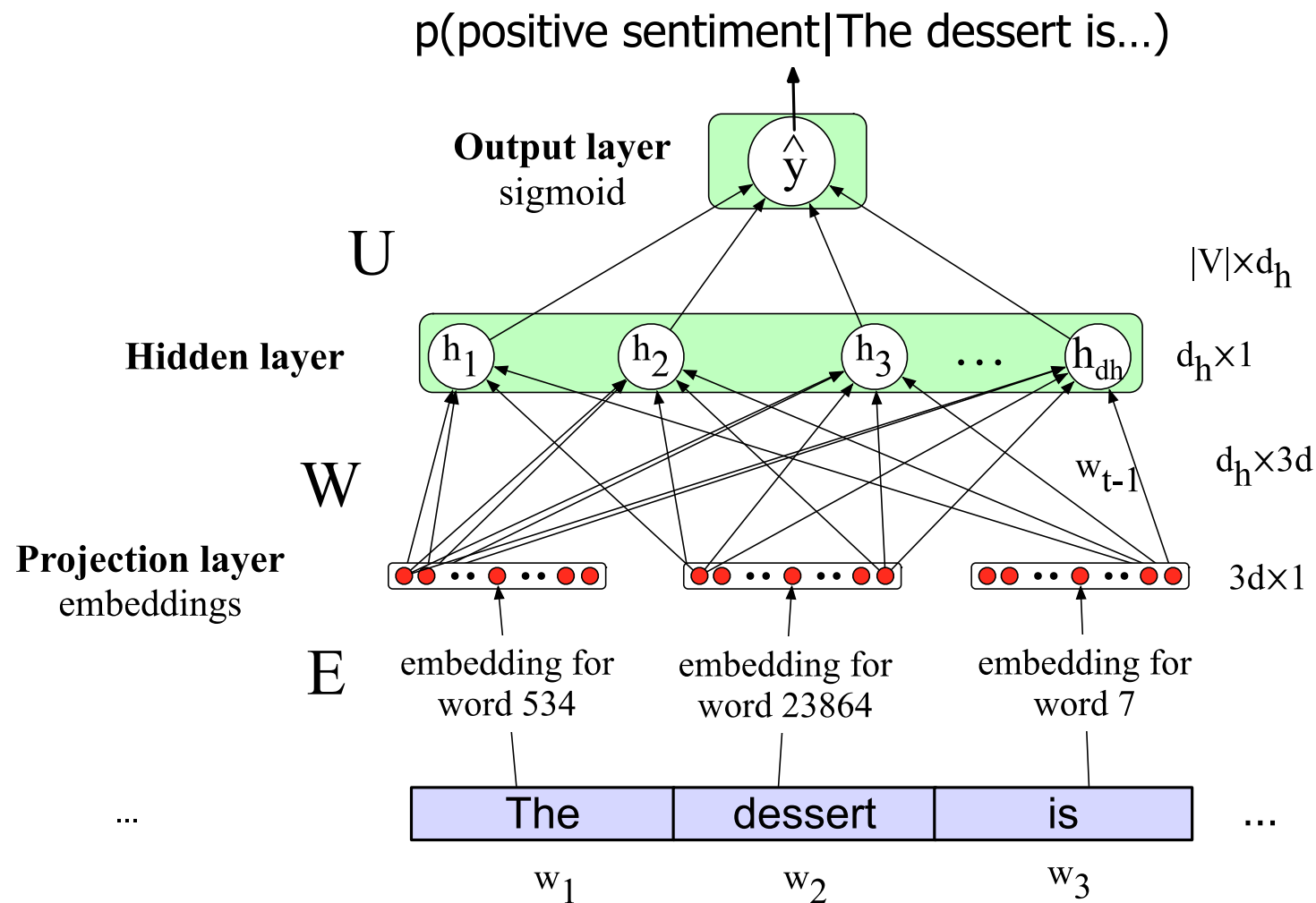# Feedforward nets for simple classification



- Just adding a hidden layer to logistic regression

  o allows the network to use non-linear interactions between features

  o which may (or may not) improve performance.

- The real power of deep learning comes from the ability to **learn** features from the data

- Instead of using hand-built human-engineered features for classification
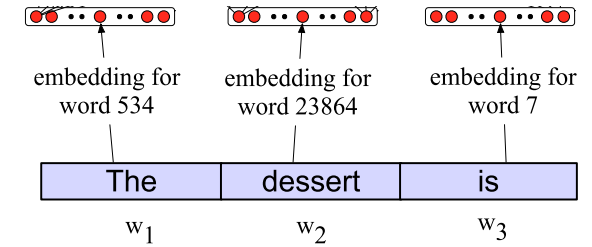
- Use learned representations like embeddings!

$$U$$

$$\sigma$$

$$W$$

$$x_1 \quad\quad\quad\quad x_n$$

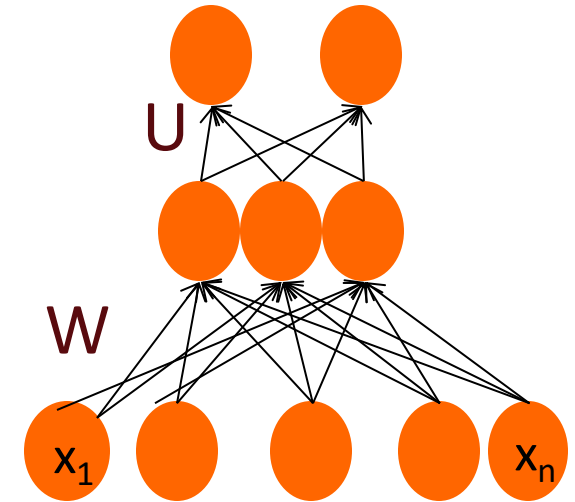$$e_1 \quad e_2 \quad\quad\quad e_n$$

- This assumes a fixed size length (3)!
- Kind of unrealistic.
- Some simple solutions (more sophisticated solutions later)
1. Make the input the length of the longest review
   - If shorter then pad with zero embeddings
   - Truncate if you get longer reviews at test time
2. Create a single "sentence embedding" (the same dimensionality as a word) to represent all the words
   - Take the mean of all the word embeddings
   - Take the element-wise max of all the word embeddings
     - For each dimension, pick the max value from all words

embedding for word 534    embedding for word 23864    embedding for word 7

| The | dessert | is |

$w_1$        $w_2$        $w_3$

- # What if you have more than two output classes?

  o Add more output units (one for each class)

  o And use a "softmax layer"

$$\text{softmax}(z_i) \;=\; \frac{e^{z_i}}{\sum_{j=1}^{k} e^{z_j}} \quad 1 \leq i \leq D$$

# Neural Language Models (LMs)

- **Language Modeling**: Calculating the probability of the next word in a sequence given some history.
  - ○ We've seen N-gram based LMs
  - ○ But neural network LMs far outperform n-gram language models
- State-of-the-art neural LMs are based on more powerful neural network technology like Transformers
- But **simple feedforward LMs** can do almost as well!

