

Week 1 - Recurrent Neural Networks

Wednesday, September 30, 2020 3:07 PM

① Why Sequence Models

- Used for sequential data

Examples of sequence data

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

- All of these can be addressed using supervised learning with label x, y as the training set

② Notation

- Motivating example

\rightarrow Named entity recognition

- Outputs the names in a sentence

x : Harry Potter and Hermoine Granger invented a new spell

y : 1 1 0 ... 1 1 0 0 0 0

$T_x = 9 \rightarrow$ length of input

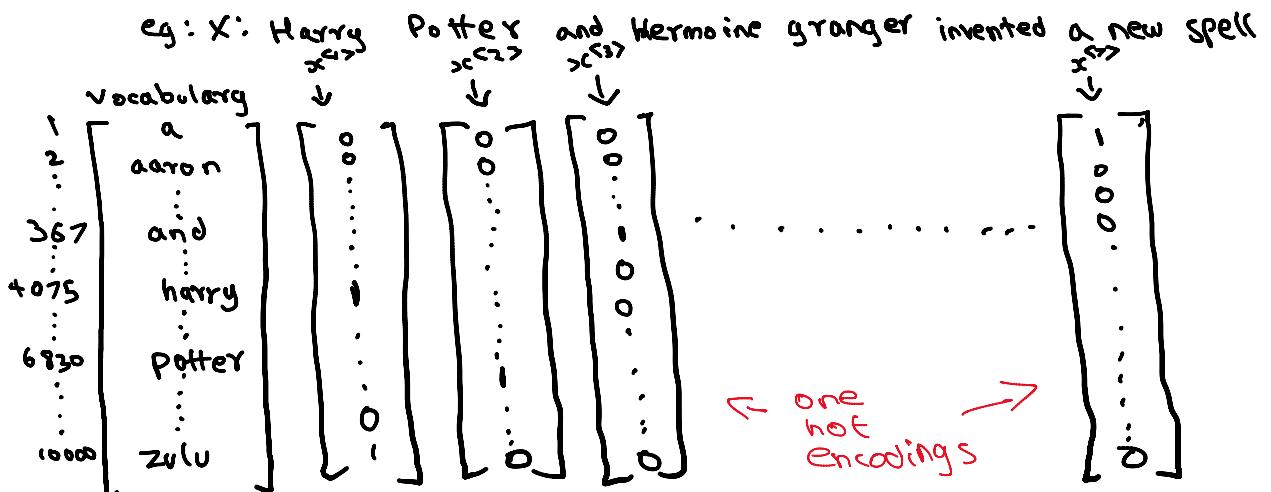
$T_y = 9 \rightarrow$ length of output

$\{x^{(i)} < t > T_x^{(i)}\}$ for i training examples

- Representing words

In NLP, we need a vocabulary list that contains all the words in our target sets

e.g.: x : Harry Potter and Hermoine Granger invented a new spell

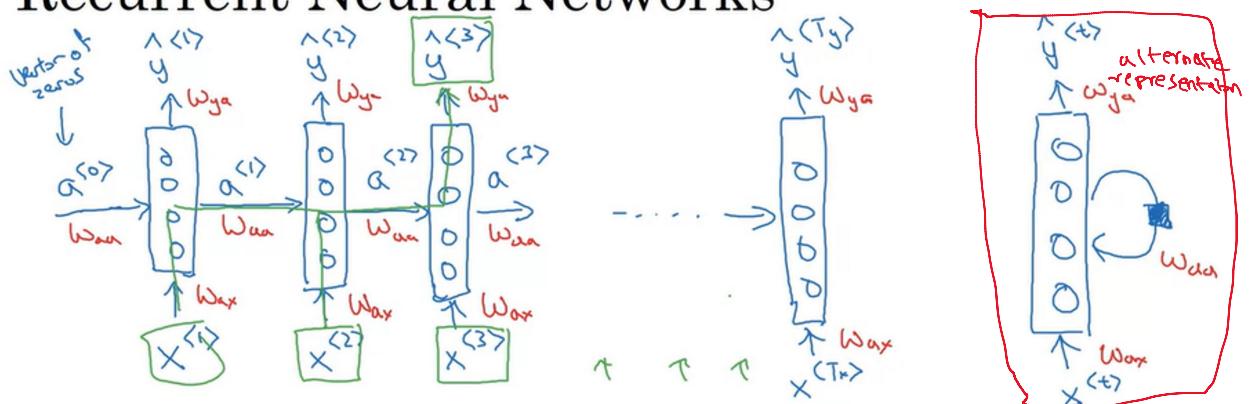


Unknown words could be represented by the same field in the vocabulary

③ Recurrent Neural Network model

- Why not a standard network
 - Inputs and outputs differ in length
 - No shared features learned across different positions of text
 - One-hot encoded vectors are high dimensional, requiring lot of parameters
- Recurrent Neural Networks

