## Why sequence models

### Examples of sequence data

| | | | |
|---|---|---|---|
| Speech recognition | [audio waveform] $x$ | → | "The quick brown fox jumped over the lazy dog." $y$ |
| Music generation | ∅ | → | [musical notation] |
| Sentiment classification | "There is nothing to like in this movie." | → | ★☆☆☆☆ |
| DNA sequence analysis | AGCCCCTGTGAGGAACTAG | → | AGCCCCTGTGAGGAACTAG |
| Machine translation | Voulez-vous chanter avec moi? | → | Do you want to sing with me? |
| Video activity recognition | [video frames] | → | Running |
| Name entity recognition | Yesterday, Harry Potter met Hermione Granger. | → | Yesterday, Harry Potter met Hermione Granger. |

- Can be addressed as supervised learning problems
- In some, both the input X and the output Y are sequences, and in that case, sometimes X and Y can have different lengths,
- In some of these examples only either X or only the opposite Y is a sequence.

**Named entity recognition:** Identify names, times, locations etc in a sequence (used by search engines)
How would we frame X and y?

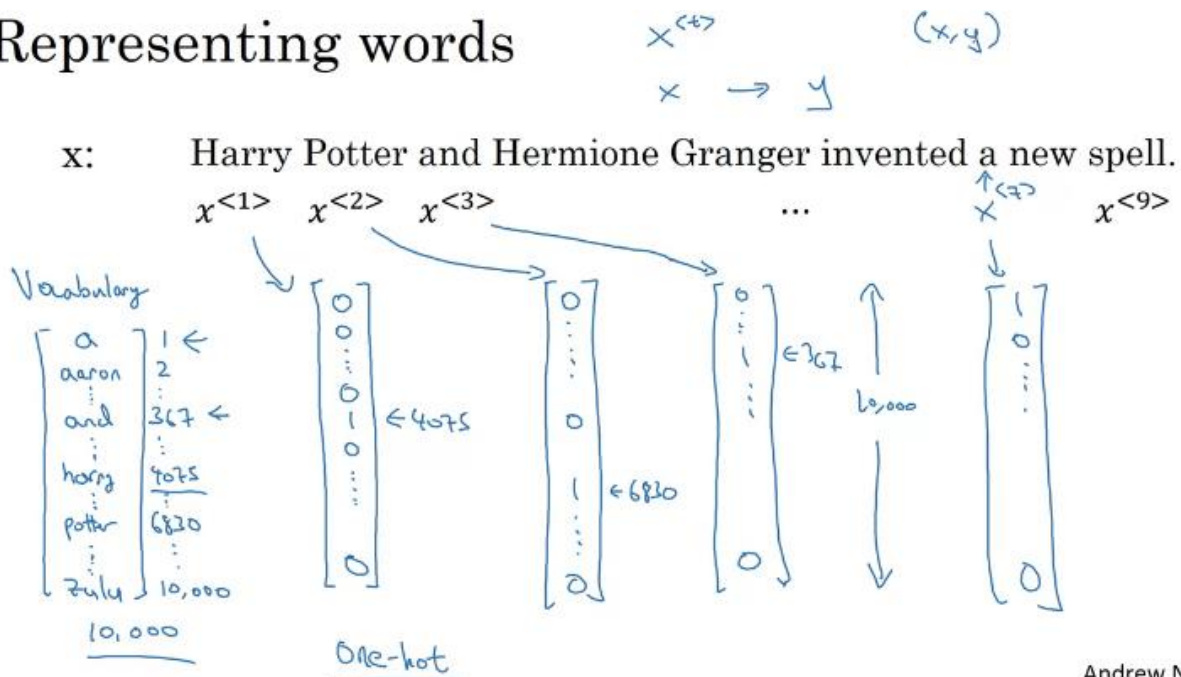x:       Harry Potter and Hermione Granger invented a new spell.

y:       1    1    0    1    1    0    0    0    0

Maybe the above isn't the best output representation, there are some more sophisticated output representations that tells you not just whether a word is part of a person's name, but tells you where are the start and ends of people's names are in the sentence.