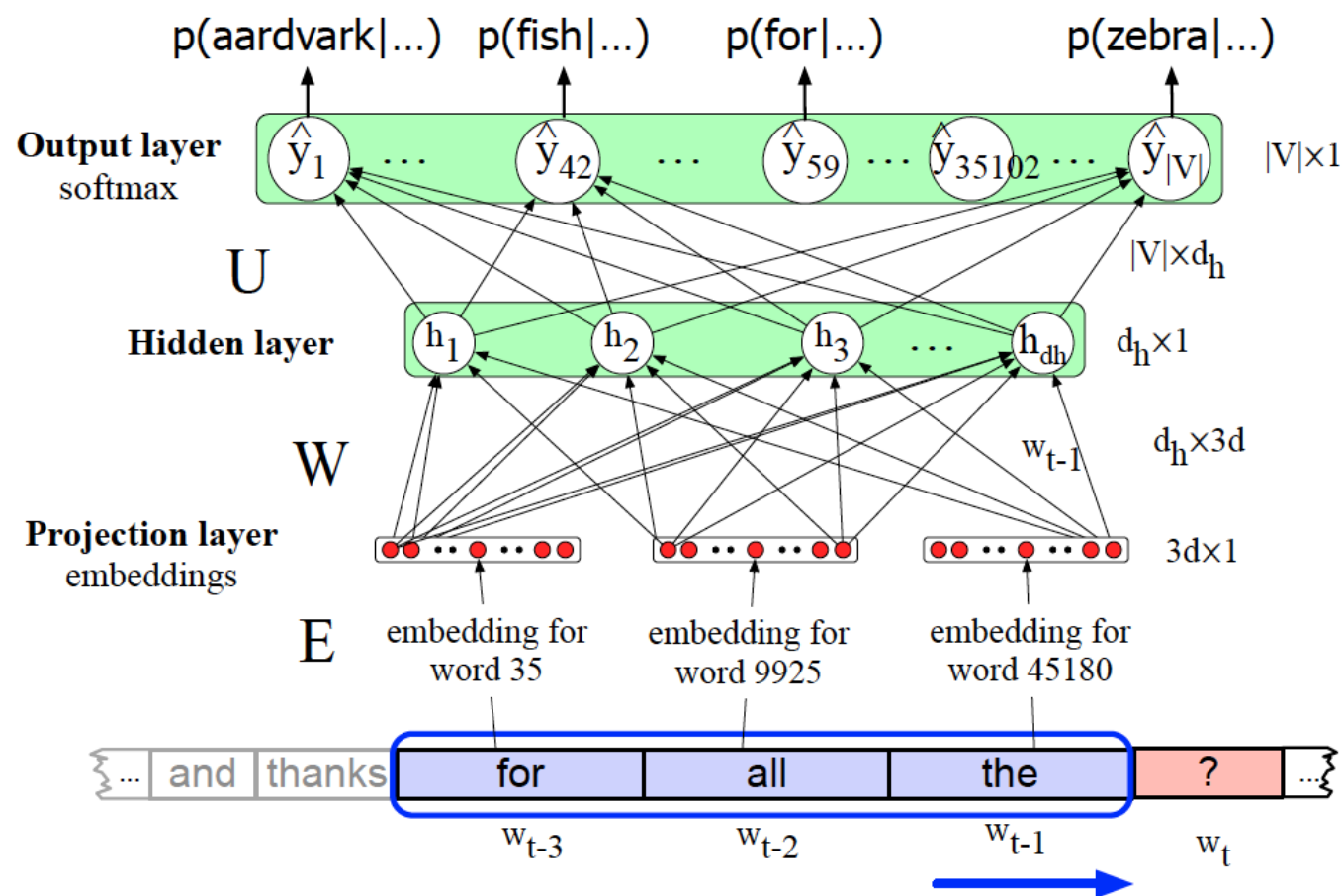**Task**: predict next word $w_t$

given prior words $w_{t-1}$, $w_{t-2}$, $w_{t-3}$, …

**Problem**: Now we're dealing with sequences of arbitrary length.

**Solution**: Sliding windows (of fixed length)

$$P(w_t|w_1^{t-1}) \approx P(w_t|w_{t-N+1}^{t-1})$$

# Why Neural LMs work better than N-gram LMs

- **Training data:**

  ○ We've seen:  I have to make sure that the cat gets fed.

  ○ Never seen:  dog gets fed

- **Test data:**

  ○ I forgot to make sure that the dog gets ___

- N-gram LM can't predict "fed"!

- Neural LM can use similarity of "cat" and "dog" embeddings to generalize and predict "fed" after dog

# SIT330-770: Natural Language Processing

## Week 6 - Neural Networks and Neural LMs

**Dr. Mohamed Reda Bouadjenek**

School of Information Technology, Faculty of Sci Eng & Built Env

reda.bouadjenek@deakin.edu.au
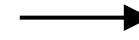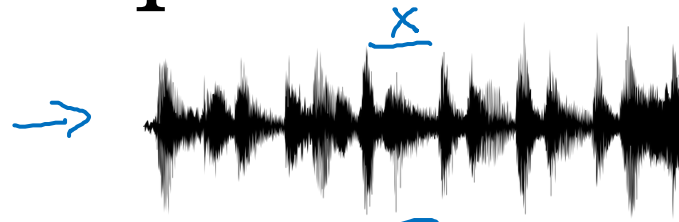
Andrew Ng

Neural Networks and Deep Learning

# Recurrent Neural Networks

## Why sequence models?

deeplearning.ai

# Examples of sequence data

Speech recognition → [audio waveform $x$] → $y$ "The quick brown fox jumped over the lazy dog."