Recurrent Neural Networks



→ $a^{<0>}$ is initialized with zeros, but some may initialize randomly

→ The weights of previous layers are passed on to the subsequent layers

→ Three weight matrices are involved
 \rightarrow W_{aa} -> conventional input

→ Three weight matrices are involved

- W_{ax} → sequential input
shape = (No. of hidden neurons, n_x)
memory of the NN
- W_{aa} → computation of previous layers
shape = (No. of hidden neurons, No. of hidden neurons)
- W_{ya} → sequential output
shape = (n_y , No. of hidden neurons)

→ Output $\hat{y}^{<t>}$ depends on previous inputs and activations

→ Downside is that it can't learn from the elements later in the sequence

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

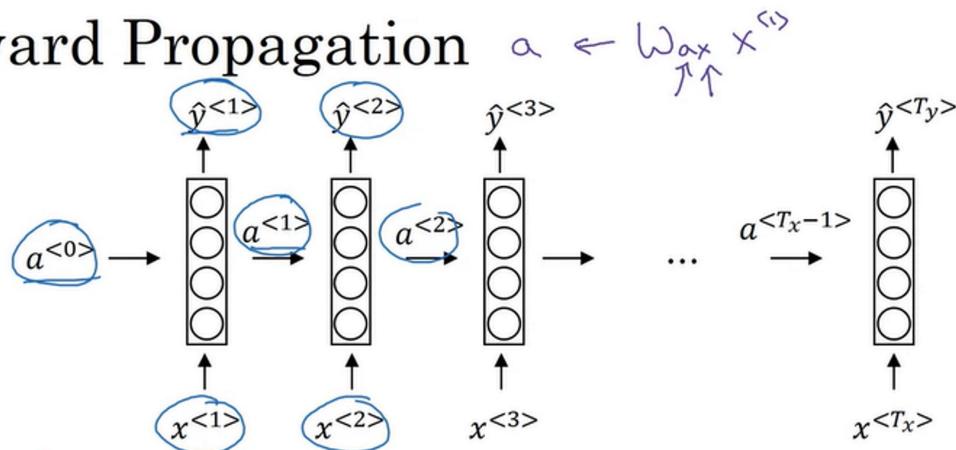
He said, "Teddy bears are on sale!"

→ Cannot output if Teddy is the name of the president

→ This can be addressed using Bidirectional RNN (BRNN)

• Forward Propagation

Forward Propagation



$$a^{<0>} = 0$$

$$\begin{aligned} a^{<t>} &= g_1(w_{aa}a^{<t-1>} + w_{ax}x^{<t>} + b_a) \quad \rightarrow \text{usually tanh/ReLU} \\ \hat{y}^{<t>} &= g_2(w_{ya}a^{<t>} + b_y) \quad \rightarrow \text{depends on output (sigmoid for 2 or p classes)} \end{aligned}$$

$$\boxed{\begin{aligned} a^{<t>} &= g(w_{aa}a^{<t-1>} + w_{ay}a^{<t>} + b_a) \\ \hat{y}^{<t>} &= g(w_{ya}a^{<t>} + b_y) \end{aligned}}$$

• Simplified RNN notation

$$\begin{aligned} a^{<t>} &= g(w_{aa}a^{<t-1>} + w_{ax}x^{<t>} + b_a) \Rightarrow a^{<t>} = g(w_a[a^{<t-1>}, x^{<t>}] + b_a) \\ \hat{y}^{<t>} &= g(w_{ya}a^{<t>} + b_y) \quad w_a = [w_{aa} \quad w_{ax}] \\ \boxed{\hat{y}^{<t>} = g(w_y a^{<t>} + b_y)} \end{aligned}$$

eg: if $w_{aa} \rightarrow (100, 100)$ $a^{<t-1>} \rightarrow (100)$
 $w_{ax} \rightarrow (100, 10000)$ $x^{<t>} \rightarrow (10000)$

