

SEQUENCE MODELS

COURSE I

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

AGCCCTGTGAGGAACCTAG

AGCCC~~T~~GTGAGGAACCTAG

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.

NOTATION

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

$x^{(i)}$ $x^{(i)}$ $x^{(i)}$ $x^{(i)}$

$x^{(i)}$

y: 1 1 0 1 1 0 0 0 0

$y^{(i)}$ $y^{(i)}$ $y^{(i)}$ $y^{(i)}$ $y^{(i)}$

LENGTH OF INPUT SEQUENCE

$T_x = 9$

- $x^{(i)}$ → t^{th} element in the input sequence of training example i.

- T_x → length of a sequence

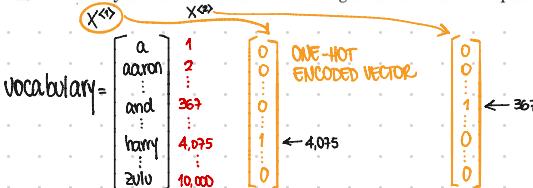
- $T_x^{(i)}$ → input sequence length for training example i.

- $y^{(i)}$ → t^{th} element in the output sequence of training example i.

- $T_y^{(i)}$ → output sequence length for training example i.

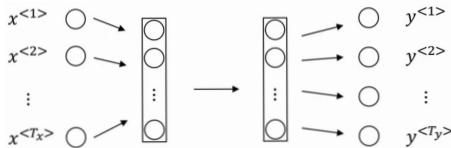
* for representing words you first need to come up with an array of all words that will be used in your model.

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



RECURRENT NEURAL NETWORK

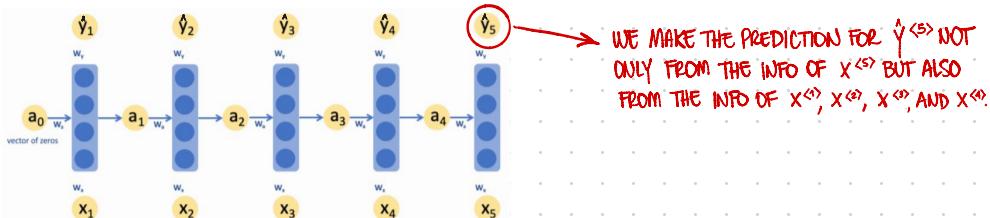
* why not a standard NN?



#1. inputs and outputs can have \oplus lengths in \oplus examples.

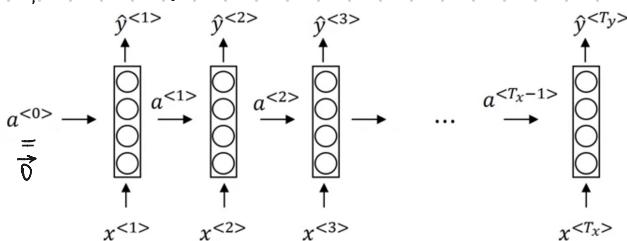
#2. doesn't share features learned across \oplus positions in the text.

* RNN architecture:



* weakness of this NN → it only uses past info to make a prediction.

* forward propagation:



$$a^{<t>} = g(w_a a^{<t-1>} + w_x x^{<t>} + b_a) \rightarrow \text{activation functions: } \underline{\text{tanh or ReLU}}$$

$$\hat{y}^{<t>} = g(w_y a^{<t>} + b_y) \rightarrow \begin{aligned} &\text{sigmoid (binary)} \\ &\text{softmax (multiclass)} \end{aligned}$$

$$a^{<t>} = g(w_a a^{<t-1>} + w_x x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(w_y a^{<t-1>} + b_y)$$