Music generation ∅ → [musical notation]

Sentiment classification "There is nothing to like in this movie." → ★☆☆☆☆

DNA sequence analysis → AGCCCCTGTGAGGAACTAG → AG CCCCTGTGAGGAACT AG

Machine translation Voulez-vous chanter avec moi? → Do you want to sing with me?

Video activity recognition [sequence of images] → Running

Name entity recognition → Yesterday, Harry Potter met Hermione Granger. → Yesterday, Harry Potter met Hermione Granger.

Andrew Ng

# Recurrent Neural Networks

---

# Notation

# Motivating example

NLP

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

$\rightarrow$   $x^{\langle 1 \rangle}$   $x^{\langle 2 \rangle}$   $x^{\langle 3 \rangle}$   $- - \cdots$   $x^{\langle t \rangle}$   $- - \cdots$   $x^{\langle 9 \rangle}$

$T_x = 9$

$\rightarrow$ y:   1   1   0   1   1   0   0   0   0

$y^{\langle 1 \rangle}$   $y^{\langle 2 \rangle}$   $y^{\langle 3 \rangle}$   $- - \cdots$   $y^{\langle 9 \rangle}$

$T_y = 9$

$x^{(i)\langle t \rangle}$   $T_x^{(i)} = 9$   15

$y^{(i)\langle t \rangle}$   $T_y^{(i)}$

# Representing words

$x^{<t>}$          $(x, y)$

$x \longrightarrow y$

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

$x^{<1>}$  $x^{<2>}$  $x^{<3>}$         ...         $x^{<9>}$

$x^{<t>}$

Vocabulary

| | |
|---|---|
| a | 1 ← |
| aaron | 2 |
| ... | ... |
| and | 367 ← |
| ... | ... |
| harry | 4075 |
| ... | ... |
| potter | 6830 |
| ... | ... |
| zulu | 10,000 |

<UNK>  10,000

$\in 4075$

$\in 6830$

$\in 367$

10,000

One-hot

Andrew Ng

# Representing words

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

$x^{<1>}$    $x^{<2>}$    $x^{<3>}$                    ...                    $x^{<9>}$

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

# Recurrent Neural Networks
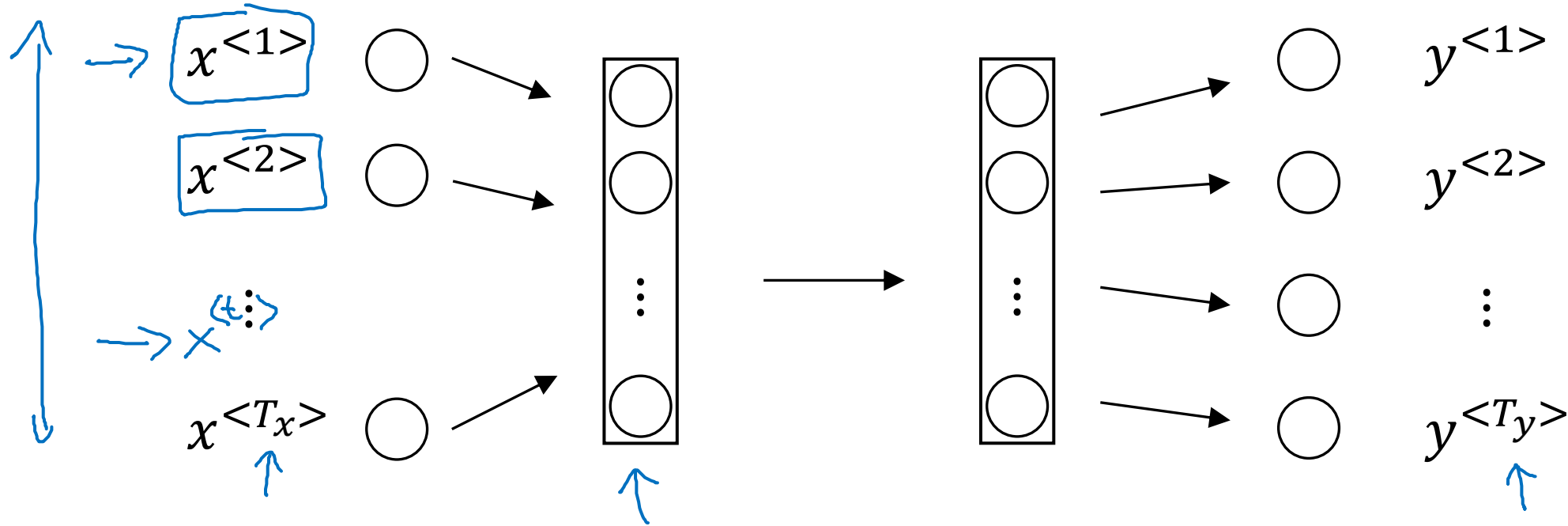
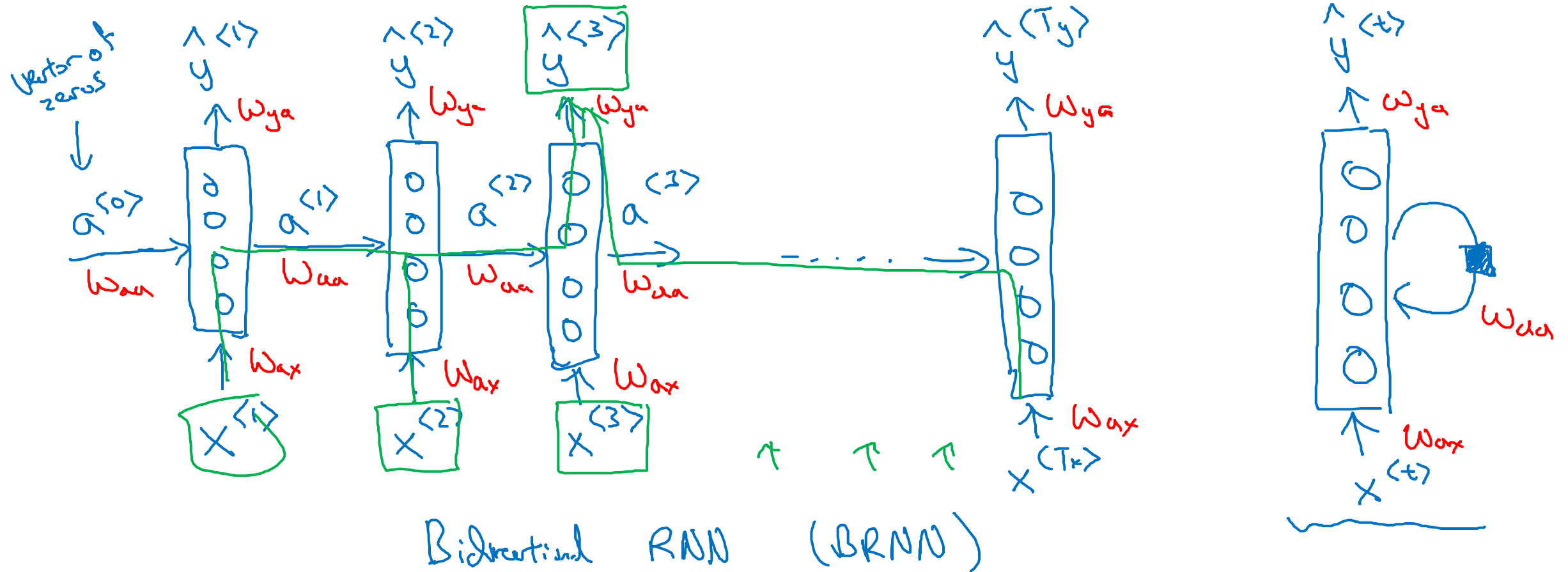deeplearning.ai