$$\text{L} \geq \boxed{\hat{y}^{(t)} = g(w_y a^{(t)} + b_y)}$$

eg: if $w_{aa} \rightarrow (100, 100)$ $a^{(t-1)} \rightarrow (100)$
 $w_{ax} \rightarrow (100, 10000)$ $x^{(t)} \rightarrow (10000)$

$$100 \uparrow [w_{aa} \quad w_{ax}] = w_a$$

$\xleftarrow[10000]{10000} (100, 10000)$

↳ stack both vectors

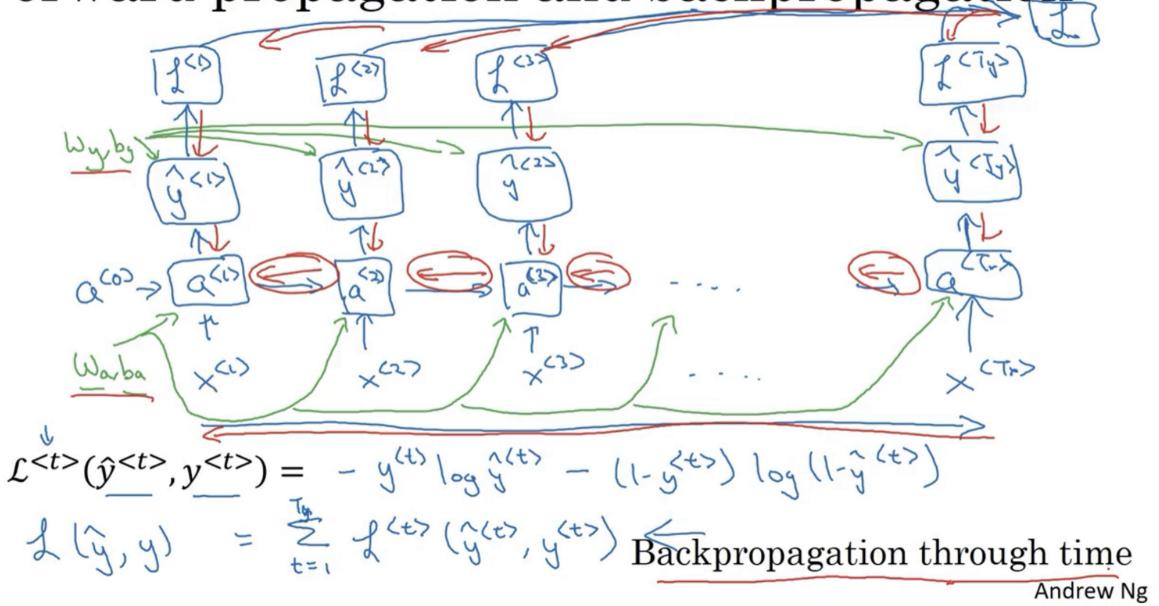
$$\Rightarrow [w_{aa} \quad w_{ax}] \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} = w_{aa} a^{(t-1)} + w_{ax} x^{(t)}$$

↳ original eqn.