Now lets index the position of each word and its label, because these are temporal sequences.
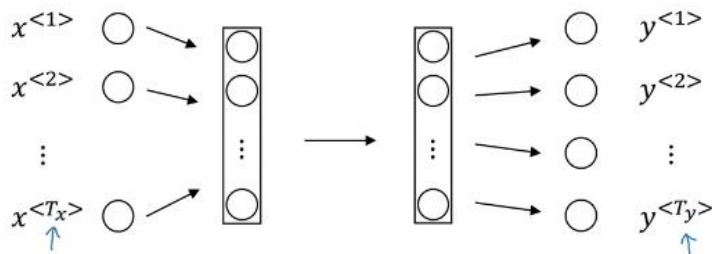
$$x: \quad \underline{\text{Harry Potter}} \text{ and } \underline{\text{Hermione Granger}} \text{ invented a new spell.}$$

$$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \cdots \quad x^{<t>} \quad \cdots \quad x^{<9>}$$

$$T_x = 9$$

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

$$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad \cdots \quad y^{<9>}$$

$$T_y = 9$$

**Notation:**

$x^{(i)}$: $i^{th}$ training example

$x^{(i)<t>}$: $t^{th}$ element of the sequence of training example i

$T_x^{(i)}$: : Length of the input sequence for training example i

$y^{(i)<t>}$: $t^{th}$ element of output sequence of $i^{th}$ training example

$T_y^{(i)}$: Length of output sequence of ith training example

In our above example, $T_x^{(i)}$ is 9 since we have 9 words

<span style="color:red">Representing words</span>

We use a dictionary of vocabulary. e.g. lets consider a dictionary of 10,000 words

We'll use one-hot representation with this vocabulary

## Representing words

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

$$x \rightarrow y$$

$$x: \quad \text{Harry Potter and Hermione Granger invented a new spell.}$$

$$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \cdots \quad x^{<t>} \quad x^{<9>}$$

Vocabulary

$$\begin{bmatrix} a \\ aaron \\ \vdots \\ and \\ \vdots \\ harry \\ potter \\ \vdots \\ zulu \end{bmatrix} \begin{matrix} 1 \leftarrow \\ 2 \\ \\ 367 \leftarrow \\ \\ 4075 \\ 6830 \\ \\ 10,000 \end{matrix}$$

10,000

$$x^{<1>} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 4075 \qquad \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 6830 \qquad \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 367 \quad 10,000 \qquad \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

One-hot

Andrew Ng

# RNN

## Why not a standard network?

$x^{<1>}$ $x^{<2>}$ ⋮ $x^{<T_x>}$ → ... → $y^{<1>}$ $y^{<2>}$ ⋮ $y^{<T_y>}$

## Problems:

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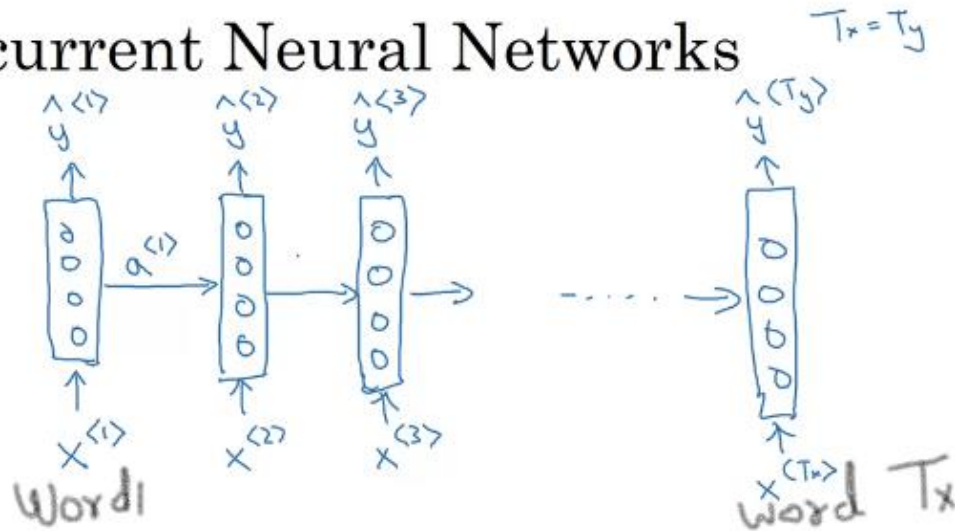- Doesn't share features learned across different positions of text.

Problems:

1. Inputs and outputs can be different lengths for different examples. Maybe we could identify the maximum length of a sentence, and pad other smaller sentences upto that length, but this still wouldn't look like a good representation
2. Doesn't share features learned across different positions of texts. e.g. if the NN figures out that the word Harry appearing in position 1 is a sign that it is a person's name, it would be nice if it could automatically figure out that Harry appearing in another position is also a name. This is similar to CNN where we'd like things learnt in one part of an image to generalize well to other parts of the image as well.