---

# Recurrent Neural Network Model

# Why not a standard network?



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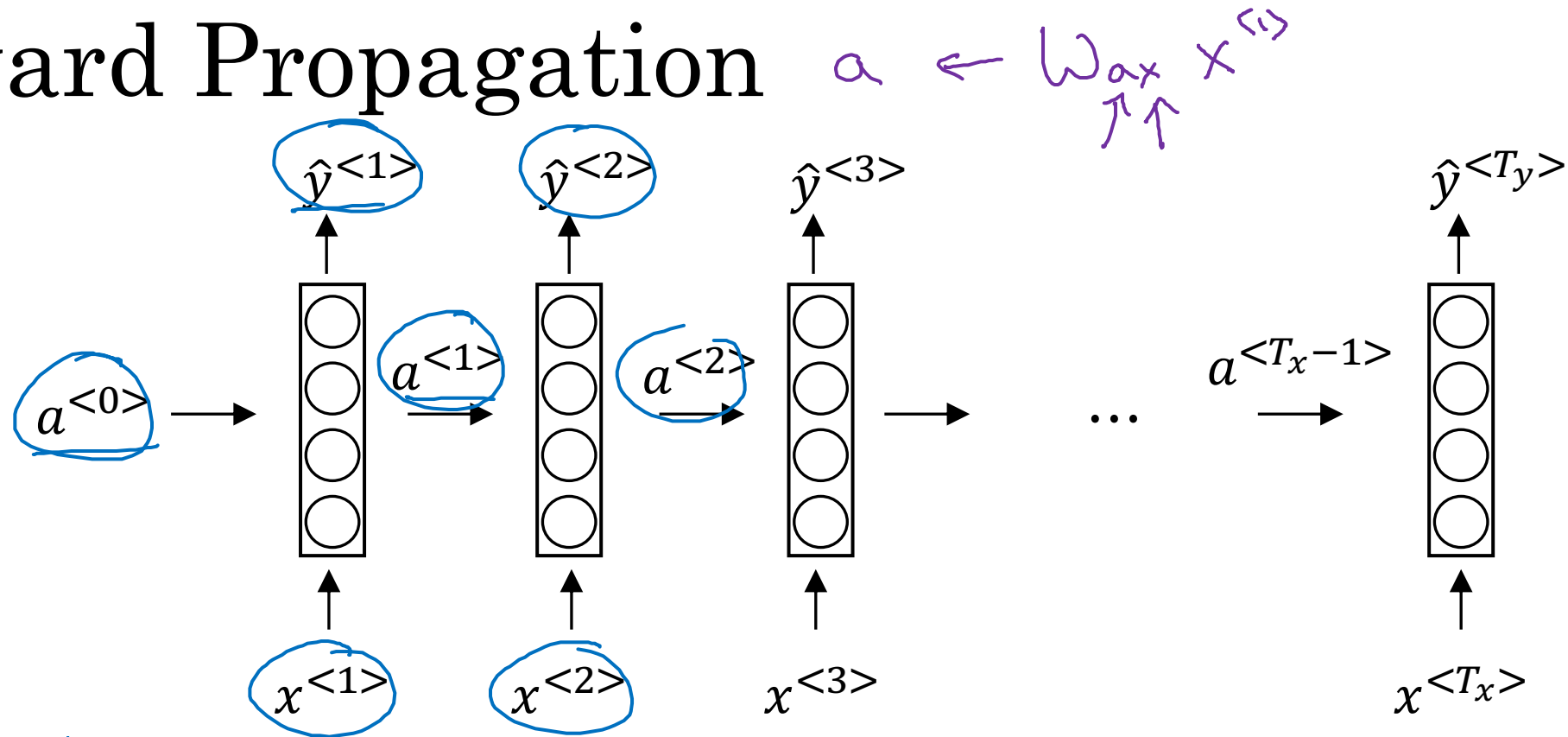- Doesn't share features learned across different positions of text.

Andrew Ng

# Recurrent Neural Networks



Bidirectional RNN (BRNN)

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

He said, "Teddy bears are on sale!"