④ Backpropagation through time

- Implemented automatically by the framework K
- Computation graph

Forward propagation and backpropagation

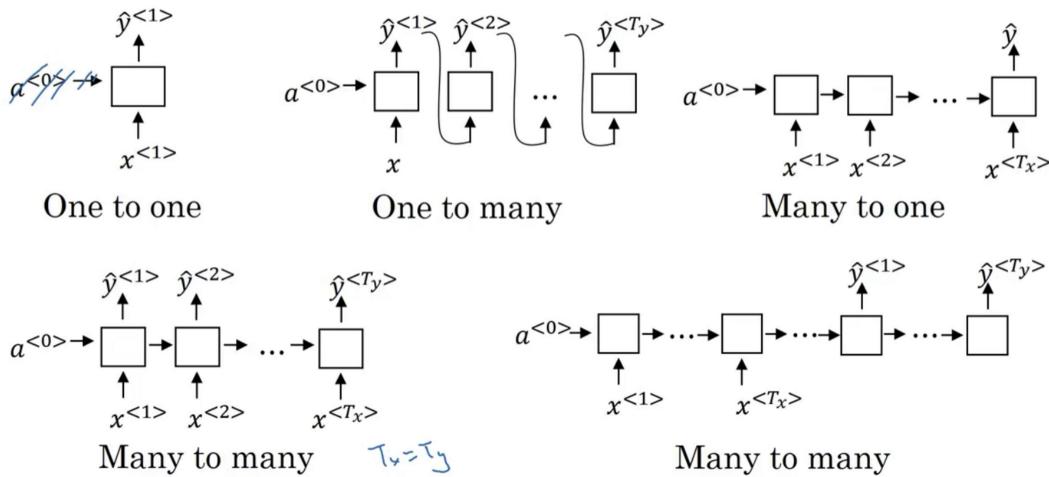


- w_a, b_a, w_y, b_y are shared across each element in a sequence
- cross entropy function is used to compute loss
- it's called "backpropagation through time" because we pass activation from one sequence element to another backwards in time.

⑤ Different types of RNNs

- we've seen RNN for $T_x = T_y$
- Other cases are

Summary of RNN types



Andrew Ng

- Attention models also exist which is more complex than above architectures

⑥ Language model and sequence generation

- Language model
 - It estimates the probability of that particular sequence of words
 - eg: In speech recognition

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

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

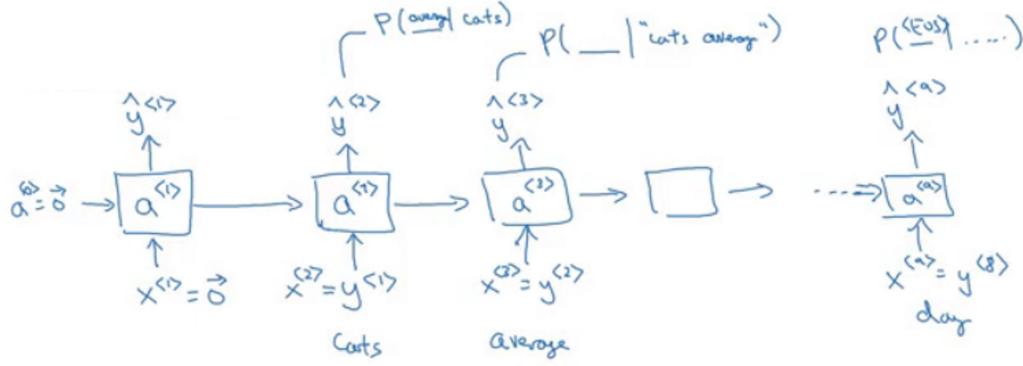
Even though these sentences sound alike, the second sentence is 3 orders more probable to occur.

$$\Rightarrow P(\text{sentence}) = ? \Rightarrow P(y^{<1>}; y^{<2>}; \dots; y^{<T_y>})$$

- Language modelling with an RNN
 - Training set: Large corpus of English text
 - We tokenize the input sentence by mapping words to vocabulary set (one-hot encoding)
 - Add an `<EOS>` (end of sentence token) and `<UNK>` token for unknown words
 - We also set $x^{<t>} = y^{<t-1>}$

• RNN model

→ For the sentence "cats average 15 hours of sleep a day"



→ Loss function is the cross entropy loss

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

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

$i \rightarrow$ all elements in corpus
 $t \rightarrow$ for the entire sequence

→ To use the model

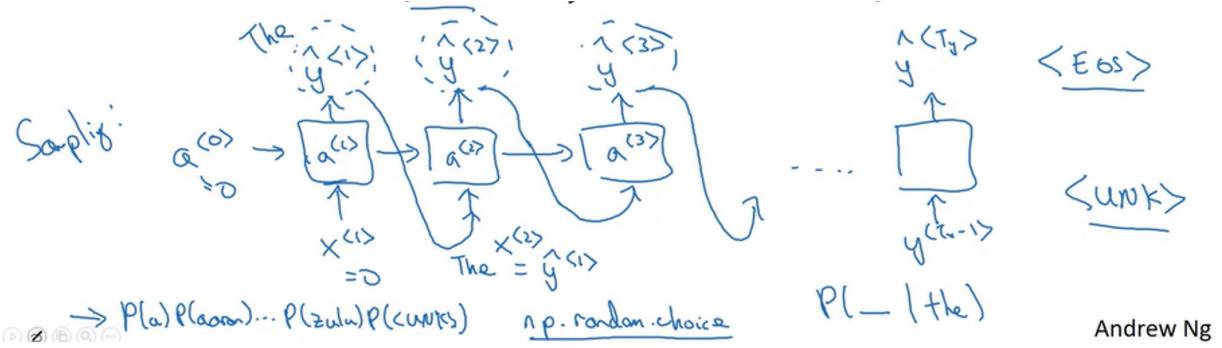
- To predict next word given the previous words, we feed the sentence and get the final $\hat{y}^{<t>}$ hot vector and sort by maximum probability
- To compute probability of a sentence with 3 words