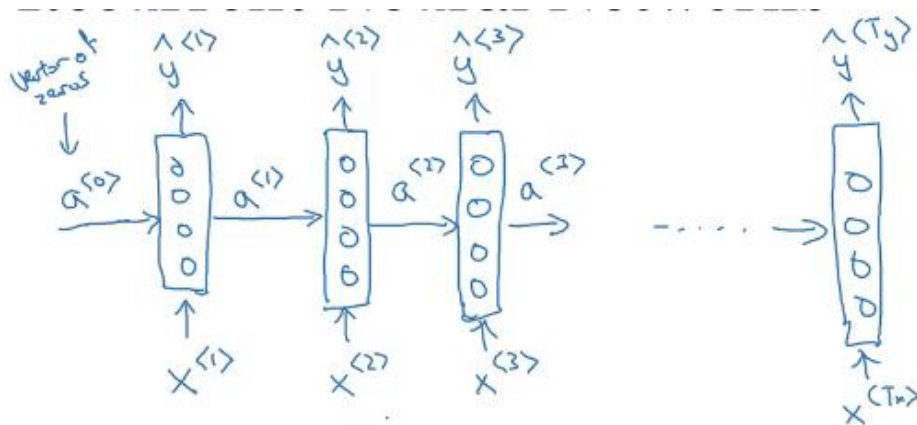**An RNN addresses both these limitations**

What is an RNN?

# Recurrent Neural Networks

$T_x = T_y$



Word1 ... word $T_x$

In the above, we first predict $y^{<1>}$ by feeding the first word $X^{<1>}$ into an NN layer. Then we read the second word $X^{<2>}$, and also use some information from what was computed at 1st step, to predict $y^{<2>}$.

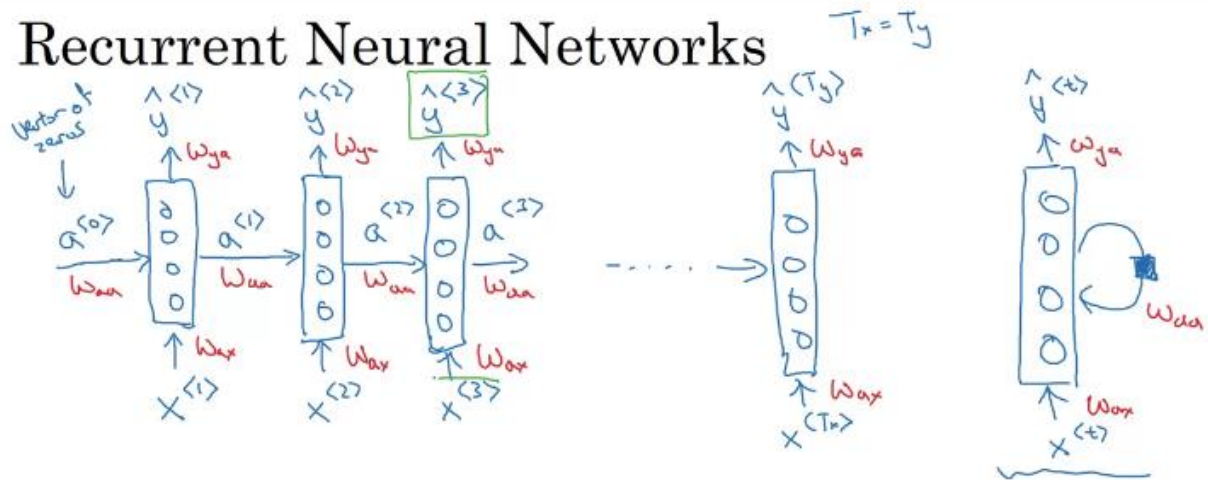In this architecture, $T_x$ and $T_y$ are identical, i.e. input and output sequence lengths are the same.

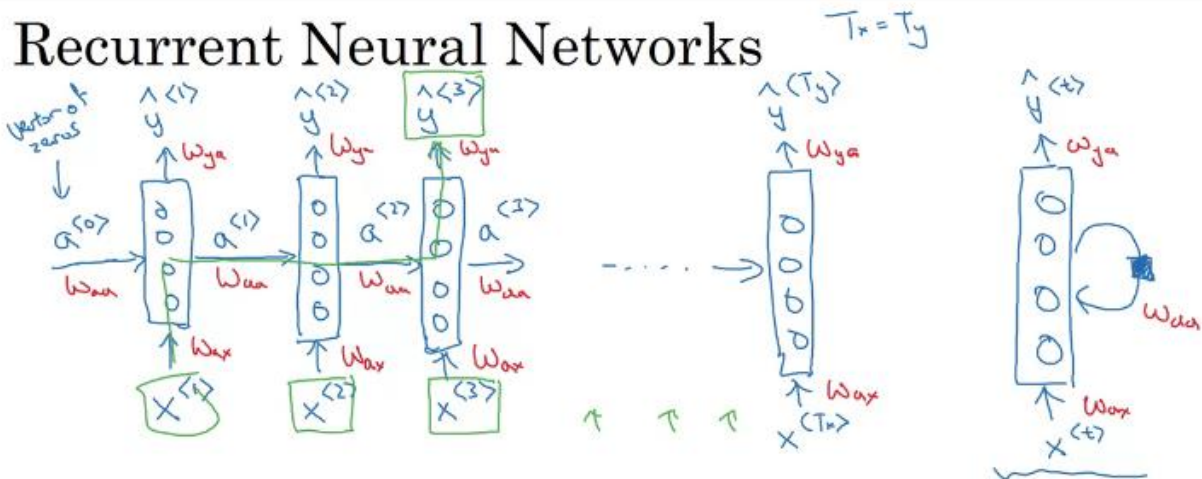We also add a vector of activations at time 0, usually a vector of 0s.



At every time step:

1. $W_{ax}$ (parameters governing connection from X to hidden layer) are identical
2. $W_{aa}$ (parameters governing horizontal connections) are identical
3. $W_{ya}$ (parameters governing output predictions) are identical

# Recurrent Neural Networks



$T_x = T_y$

So prediction $y^{<3>}$ uses info from not only $X^{<3>}$ but also $X^{<1>}$ and $X^{<2>}$

# Recurrent Neural Networks



$T_x = T_y$

Drawback: Doesn't use $X^{<4>}$, $X^{<5>}$ etc i.e. later words in the sentence.

In below sentence, not possibly to decide whether Teddy is part of a person's name by **just** looking at the first three words.
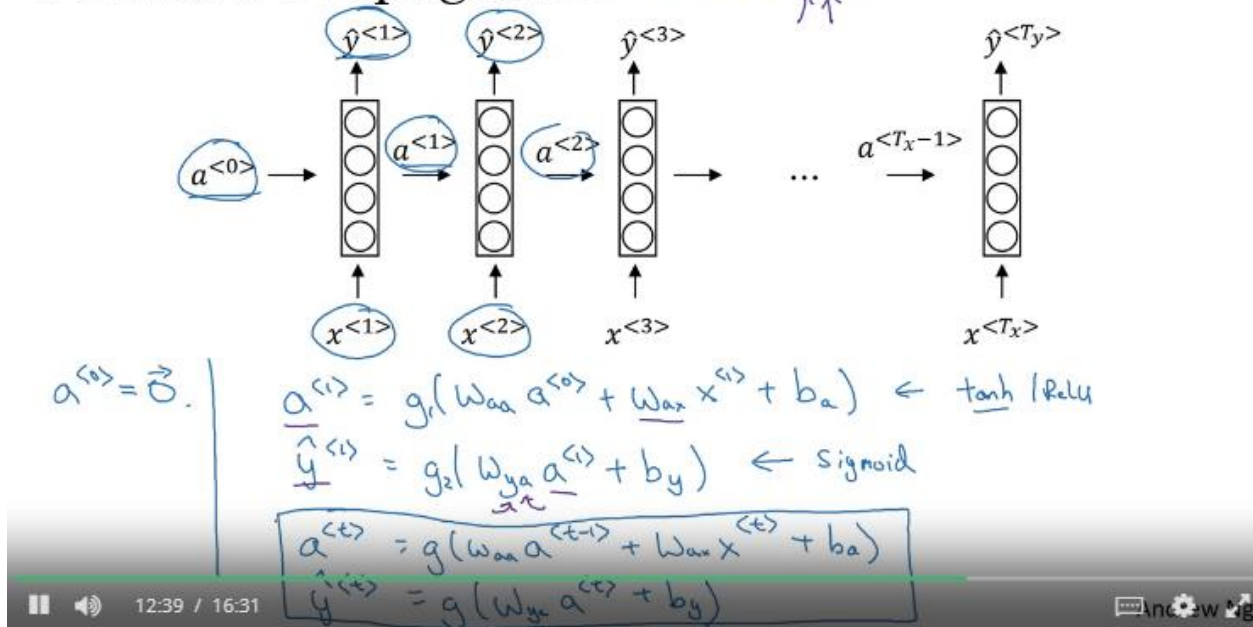
He said, "Teddy Roosevelt was a great President."

He said, "Teddy bears are on sale!"