Andrew Ng

# Forward Propagation

$a \leftarrow W_{ax} x^{(1)}$



$\hat{y}^{<1>}$  $\hat{y}^{<2>}$  $\hat{y}^{<3>}$  $\hat{y}^{<T_y>}$

$a^{<0>}$  $a^{<1>}$  $a^{<2>}$  $a^{<T_x-1>}$  ...

$x^{<1>}$  $x^{<2>}$  $x^{<3>}$  $x^{<T_x>}$

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

$a^{<1>} = g_1(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a) \quad \leftarrow \quad \tanh / \text{ReLU}$

$\hat{y}^{<1>} = g_2(W_{ya} a^{<1>} + b_y) \quad \leftarrow \quad \text{Sigmoid}$

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

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

# Simplified RNN notation

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

(100,100)   100   (100,10,000)   10,000

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

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

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

$$100 \uparrow [W_{aa} \vdots W_{ax}] = W_a$$
$$\underset{100}{\longleftrightarrow} \quad \underset{10\,000}{\longleftrightarrow}$$

(100, 10100)

$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} \begin{matrix} \updownarrow 100 \\ \updownarrow 10000 \end{matrix} \quad \updownarrow 10100$$

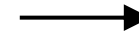$$[W_{aa} \vdots W_{ax}] \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = W_{aa}a^{<t-1>} + W_{ax}x^{<t>}$$

# Recurrent Neural Networks

---

# Different types of RNNs

# Examples of sequence data

$T_x$  $T_y$

$X$  $y$

Speech recognition  ⟶  "The quick brown fox jumped over the lazy dog."

Music generation  ∅  ⟶  (musical notation)

Sentiment classification  "There is nothing to like in this movie."  ⟶  ★☆☆☆☆

DNA sequence analysis  AGCCCCTGTGAGGAACTAG  ⟶  AG**CCCCTGTGAGGAACT**AG

Machine translation  Voulez-vous chanter avec moi?  ⟶  Do you want to sing with me?

Video activity recognition  (images)  ⟶  Running

Name entity recognition  Yesterday, Harry Potter met Hermione Granger.  ⟶  Yesterday, **Harry Potter** met **Hermione Granger**.

Andrew Ng

# Examples of RNN architectures

$T_x = T_y$

Sentiment classification
$x = text$
$y = 0/1 \qquad 1 \cdots 5$



Many-to-many

Many-to-one

one-to-one

# Examples of RNN architectures



Music generation

$x \rightarrow y^{<1>} y^{<2>} \ldots y^{<T_y>}$

One-to-many

$x = \phi$

Machine translation

encoder

decoder

Many-to-many

# Summary of RNN types



One to one

One to many

Many to one

Many to many

$T_x = T_y$

Many to many

Andrew Ng

Recurrent Neural Networks

Language model and sequence generation

deeplearning.ai

# What is language modelling?

Speech recognition

The apple and <u>pair</u> salad.

$\rightarrow$ The apple and <u>pear</u> salad.

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

$P$(The apple and pear salad) = $5.7 \times 10^{-10}$

$P(\text{Sentence}) = ?$ $\qquad P\left(y^{\langle 1 \rangle}, y^{\langle 2 \rangle}, \ldots, y^{\langle T_y \rangle}\right)$

Andrew Ng

# Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

Cats average 15 hours of sleep a day. ↙ &lt;EOS&gt;

$y^{\langle 1 \rangle}$   $y^{\langle 2 \rangle}$   $y^{\langle 3 \rangle}$   $\cdots$   $y^{\langle 8 \rangle}$   $y^{\langle 9 \rangle}$

$x^{\langle t \rangle} = y^{\langle t-1 \rangle}$

The Egyptian M̶a̶u̶ is a bread of cat. &lt;EOS&gt;

&lt;UNK&gt;

10,000

Andrew Ng

# RNN model

$P(a) \; P(aaron) \cdots P(cats) \cdots P(zulu)$
$P(\langle UNK \rangle)$
$P(\langle EOS \rangle)$

$P(\text{average} \mid \text{cats})$

$P( \underline{\hspace{1cm}} \mid \text{"cats average"})$

$P(\langle EOS \rangle \ldots .)$



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

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

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

$$P(y^{<1>}, y^{<2>}, y^{<3>}) \Leftarrow$$
$$= P(y^{<1>}) \, P(y^{<2>} \mid y^{<1>})$$
$$P(y^{<3>} \mid y^{<1>}, y^{<2>})$$

Andrew Ng

# Recurrent Neural Networks

## Sampling novel sequences

deeplearning.ai

# Sampling a sequence from a trained RNN

$P(y^{<1>}, \dots, y^{(T_x)})$

Training:

$\hat{y}^{<1>}$  $\hat{y}^{<2>}$  $\hat{y}^{<3>}$  $\hat{y}^{<T_y>}$

$a^{<0>} \rightarrow \boxed{a^{<1>}} \rightarrow \boxed{a^{<2>}} \rightarrow \boxed{a^{<3>}} \rightarrow \cdots \rightarrow \boxed{a^{<T_y>}}$

$x^{<1>}$  $y^{<1>}$  $y^{<2>}$  $y^{<T_x-1>}$

Sampling:

The $\hat{y}^{<1>}$  $\hat{y}^{<2>}$  $\hat{y}^{<3>}$  $\hat{y}^{<T_y>}$

$a^{<0>} = 0 \rightarrow \boxed{a^{<1>}} \rightarrow \boxed{a^{<2>}} \rightarrow \boxed{a^{<3>}} \cdots \boxed{\phantom{a}}$

$x^{<1>} = 0$  The $x^{<2>} = \hat{y}^{<1>}$  $y^{<T_x-1>}$

$\rightarrow P(a)P(aaron)\dots P(zulu)P(<UNK>)$   n.p.random.choice

$<EOS>$

$<UNK>$

$P(\_\_ | the)$

Andrew Ng

# Character-level language model

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

Vocabulary = [ a, b, c, ..., z, ⌣, ., ,, ;, 0, ..., 9, A, ..., Z]

$y^{<1>} y^{<2>} y^{<3>}$  $y^{<4>}$

Cat    average
↑ ↑ ↑ ↑  ...