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

⇒ We multiply probability of occurrence of each word to get probability of occurrence of the sentence

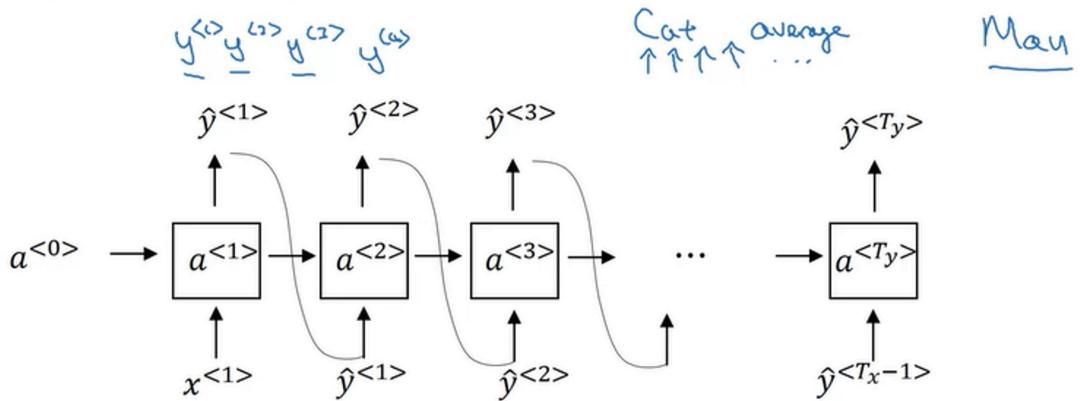
⑦ Sampling novel sequences

- To check what a language model has learned, we can try to sample novel sequences.
- To sample a novel sequence



- we pass $a^{<>} = \text{zeros}$ & $x^{<>} = \text{zeros}$
 - Random word from output $y^{<>}$ is chosen as the first word
 - The predicted word is passed as $x^{<2>}$
 - The above 2 steps are repeated until we get $\langle \text{EOS} \rangle$ or after r the required no. of iterations
 - $\langle \text{UNK} \rangle_s$ can be removed

- Character-level language model



- Pros
 - No unknown token
 - Cons
 - Much longer sequences
 - Not good at long range dependencies in a sentence
 - Computationally more expensive and harder to train

⑧ Vanishing gradients with RNNs

- RNNs turn into vanishing gradient problems because of extremely long term dependencies in NLP.
eg: The cat, which already ate was full
The cats, which already ate ., ., were full
The rat/cats affects was/were over arbitrary no. of words

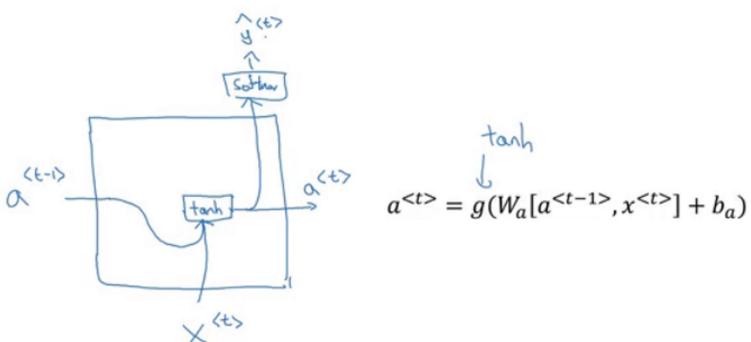
The cats, which already ate. . . were full