Later in this course: Bidirectional RNNs (BRNNs) address this

## Forward propagation

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



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

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

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

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

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

⏸ 🔊   12:39 / 16:31                                     Andrew Ng

# Simplified RNN notation

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

$\begin{array}{cc} \uparrow & \uparrow \\ 100 & 10,000 \\ (100,100) & (100,10,000) \end{array}$

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

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

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

$100 \uparrow \left[ W_{aa} \mid W_{ax} \right] = W_a$

$\underset{100}{\longleftrightarrow} \underset{10\,000}{\longleftrightarrow}$   $(100, 10100)$

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

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

See my handwritten notes for the math

## Backpropagation through Time

# Forward propagation and backpropagation



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

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

Backpropagation through time

Andrew Ng

## Types of RNN

# Examples of sequence data



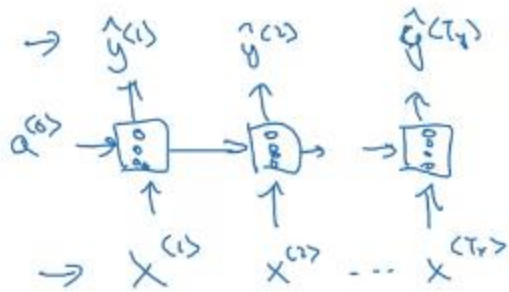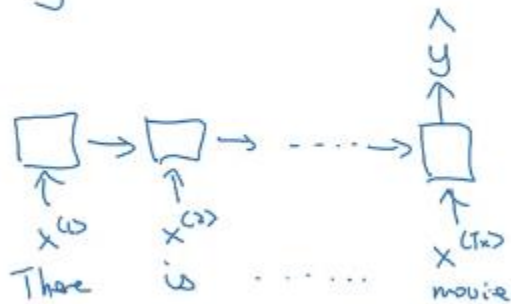| | | |
|---|---|---|
| Speech recognition | (audio) | "The quick brown fox jumped over the lazy dog." |
| Music generation | ∅ | (musical notes) |
| Sentiment classification | "There is nothing to like in this movie." | ★☆☆☆☆ |
| DNA sequence analysis | AGCCCCTGTGAGGAACTAG | AGCCCCTGTGAGGAACTAG |
| 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

Many to many

$$T_x = T_y$$



Many-to-many

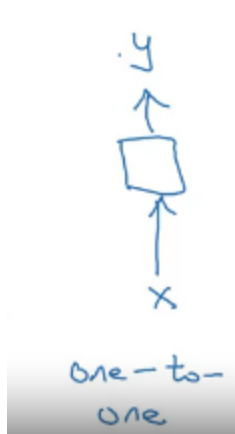## Many to one
Sentiment classification
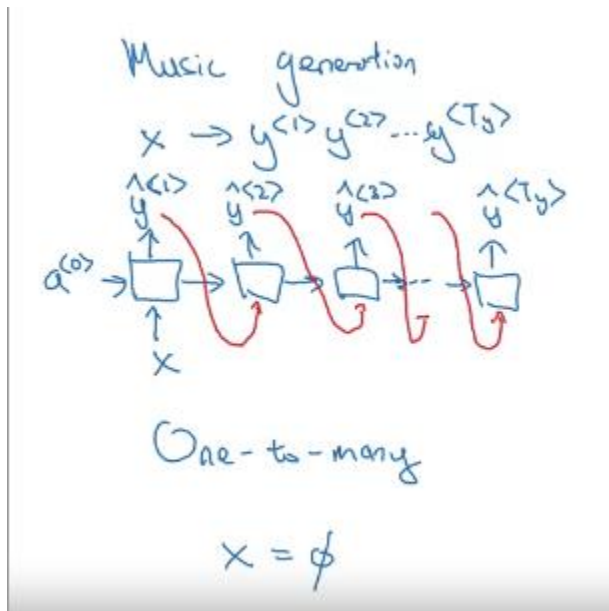
Sentiment classification
x = text
y = 0/1    1...5



There    is    . . . . . . . movie

Many-to-one

## One to one

One to many

# One to many
Music generation



Music generation

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

One-to-many

$$x = \phi$$

# Many to many with different input and output lengths
e.g. machine translation

Machine translation

encoder

decoder

Many-to-many

# Summary of RNN types



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

$a^{<0>}$

$x^{<1>}$

One to one

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

$a^{<0>} \rightarrow$

$x$

One to many

$$\hat{y}$$

$a^{<0>} \rightarrow$

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

Many to one

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

$a^{<0>} \rightarrow$

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

Many to many   $T_x = T_y$

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

$a^{<0>} \rightarrow$

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

Many to many

Andrew Ng

## Language model & sequence generation
Estimates the probability of that particular sequence of words.

# What is language modelling?

Speech recognition

The apple and pair salad.