May



Andrew Ng

# Sequence generation

## News

## Shakespeare

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.

The mortal moon hath her eclipse in love.

And subject of this thou art another this fold.

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

For whose are ruse of mine eyes heaves.

# Recurrent Neural Networks

---

# Vanishing gradients with RNNs

deeplearning.ai

# Vanishing gradients with RNNs

The (cat), which ate —————, was full

The cats, Wh — ————— ..... were full

$\hat{y}^{<1>}$  $\hat{y}^{<2>}$  $\hat{y}^{<3>}$  $\hat{y}^{<T_y>}$

$a^{<0>}$  $a^{<1>}$  $a^{<2>}$  $a^{<3>}$  ...  $a^{<T_y>}$

$x^{<1>}$  $x^{<2>}$  $x^{<3>}$  $x^{<T_x>}$

$x$  ...  $\hat{y}$

100

Exploding gradients.

NaN    Gradient clipping

# Recurrent Neural Networks

## Gated Recurrent Unit (GRU)

deeplearning.ai

# RNN unit

$y^{<t>}$

Softmax

$a^{<t-1>}$

$a$

tanh

$a^{<t>}$

$x^{<t>}$

tanh

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

Andrew Ng

# GRU (simplified)

$c^{\langle t-1 \rangle}$
$= a^{\langle t-1 \rangle}$

softmax $\rightarrow y^{\langle t \rangle}$

$\rightarrow c^{\langle t \rangle}$
$= a^{\langle t \rangle}$

$\tilde{c}^{\langle t \rangle}$  $\Gamma_u$

tanh  $\sigma$

$x^{\langle t \rangle}$

$\Gamma_u = 1$   $\Gamma_u = 0$  $\Gamma_u = 0$  $\Gamma_u = 0$  ......  $\downarrow = 1$
$c^{\langle t \rangle} = 1$

The cat, which already ate ..., was full.

$C$ = memory cell

$\rightarrow \boxed{C^{\langle t \rangle}} = a^{\langle t \rangle}$

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

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

"update"

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

element-wise

$\Gamma_u = 0.00001$

Gate

[Cho et al., 2014. On the properties of neural machine translation: Encoder-decoder approaches]  $\leftarrow$
[Chung et al., 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling]  $\leftarrow$

Andrew Ng

# Full GRU

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

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

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

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

LSTM

The cat, which ate already, was full.

# Recurrent Neural Networks

deeplearning.ai

## LSTM (long short term memory) unit

# GRU and LSTM

## GRU

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

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

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

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

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

$\Gamma_f$

## LSTM

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

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

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

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

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

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

Andrew Ng

# LSTM in pictures

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

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

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

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

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

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

$c^{(t-1)}$

peephole connection

Andrew Ng

Recurrent Neural Networks

Bidirectional RNN

deeplearning.ai

# Getting information from the future

He said, "Teddy bears are on sale!"

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



RNN
GRU
LSTM

# Bidirectional RNN (BRNN)



$$\hat{y}^{\langle t \rangle} = g(W_y[\overrightarrow{a}^{\langle t \rangle}, \overleftarrow{a}^{\langle t \rangle}] + b_y)$$

Acyclic graph

He said "Teddy Roosevelt ..."

BRNN w/ LSTM

Andrew Ng

deeplearning.ai

# Recurrent Neural Networks

Deep RNNs

# Deep RNN example

$\Rightarrow y^{<1>}$

$\Rightarrow \hat{y}^{(2)}$

$a^{[2]<t>}$

$y^{<1>}$       $y^{<2>}$       $y^{<3>}$       $y^{<4>}$

| $a^{[3]<0>}$ → | $a^{[3]<1>}$ | → | $a^{[3]<2>}$ | → | $a^{[3]<3>}$ | → | $a^{[3]<4>}$ |

| $a^{[2]<0>}$ → | $a^{[2]<1>}$ | → | $a^{[2]<2>}$ | → | $a^{[2]<3>}$ | → | $a^{[2]<4>}$ |

$W_a^{[2]}, b_a^{[2]}$

$W_a^{[1]}, b_a^{[2]}$

| $a^{[1]<0>}$ → | $a^{[1]<1>}$ | → | $a^{[1]<2>}$ | → | $a^{[1]<3>}$ | → | $a^{[1]<4>}$ |

$a^{[0]}$

$W_a^{[1]}, b_a^{[1]}$

$x^{<1>}$       $x^{<2>}$       $x^{<3>}$       $x^{<4>}$

RNN
GRU
LSTM     BRNN

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

Andrew Ng

Sequence to sequence models

Basic models

deeplearning.ai

# Sequence to sequence model

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

Jane   visite   l'Afrique   en   septembre

$\longrightarrow$   Jane   is   visiting   Africa   in   September.

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

[Sutskever et al., 2014. Sequence to sequence learning with neural networks]

[Cho et al., 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation]

Andrew Ng

# Image captioning

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

A cat sitting on a chair



$11 \times 11$
s = 4

$55 \times 55 \times 96$

MAX-POOL

$3 \times 3$
s = 2

$27 \times 27 \times 96$

$5 \times 5$
same

$27 \times 27 \times 256$
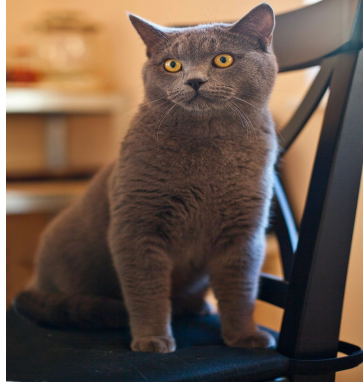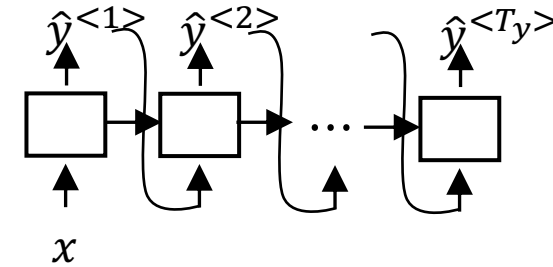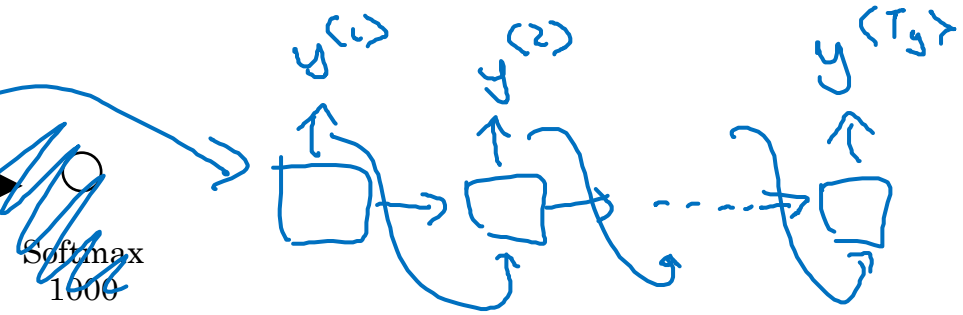
MAX-POOL