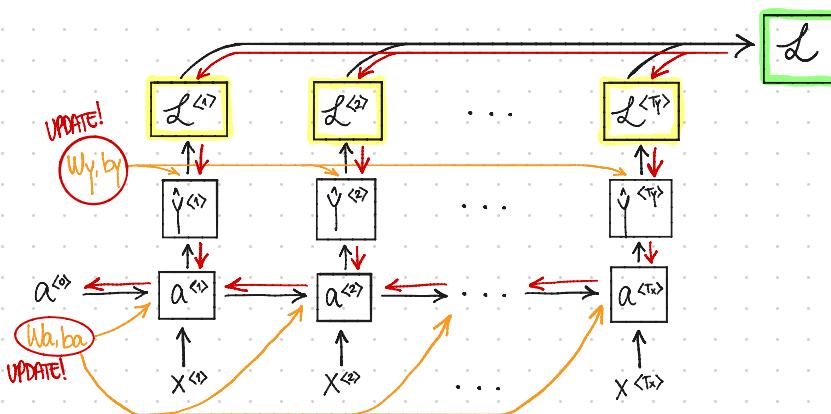
$$a^{(t)} = g(W_a [a^{(t-1)}, x^{(t)}] + b_a)$$

$[W_a; W_{ax}] \rightarrow [W_{aa} \underset{100}{\underset{\swarrow}{|}} \underset{10,000}{\underset{\searrow}{|}} | W_{ax} \underset{(100 \times 100)}{\underset{\swarrow}{|}} \underset{(10,000)}{\underset{\searrow}{|}}$

$[a^{(t-1)}, x^{(t)}] = \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} \underset{100}{\underset{\swarrow}{|}} \underset{10,000}{\underset{\searrow}{|}} \underset{10,100}{\underset{\uparrow}{|}}$

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

FORWARD AND BACKPROP THROUGH TIME



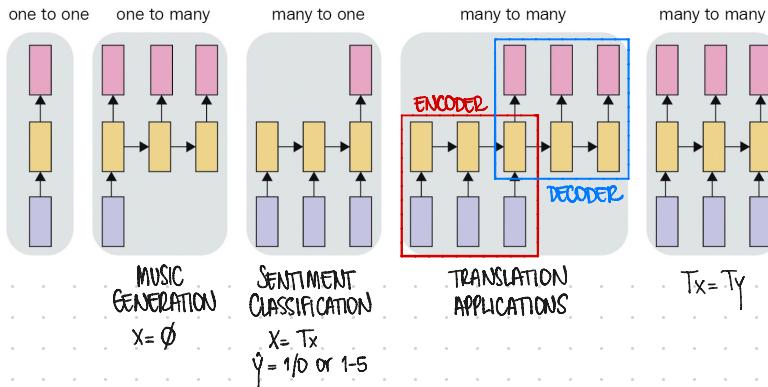
$$\mathcal{L}^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1-y^{(t)}) \log (1-\hat{y}^{(t)})$$

LOSS OF A SINGLE PREDICTION

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

LOSS OF ENTIRE SEQUENCE

DIFFERENT TYPES OF RNNs



LANGUAGE MODEL AND SEQUENCE GENERATION

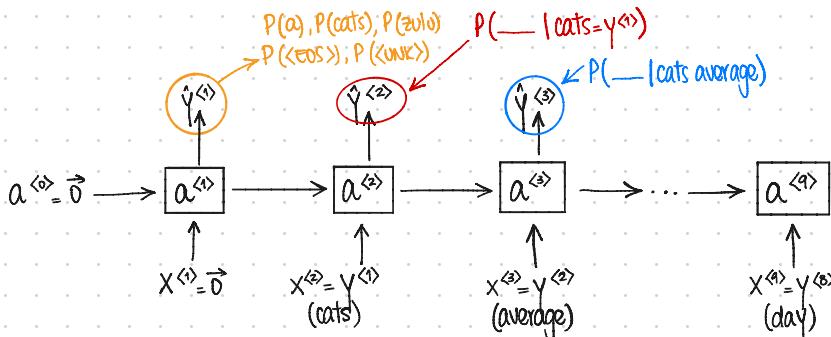
what a language model does is it tells you what the probability of a sentence is:

it inputs a sentence and estimates its probability

$$P(y^{(1)}, y^{(2)}, \dots, y^{(T)})$$

- #1 get a large corpus of text (training set)
- #2 tokenize the text
- #3 map each token to a one-hot vector
- #4 build the RNN model