The rat/cats affects was/were over a long range

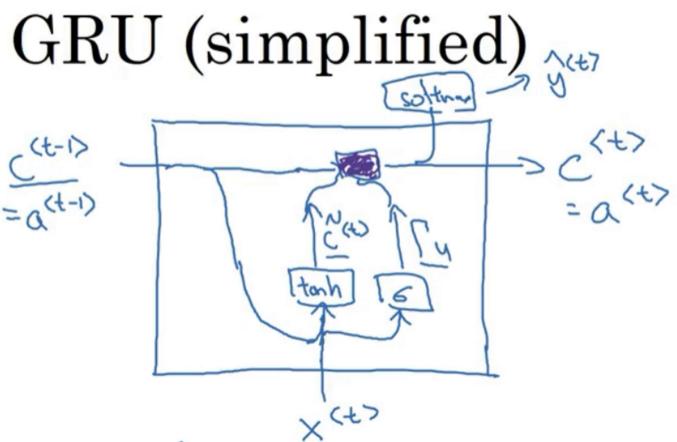
- RNNs are not very good at learning long term dependencies
- Similar to deep neural networks, we run into vanishing gradient problems with deeper networks or RNNs with large sequence size.
- Exploding gradients also happen but can be easily identified as they cause a lot of NaNs. They are easily fixed by gradient clipping

⑨ Gated Recurrent Unit (GRU)

- Basic RNN unit



- Each layer in a GRU has a variable memory cell which is the $c^{(t)} = a^{(t)}$
 \hookrightarrow memory cell
- $\tilde{c}^{(t)} = \tanh(W_c[c^{(t-1)}, x^{(t)}] + b_u)$
- $\Gamma_u = \sigma(W_u[c^{(t-1)}, x^{(t)}] + b_u)$
 \hookrightarrow update gate
- $c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$
- update gate is from 0 to 1 (mostly 0 or 1)
- we update c based on update cell and previous cell
- It updates when the memory cell is no longer needed



Shapes:

- $a^{<t>} \rightarrow (\text{No of hidden Neurons})$
- $c^{<t>} \approx a^{<t>} \text{ same as } a^{<t>}$
- $y^{<t>} \text{ also same as } a^{<t>}$

- FULL GRU

Full GRU

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

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

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

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

- it contains an extra gate that is used to calculate candidate
- tells how relevant $c^{<t-1>}$ is to compute $c^{<t>}$

⑩ Long Short Term Memory (LSTM)

- more powerful and generalised than GRU
- In LSTM $c^{<t>} = a^{<t>}$
- Equations

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

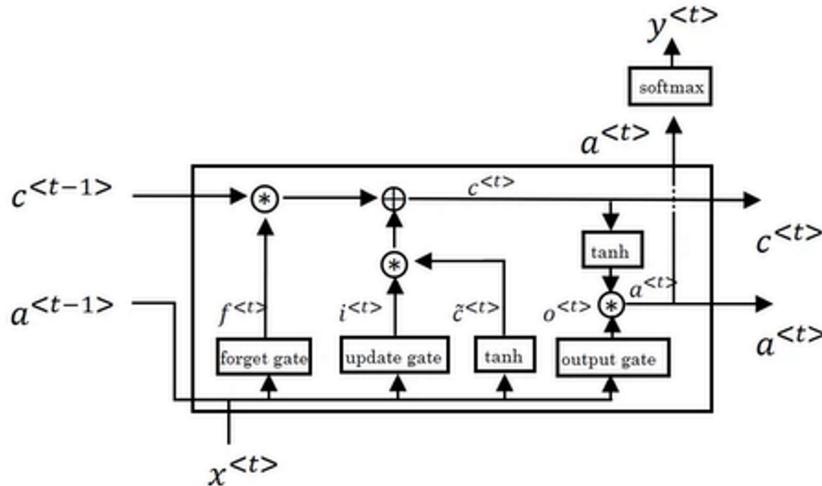
$$\text{update gate} \leftarrow \Gamma_u = \sigma(W_u [a^{<t-1>}, x^{<t>}] + b_u)$$