$\rightarrow$ The apple and pear salad.

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

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

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

First we tokenize the sentence, and one-hot encode

We also add a token for end-of-sentence

# Language modelling with an RNN

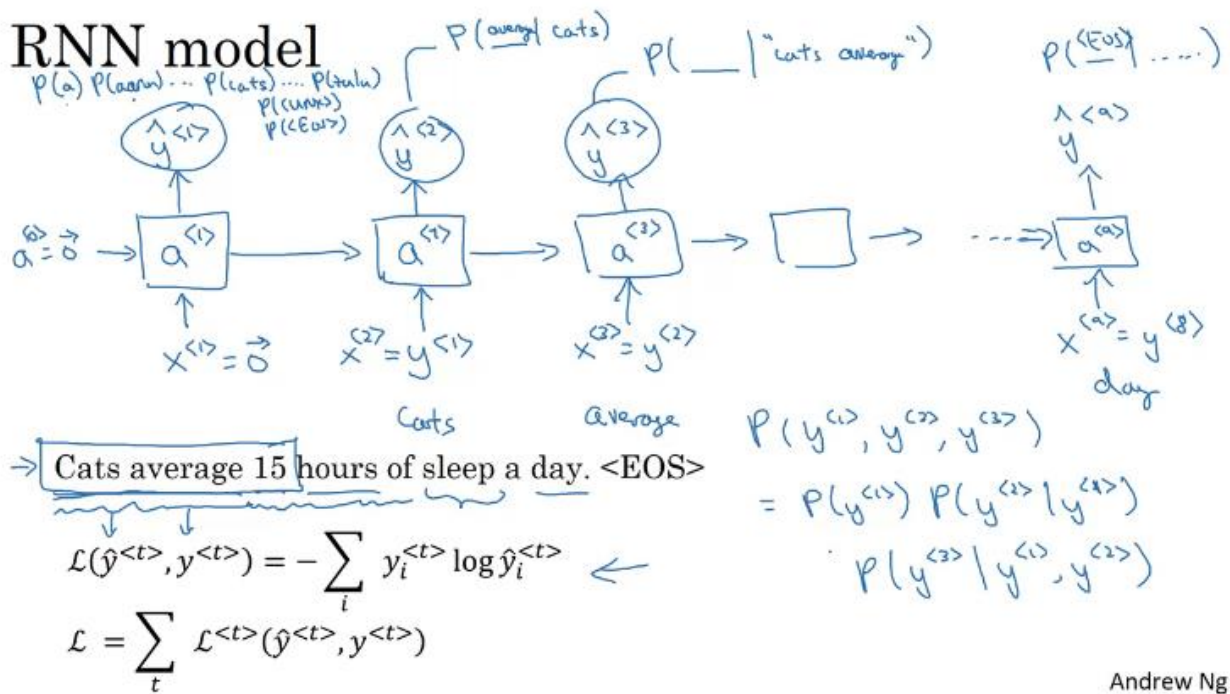Training set: large corpus of english text.

Tokenize

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

$y^{<1>}$   $y^{<2>}$   $y^{<3>}$   $\ldots$   $y^{<8>}$   $y^{<9>}$

The Egyptian Mau is a bread of cat. $<EOS>$

$<UNK>$

10,000

RNN model

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

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

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

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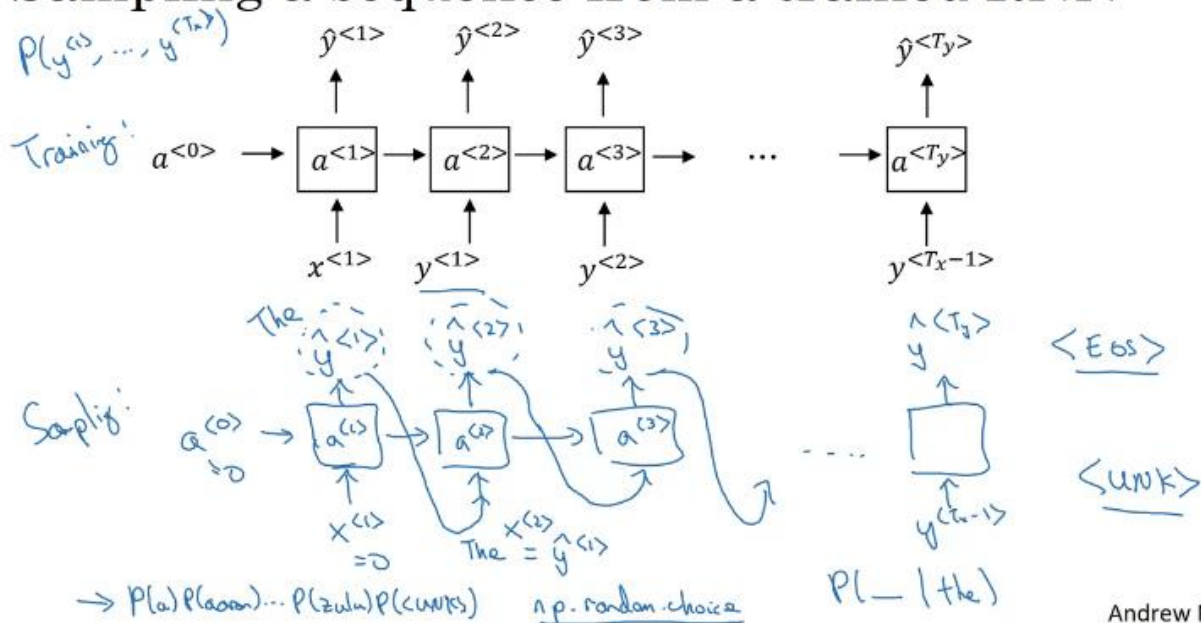$$\cdot P(y^{<3>} | y^{<1>}, y^{<2>})$$

Andrew Ng

## Sampling novel sequences

Sequence models model the chance of any particular sequence of words, so what we like to do is sample from this distribution to generate novel sequences of words.

At the first timestep, we sample from the vocabulary of our words using softmax function.

At the second timestep, we pass in the sampled word from the first timestep as an input i.e. $x^{<2>}$ = yhat$_{<1>}$.

And so on till end of all timesteps.

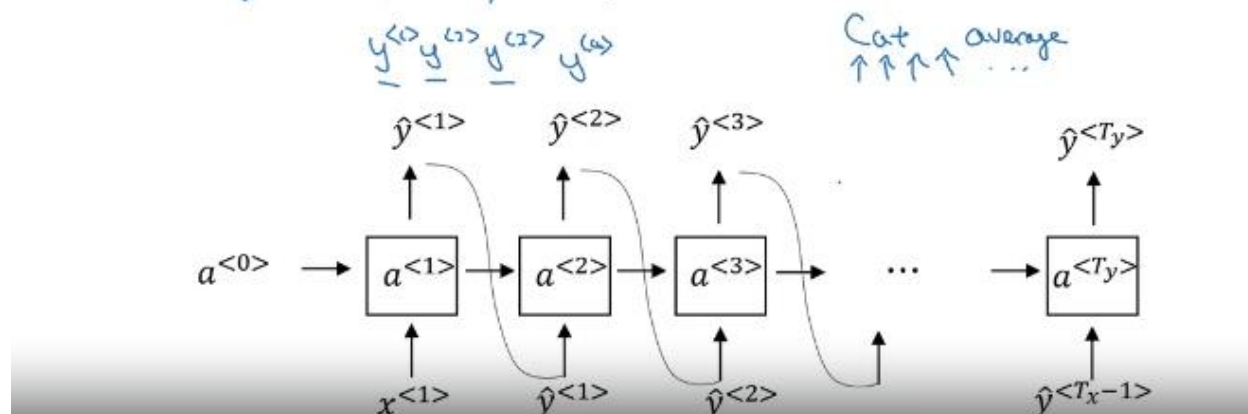

Sampling a sequence from a trained RNN

Andrew Ng

We can also build a character level RNN

# Character-level language model

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

$$Vocabulary = [a, b, c, \cdots, z, \sqcup, ., ,, ;, 0, \cdots, a, A, \cdots, z]$$

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

Cat average
↑↑↑↑ ...

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

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

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

**Pros:**
1. Don't have to worry about unknown word tokens

**Cons**
1. End up with much longer sequences
2. Computationally expensive

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

For whose are ruse of mine eyes heaves.

## Vanishing gradients with RNN

We should match 'cats' with 'were' and 'cat' with 'was'.

But basic RNNs are not good at capturing long term dependency.

The Cat, which oteally ate --------- , was full

The Cats, were full
wh

$\wedge^{<1>}$  $\wedge^{<2>}$  $\wedge^{<3>}$

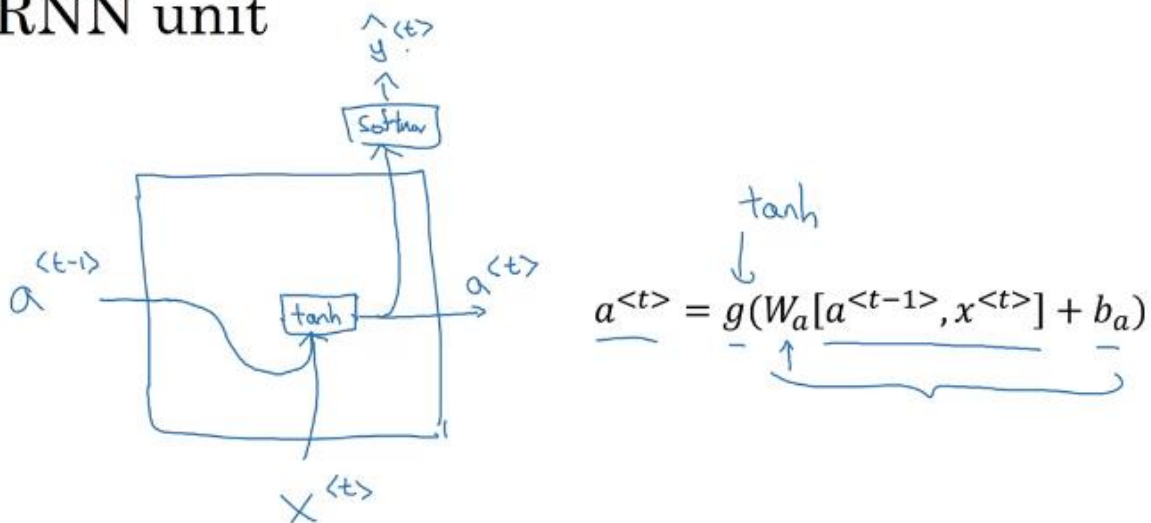Outputs have local influences, ie values closer in the sequence

Recall that for exploding gradients, we use gradient clipping.

## Gated Recurrent Units (GRU)

Helps capture long term dependencies

Below is the visualization of the RNN unit of the hidden layer of the RNN



RNN unit

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

GRU units will have a new variable C which is a memory cell

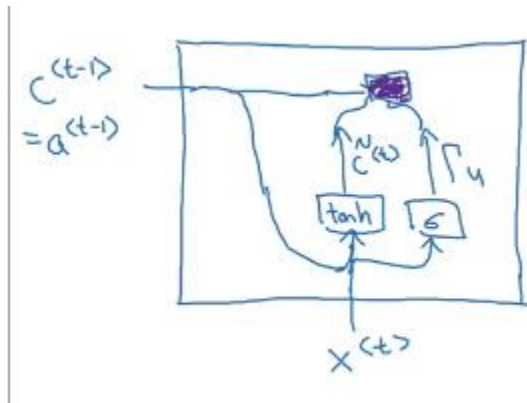At every timestep, $\tilde{C}^{<t>}$ will be a candidate value for replacing the memory cell

$$C = memory\ cell$$

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

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

The important idea of GRU is a gate, $\Gamma_u$ where the u refers to update. Can think of it as always 0 or 1, though computed by sigmoid function.



Gate

The new $C^{<t>}$ will be $\Gamma_u$ times the candidate value, plus $1-\Gamma_u$ times the old value.
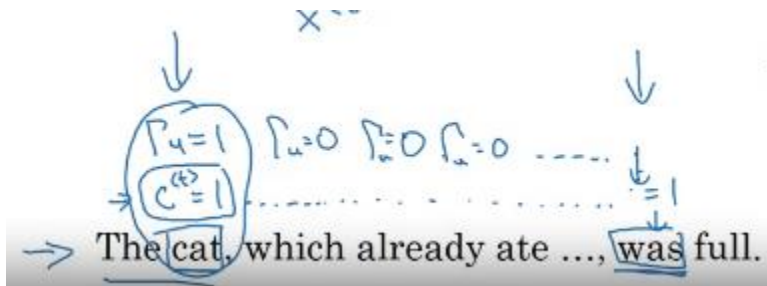
So when $\Gamma_u$ is 1, we update to candidate value. If 0, then don't update, hang on to old value.



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

The purple box above is just

So the word 'cat' is read early in the sentence, so the update gate is set to , and then is not updated at successive. Finally it gets used to decide the word "was"



With tiny value of $\Gamma_u$, the previous value of $C_t$ basically continues to persist without any update.

# Full GRU

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

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

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

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

The cat, which ate already, was full.

LSTM