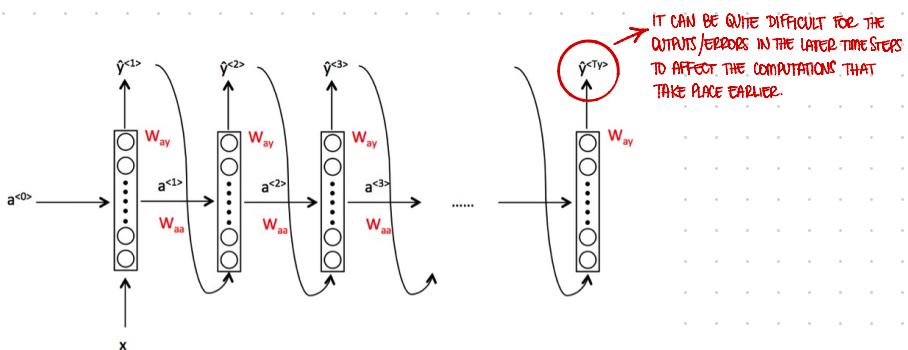
Cats average 15 hours of sleep a day.



VANISHING GRADIENTS WITH RNNs

* language can have very long-term dependencies, where a word can affect what needs to come much later in a sentence.

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



* it is difficult to get the (NN) to realize that it needs to memorize or to just see a singular noun or plural noun, and then later in the sequence generate the corresponding verb form that goes with the noun it previously saw.

GATED RECURRENT UNIT (GRU)

* it has a unit \mathbb{C} called memory cell

it remembers whether "cat" is Singular or plural.

- $C^{<t>} = a^{<t>} \rightarrow C^{<t>} \text{ will output an activation value}$
- at every time step, we will consider over-writing the memory cell with a value $\tilde{C}^{<t>}.$
- $\tilde{C}^{<t>} = \tanh(W_C [C^{<t-1>}, x^{<t>}] + b_C) \rightarrow \tilde{C}^{<t>} \text{ will be computed using a tanh activation function}$
- $\Gamma_u = \sigma(W_u [C^{<t-1>}, x^{<t>}] + b_u) \rightarrow \text{gate of the GRU (it's either 0 or 1)}$

UPDATE

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

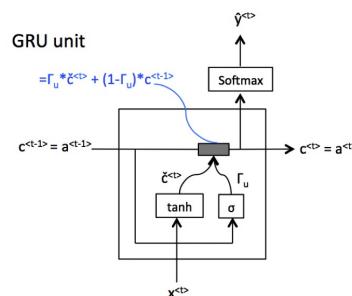
The Cat, which already ate , was full.

$$C^{<t>} = 1$$

$$\Gamma_u = 1$$

$$\Gamma_u = 0 \quad \tilde{C}^{<t>} = 0$$

you memorize "cat" was singular



* $C^{<t>}$, $\tilde{C}^{<t>}$, Γ_u will all have the same dimension

* for a full GRU you add Γ_f which tells you how relevant $C^{<t-1>}$ is to compute the next candidate for $C^{<t>}$.

$$\begin{aligned}\tilde{C}^{<t>} &= \tanh(W_c [\Gamma_f * C^{<t-1>} + x^{<t>}] + b_c) \\ \Gamma_u &= \sigma(W_u [C^{<t-1>} + x^{<t>}] + b_u) \\ \Gamma_r &= \sigma(W_r [C^{<t-1>} + x^{<t>}] + b_r) \\ C^{<t>} &= \Gamma_u * \tilde{C}^{<t>} + (1 - \Gamma_u) * C^{<t-1>}\end{aligned}$$

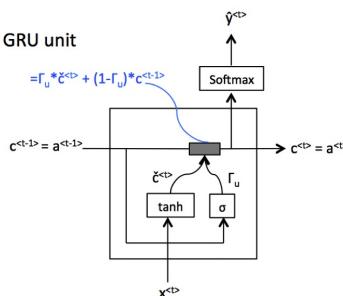
FULL GRU

LONG SHORT TERM MEMORY

* more powerful than GRU, and more general, too.

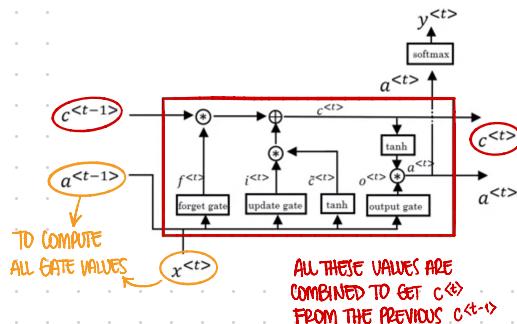
GRU

$$\begin{aligned}\tilde{C}^{<t>} &= \tanh(W_c [\Gamma_f * C^{<t-1>} + x^{<t>}] + b_c) \\ \Gamma_u &= \sigma(W_u [C^{<t-1>} + x^{<t>}] + b_u) \\ \Gamma_r &= \sigma(W_r [C^{<t-1>} + x^{<t>}] + b_r) \\ C^{<t>} &= \Gamma_u * \tilde{C}^{<t>} + (1 - \Gamma_u) * C^{<t-1>} \\ a^{<t>} &= C^{<t>}\end{aligned}$$



LSTM

$$\begin{aligned}\tilde{C}^{<t>} &= \tanh(W_c [a^{<t-1>} + x^{<t>}] + b_c) \\ \Gamma_f &= \sigma(W_f [a^{<t-1>} + x^{<t>}] + b_f) \text{ FORGET GATE} \\ \Gamma_o &= \sigma(W_o [a^{<t-1>} + x^{<t>}] + b_o) \text{ OUTPUT GATE} \\ C^{<t>} &= \Gamma_u * \tilde{C}^{<t>} + \Gamma_f * C^{<t-1>} \\ a^{<t>} &= \Gamma_o * \tanh(C^{<t>})\end{aligned}$$



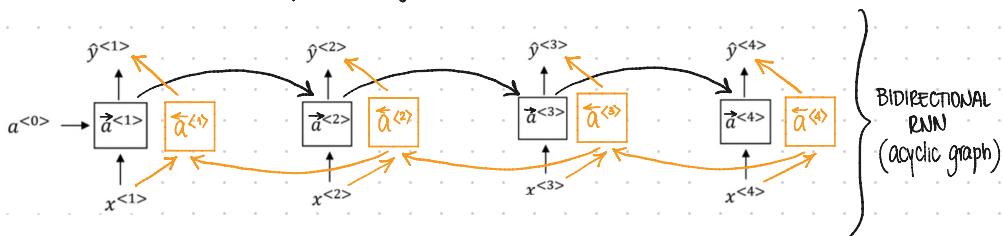
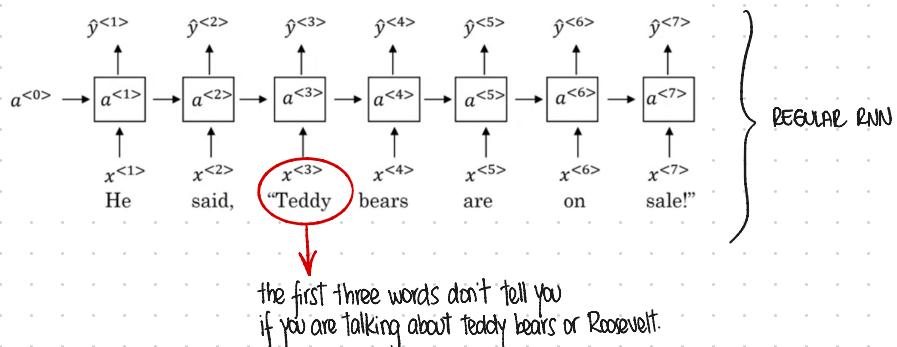
* when to use each? → GRUs are simpler and easier to build a much bigger NN. Computationally, they even run faster.

LSTMs are more powerful and flexible. Historically, more proven choice.

BIDIRECTIONAL RNN

He said, "Teddy bears are on sale!"

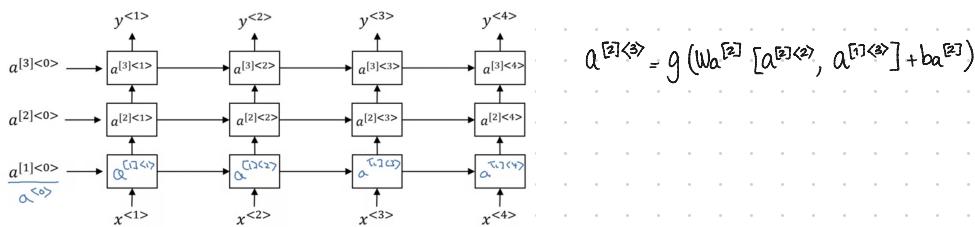
He said, "Teddy Roosevelt was a great President!"



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

* disadvantage → you need the full sequence of data before you start making predictions.

DEEP RNNs



WORD REPRESENTATION

* a way of representing words.

* we represent words through a vocabulary vector $|V|=10,000$. And generate a one-hot representation for each of the words in $|V|$.

Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
[0]	[0]	[0]	[0]	[0]	[0]
0	0	0	0	1	0
0	0	0	0	0	0
⋮	⋮	⋮	⋮	⋮	⋮
1	⋮	⋮	⋮	0	⋮
⋮	1	0	1	0	1
0	0	0	0	0	⋮

APPLE AND ORANGE ARE MUCH MORE SIMILAR AND COULD GO WELL WITH OTHER WORDS LIKE "JUICE!"

* weakness of this representation:
it treats each word as a thing unto itself, and it doesn't allow an algorithm to generalize the cross-words.

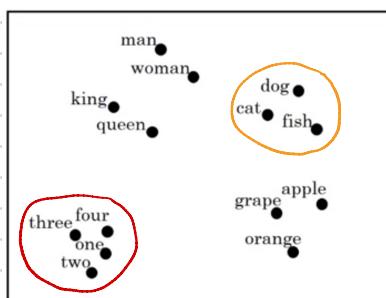
* To solve that issue we use featurized representation where we can learn some features from each of them.

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
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

↑
300 FEATURES = 300-D VECTOR
↓
THE 300-D VECTOR

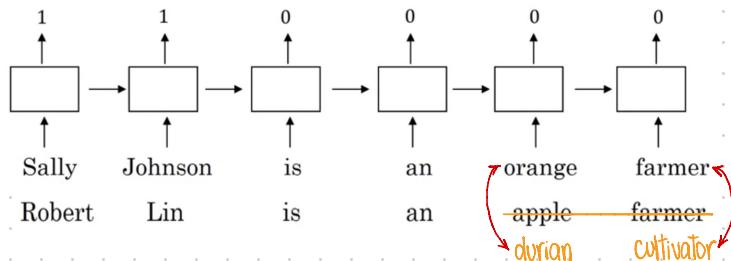
e_{456}

* visualizing word embeddings → you take the 300-D vector and visualize it in a 2-D space using t-SNE.



USING WORD EMBEDDINGS

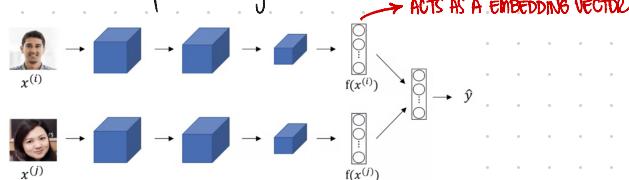
* named entity recognition example:



* because of word embedding your algorithm can learn that both apple and durian are fruits, and that farmer and cultivator are similar, too.

- #1. learn word embeddings from a large text corpus (1 - 10B words)
- #2. Transfer embedding to new task with smaller training set (100k words)
- #3. Optional: continue to finetune the word embeddings with new data.

* relation with face encoding:



PROPERTIES OF WORD EMBEDDINGS

* they can also help with analogy reasoning

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
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	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

Emman Euwoman

* Can we define an algorithm that understands that man \rightarrow woman ?
King \rightarrow queen ?