$3 \times 3$
s = 2

$13 \times 13 \times 256$

$3 \times 3$
same

$13 \times 13 \times 384$

$3 \times 3$

$13 \times 13 \times 384$

$3 \times 3$

$13 \times 13 \times 256$

MAX-POOL

$3 \times 3$
s = 2

$6 \times 6 \times 256$

= 

9216

4096

4096

Softmax
1000

$y^{<1>}$ $y^{<2>}$ $y^{<T_y>}$

$\hat{y}^{<1>}$ $\hat{y}^{<2>}$ $\hat{y}^{<T_y>}$

...

$x$

[Mao et. al., 2014. Deep captioning with multimodal recurrent neural networks]
[Vinyals et. al., 2014. Show and tell: Neural image caption generator]
[Karpathy and Li, 2015. Deep visual-semantic alignments for generating image descriptions]

Andrew Ng

Sequence to sequence models

Picking the most likely sentence

deeplearning.ai

# Machine translation as building a conditional language model



Language model:

$a^{<0>}$

$\hat{y}^{<1>}$ $\hat{y}^{<2>}$ $\hat{y}^{<T_y>}$

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

$P(y^{(1)}, \ldots, y^{(T_y)})$

Machine translation:

$a^{<0>}$

$x^{<1>}$ $x^{<T_x>}$

$\hat{y}^{<1>}$ $\hat{y}^{<T_y>}$

"Conditional language model"

$P(y^{(1)}, \ldots, y^{(T_y)} \mid x^{(1)}, \ldots, x^{(T_x)})$

Andrew Ng

# Finding the most likely translation

Jane visite l'Afrique en septembre.

English French

$$P(y^{<1>}, \ldots, y^{<T_y>} | x)$$

→ Jane is visiting Africa in September.

→ Jane is going to be visiting Africa in September.

→ In September, Jane will visit Africa.

→ Her African friend welcomed Jane in September.

$$\underset{y^{<1>},\ldots,y^{<T_y>}}{\arg\max} \; P(y^{<1>}, \ldots, y^{<T_y>} | x)$$

# Why not a greedy search?

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



$\hat{y}^{<1>}$ $\qquad$ $\hat{y}^{<T_y>}$

$a^{<0>}$

$x^{<1>}$ $\qquad$ $x^{<T_x>}$

$$\underset{y}{\arg\max} \; P(\hat{y}^{<1>}, \hat{y}^{<2>}, \ldots, \hat{y}^{<T_y>}|x)$$

$10,000$

$10$

$\dfrac{10,000^{10}}{}$

$P(y|x)$

$\longrightarrow$ Jane is visiting Africa in September.

$\longrightarrow$ Jane is going to be visiting Africa in September.

$P(\text{Jane is going}|x) > P(\text{Jane is visiting}|x)$

Andrew Ng

# Beam search algorithm

$B = 3$    (beam width)

## Step 1

$$\rightarrow P(y^{<1>} \mid x)$$

$$\begin{bmatrix} a \\ \vdots \\ in \\ \vdots \\ jane \\ \vdots \\ september \\ \vdots \\ zulu \end{bmatrix}$$

10000



$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \square \rightarrow \hat{y}^{<1>}$

$x^{<1>}$    $x^{<T_x>}$

# Beam search algorithm

$(B = 3)$

Step 1     Step 2

$$\begin{bmatrix} a \\ \vdots \\ in \\ \vdots \\ jane \\ \vdots \\ september \\ \vdots \\ zulu \end{bmatrix}$$