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

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

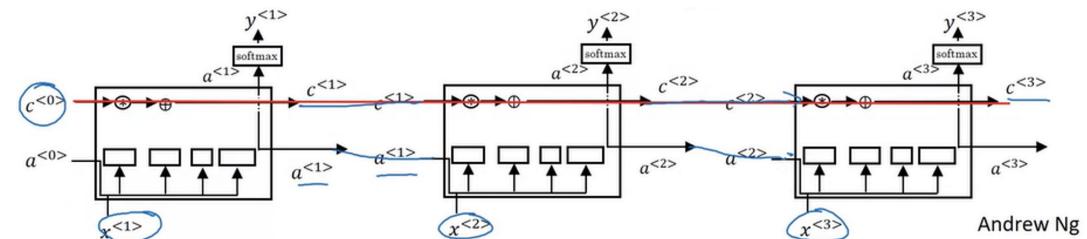
$$\begin{aligned}
 \text{gate} &\leftarrow f = \sigma(w_f[a^{t-1}, x^t] + b_f) \\
 \text{output gate} &\leftarrow o = \sigma(w_o[a^{t-1}, x^t] + b_o) \\
 c^{t-1} &= f * \tilde{c}^{t-1} + o * c^{t-1} \\
 a^{t-1} &= o * \tanh(c^{t-1})
 \end{aligned}$$

- LSTM in picture



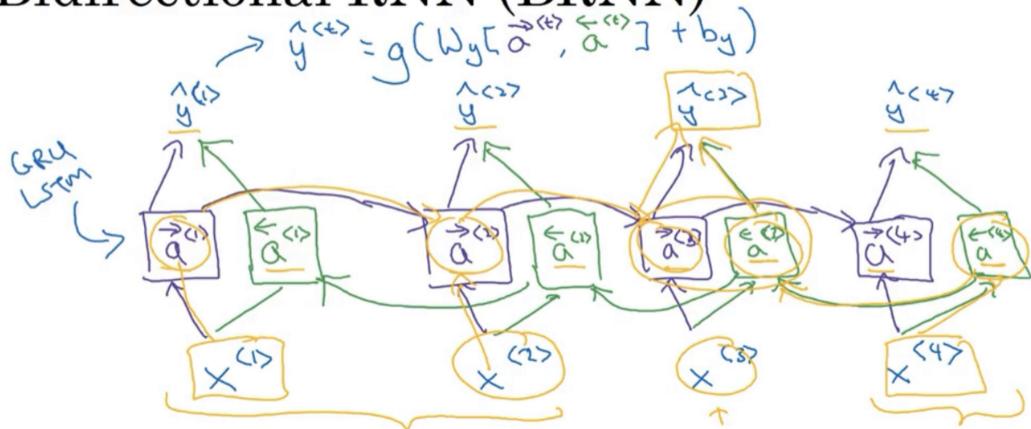
→ Some variations of LSTMs include 'peephole connections' which include c^{t-1} in every gate instead of just a^{t-1} .

- When LSTMs are chained together, a straight connection from c^{t-1} to c^t can still be made from c^{t-1} to c^t . This allows it to remember long term dependencies.



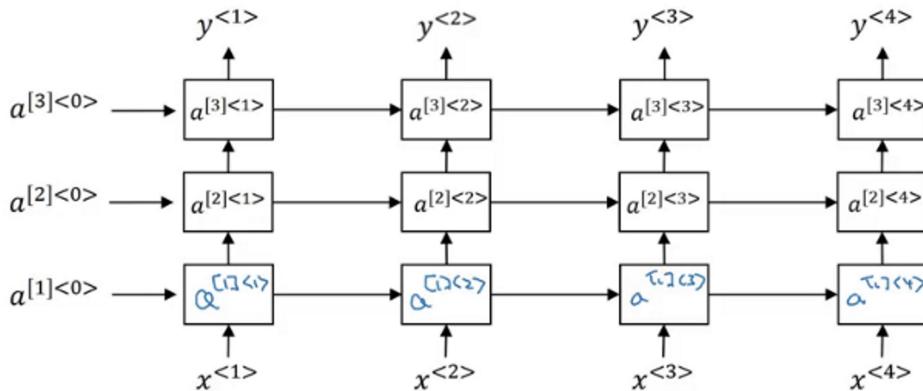
⑪ Bi-directional RNN

Bidirectional RNN (BRNN)



- Acyclic graph
forward propagation goes left to right
and right to left since it learns
from both sides
- The blocks can be basic RNN, GRU or LSTMs
- Usually BRNN with LSTM is used for
many NLP tasks
- Disadvantage is that it requires entire
sequence to make predictions. Hence,
it cannot be used in real time
tasks.

⑫ Deep RNNs



- RNN with 3 hidden layers
- Unlike feed-forward neural networks,
in deep RNNs stacking 3 layers is
considered deep and expensive to
train
- In some models, feed forward neural
networks are connected after recurrent

→ In some models, feed forward neural networks are connected after recurrent cell