10000

$y^{(1)}, y^{(2)}$

a
aaron
September
visit
zulu

a
aaron
is
visits
zulu

10,000

a
:
zulu

in    $\hat{y}^{(1)}$   $\hat{y}^{(2)}$

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow$    in

$x^{<1>}$        $x^{<T_x>}$

$P(\hat{y}^{(2)} | x, "in")$

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

jane   $\hat{y}^{(2)}$

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow$

$x^{<1>}$        $x^{<T_x>}$

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

September   $\hat{y}^{(2)}$

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow$

$x^{<1>}$        $x^{<T_x>}$

Andrew Ng

# Beam search ($B = 3$)

$B = 1 \implies$ greedy search

in september
- a
- aaron
- jane
- zulu

jane is
- a
- visits
- zulu

jane visits
- a
- africa
- zulu

$P(y^{<1>}, y^{<2>} \mid x)$

jane visits africa in september. <EOS>

$a^{<0>} \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \square$

$x^{<1>} \qquad x^{<T_x>}$

in, september, $\hat{y}^{<3>}$

jane, is, $\hat{y}^{<3>}$

$P(y^{<3>} \mid x, \text{"in september"})$

jane, visits, $\hat{y}^{<3>}$

Andrew Ng

# Sequence to sequence models

---

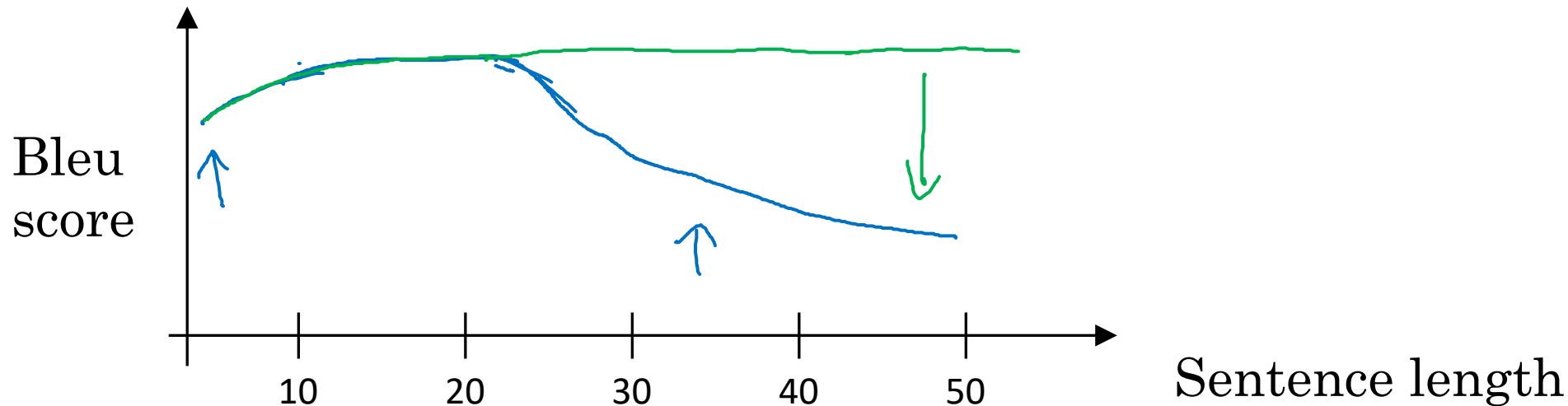# Attention model intuition

deeplearning.ai

# The problem of long sequences



Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

Jane went to Africa last September, and enjoyed the culture and met many wonderful people; she came back raving about how wonderful her trip was, and is tempting me to go too.



Andrew Ng

# Attention model intuition



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

Andrew Ng

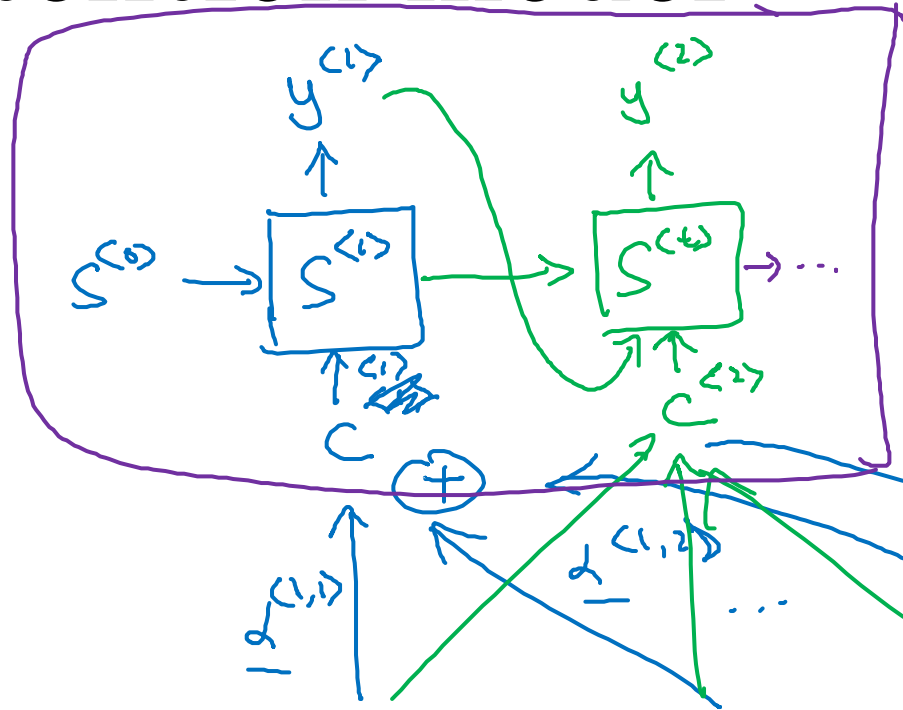Sequence to sequence models

Attention model

deeplearning.ai

# Attention model

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

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

$S^{<0>} \rightarrow S^{<1>} \rightarrow S^{<2>} \rightarrow \cdots$

$C^{<2>} = \sum_{t'} \alpha^{<2,t'>} a^{<t'>}$

$a^{<t'>} = (\overrightarrow{a}^{<t'>}, \overleftarrow{a}^{<t'>})$

$C^{<2>}$

$\alpha^{<1,2>}$    $\alpha^{<1,3>}$

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

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

$\alpha^{<1,1>}$

$C$ $\oplus$

$\overrightarrow{a}^{<0>} \rightarrow$ $\overrightarrow{a}^{<1>}$ $\overleftarrow{a}^{<1>}$    $\overrightarrow{a}^{<2>}$ $\overleftarrow{a}^{<2>}$    $\overrightarrow{a}^{<3>}$ $\overleftarrow{a}^{<3>}$    $\overrightarrow{a}^{<5>}$ $\overleftarrow{a}^{<5>}$ $\leftarrow \overleftarrow{a}^{<6>}$

$t'$

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

jane    visite    l'Afrique    en    septembre

[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

Andrew Ng

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

$T_x$    $T_y$

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

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

$s^{<t-1>}$

$a^{<t'>}$

$e^{<t,t'>}$

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

$\hat{y}^{<t-1>}$   $\hat{y}^{<t>}$

$s^{<t-1>}$   $s^{<t>}$

$a^{<0>} \rightarrow$

$x^{<1>}$    $x^{<2>}$   $\cdots$   $x^{<T_x-1>}$   $x^{<T_x>}$

[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

[Xu et. al., 2015. Show, attend and tell: Neural image caption generation with visual attention]

Andrew Ng

# Attention examples

July 20th 1969 $\longrightarrow$ $1969 - 07 - 20$

23 April, 1564 $\longrightarrow$ $1564 - 04 - 23$

Visualization of $\alpha^{<t,t'>}$: