

Sequence Models

*enlivo

Date: ___/___/___

Page: ___

Examples of Sequence data

- Speech recognition
- Music generation
- Sentiment classification
- DNA sequence analysis
- Machine translation
- Video activity recognition
- Name entity recognition

★ Recurrent Neural Network

Notation.

Motivating Example: Name entity recognition

x : Harry potter and Harmione Granger invented a new spell

y : I O I O O T Q I

~~each~~ represent As:

x : $x^{(1)}$ $x^{(2)}$ $x^{(3)}$... $x^{(t)}$... $x^{(9)}$ $T_x = 9$

y : $y^{(1)}$ $y^{(2)}$ $y^{(3)}$... $y^{(t)}$... $y^{(9)}$ $T_y = 9$

Therefore,

$x^{(i)\leftarrow t}$: t^{th} element in Sequence in i^{th} Example of training

$y^{(i)\leftarrow t}$: t^{th} ~~next~~ element in Sequence in i^{th} output of training

$T_x^{(i)}$: length of Input Sequence in i^{th} Example.

$T_y^{(i)}$: length of output Sequence in i^{th} Example.

Representing Words.

* Harry potter and Hermione Granger

P.T.O \rightarrow

P = $\begin{bmatrix} 1 & 0 & 0 & \dots & 0 \end{bmatrix}$

T = $\begin{bmatrix} 0 & 1 & 0 & \dots & 0 \end{bmatrix}$

Representing Words.

X: Harry Potter and Hermione Granger invented a new spell.
 $x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>} \quad x^{<6>} \quad x^{<7>} \quad x^{<8>} \quad x^{<9>}$

Vocabulary

a	1	0	0	0	0	1	0	0	1	0
aaron	2	0	0	0	0	1	1	0	0	0
and	367	0	0	0	1	367	0	0	0	0
harry	4075	1	4075	0	0	10,000	0	0	0	0
Potter	6830	0	1	6830	0	0	0	0	0	0
zulus	10,000	0	0	0	0	0	0	0	0	0

Size of dictionary

we did "one-hot" encoding.

= 10,000 words in Harry Potter

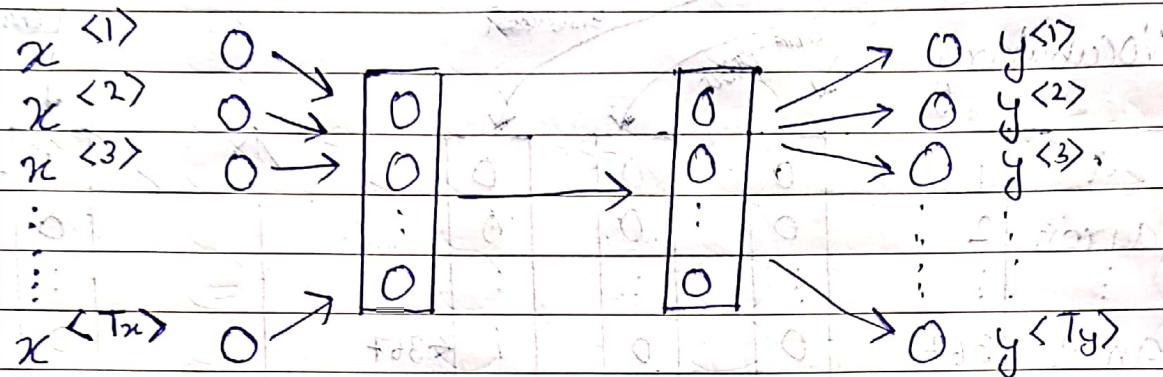
→ So, In this representation, $x^{<t>}$ will be one-hot vector.

it is

will do this as Supervised learning method
 $(X \rightarrow Y)$

RNN model

Q Why not use a Standard Network?

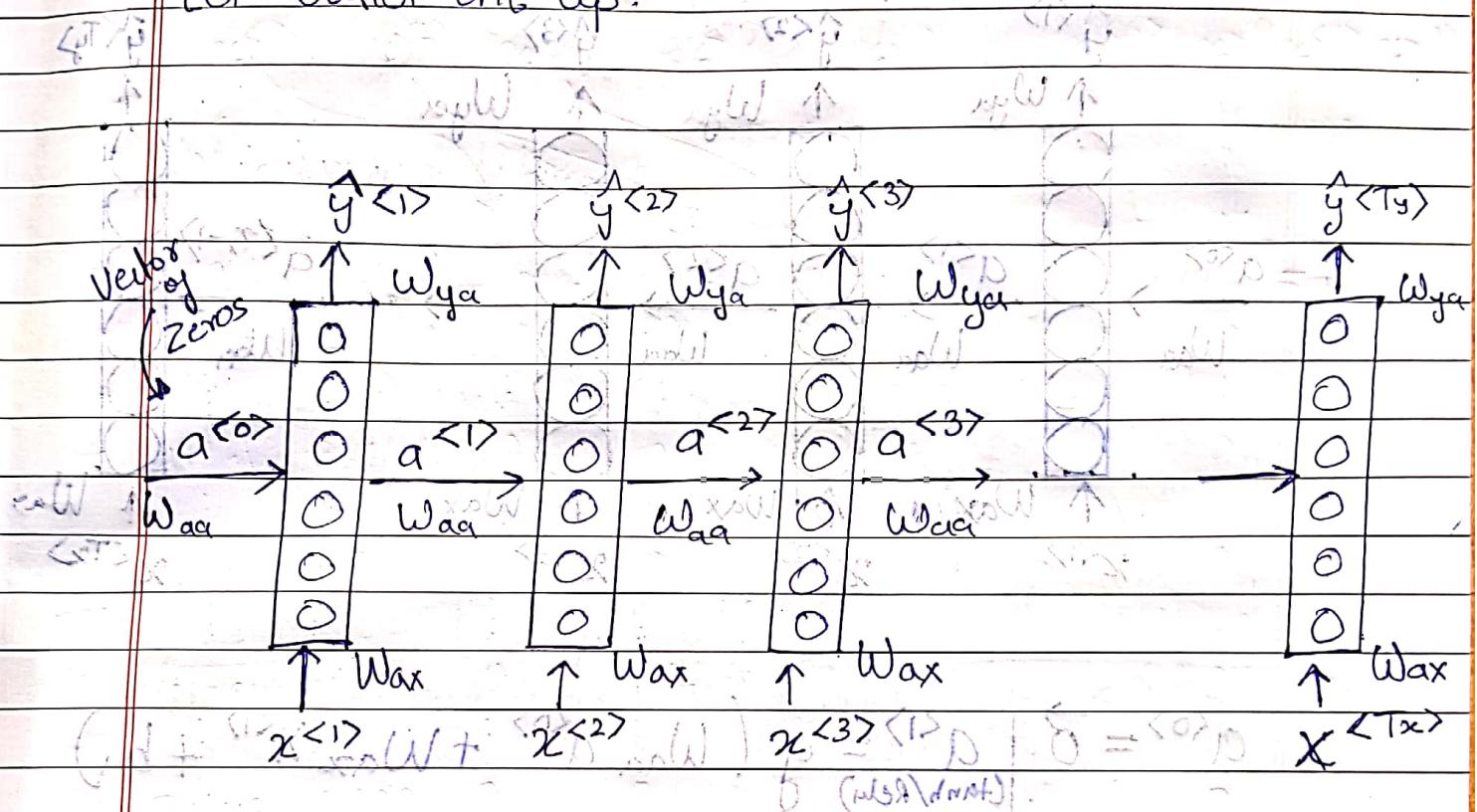


problems:

- ⇒ Input, outputs can be different length in different Examples.
- ⇒ Doesn't share features learned across different positions of text.
- ⇒ Weight matrix will end-up have enormous number of parameters.
- ∴ Standard Network won't work.

Q. What is Recurrent Neural Networks.

→ Let's build one up:

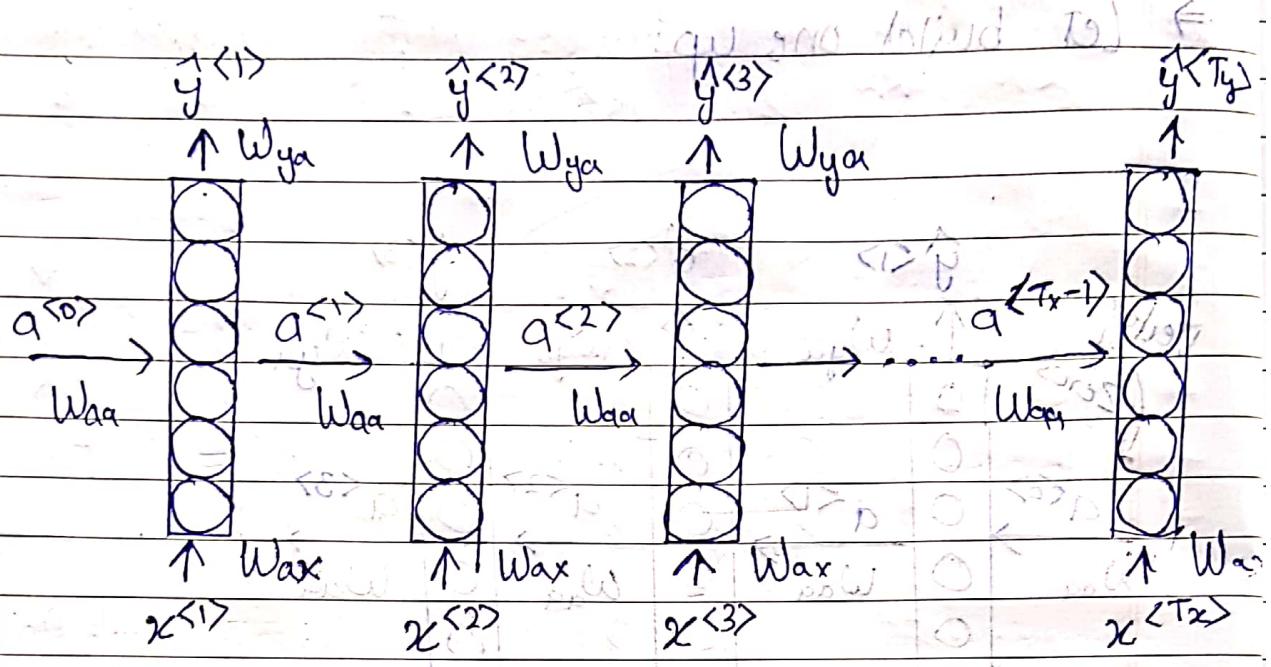


→ So, in this RNN what does it mean is that when making the prediction of $\hat{y}^{<3>}$, it gets information not only from $x^{<3>}$ but also information from $x^{<1>}$ and $x^{<2>}$ because information on $x^{<1>}$ can pass thru from $x^{<2>} \rightarrow x^{<3>}$ to help with prediction.

Limitation:

→ Only uses information that earlier in the sequence, in particular, when predicting $\hat{y}^{<3>}$ it does not use information from $x^{<4>}$, $x^{<5>}$, $x^{<6>}$, etc.

\Rightarrow forward propagation (forward pass) in terms of



$$a^{(0)} = \vec{0}$$

$$a^{(1)} = g_{(a)}(W_{aa} a^{(0)} + W_{ax} x^{(1)} + b_a)$$

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

Sigmoid/Softmax

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

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

In $a^{(t)}$ activation could be, i.e. $g(\cdot)$,

$\Rightarrow \tanh(\cdot), \text{relu}(\cdot)$

⇒ Simplified RNN Notation

$$a^{(t)} = g(W_a \cdot a^{(t-1)} + W_x \cdot x^{(t)} + b_a)$$

$$y^{(t)} = g(W_y \cdot a^{(t)} + b_y)$$

$$\rightarrow a^{(t)} = g(W_a \cdot a^{(t-1)} + W_x \cdot x^{(t)} + b_a)$$

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

i.e.

$$W_a = [W_a | W_x]$$

Take W_a and W_x and Stack them side by side as shown above that we create W_a .

For Example: If a is 100 dim, and x is 10,000 dim.

Then $W_a \in \mathbb{R}^{100 \times 100}$ and $W_x \in \mathbb{R}^{100 \times 10,000}$

∴ Stacking this two matrix will be,

$$100 \uparrow [W_a | W_x] = W_a \in \mathbb{R}^{100 \times 10,100}$$

$\xleftarrow{100} \xrightarrow{10,000}$

And same with vectors.

Stacking $a^{(t-1)}$ and $x^{(t)}$

$$[a^{(t-1)}, x^{(t)}] = \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} \quad \begin{array}{c} \uparrow 100 \\ \downarrow \\ \uparrow 10,000 \end{array} \quad \begin{array}{c} \uparrow 10,100 \\ \downarrow \end{array}$$

therefore,

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

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

which equal to:

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

∴ Now here, No change needed for this

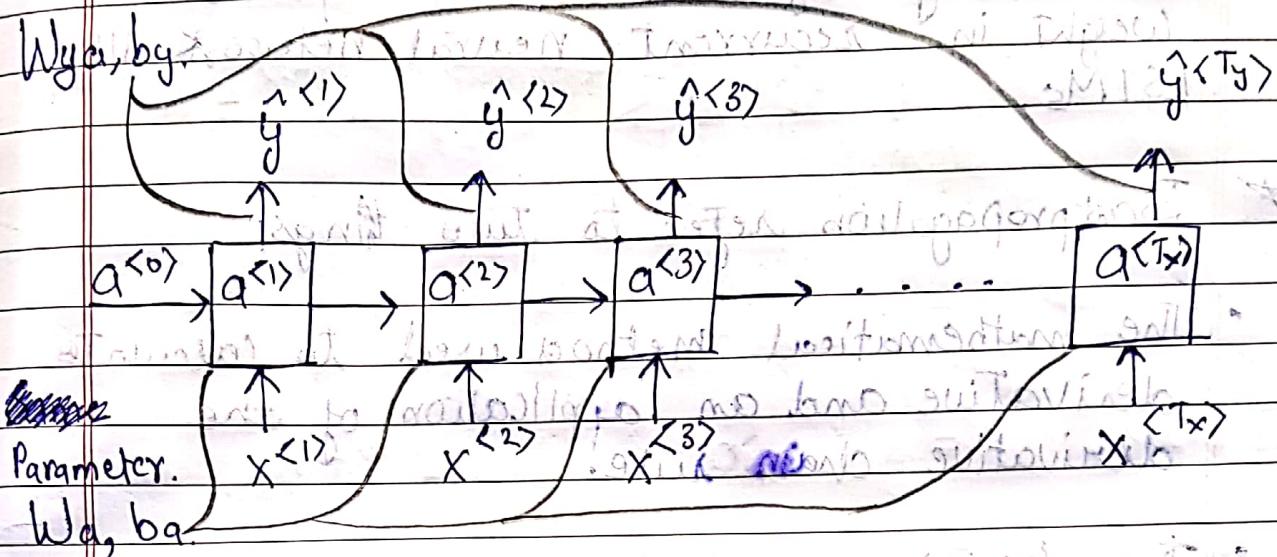
$$y^{(t)} = g(W_y \cdot a^{(t)} + b_y)$$

↑ ↑ ↑
output weight output bias.

↑
at time
step 't'
metric

x shown cross A predicted?

→ Backpropagation through time



* Loss function (Cross-Entropy)

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

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

Backpropagation Through Time, Or BPTT, is the training algorithm used to update weight in recurrent neural networks like LSTMs.

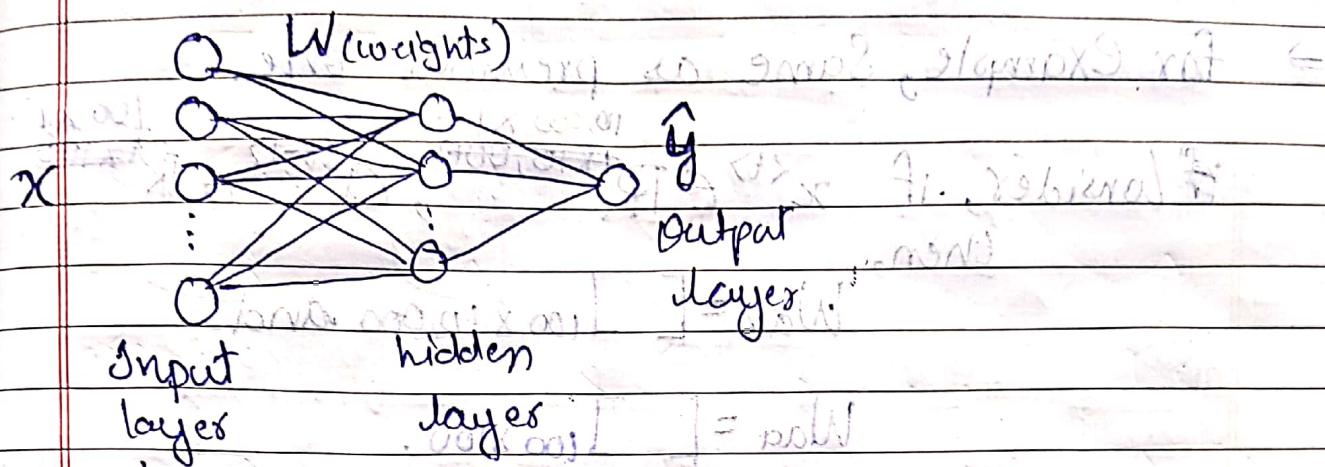
- * Backpropagation refer to two things:
 - The mathematical method used to calculate derivative and an application of the derivative chain rule.
 - The training algo. for updating network weights to minimize error.
- It is a supervised learning algo. that allows the network to be corrected with regard to the specific error made.
- * The general algorithm is as follows:
 - 1) Present a training input pattern and propagate it through the network to get output.
 - 2) Compare the prediction output to the expected and calculate the error.
 - 3) Calculate the derivatives of the error with respect to the network weights.
 - 4) Adjust the weights to minimize the error.
 - 5) Repeat.

Self-explination

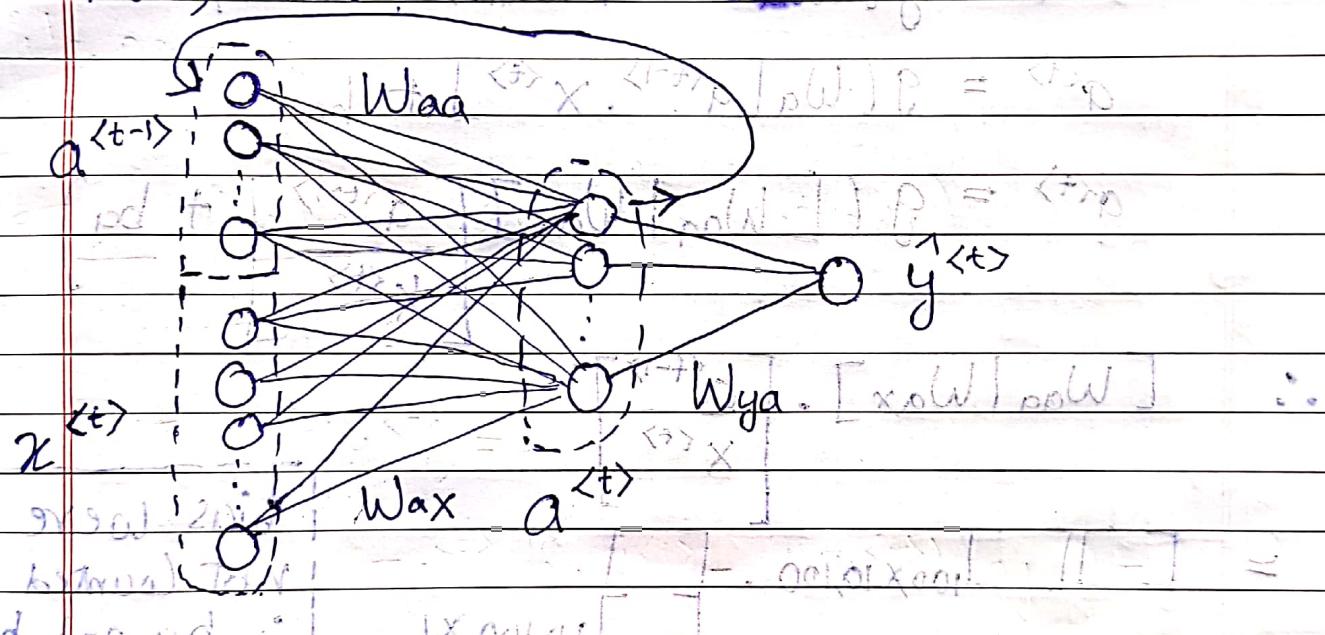
*enlivio

What RNN look like according to me?

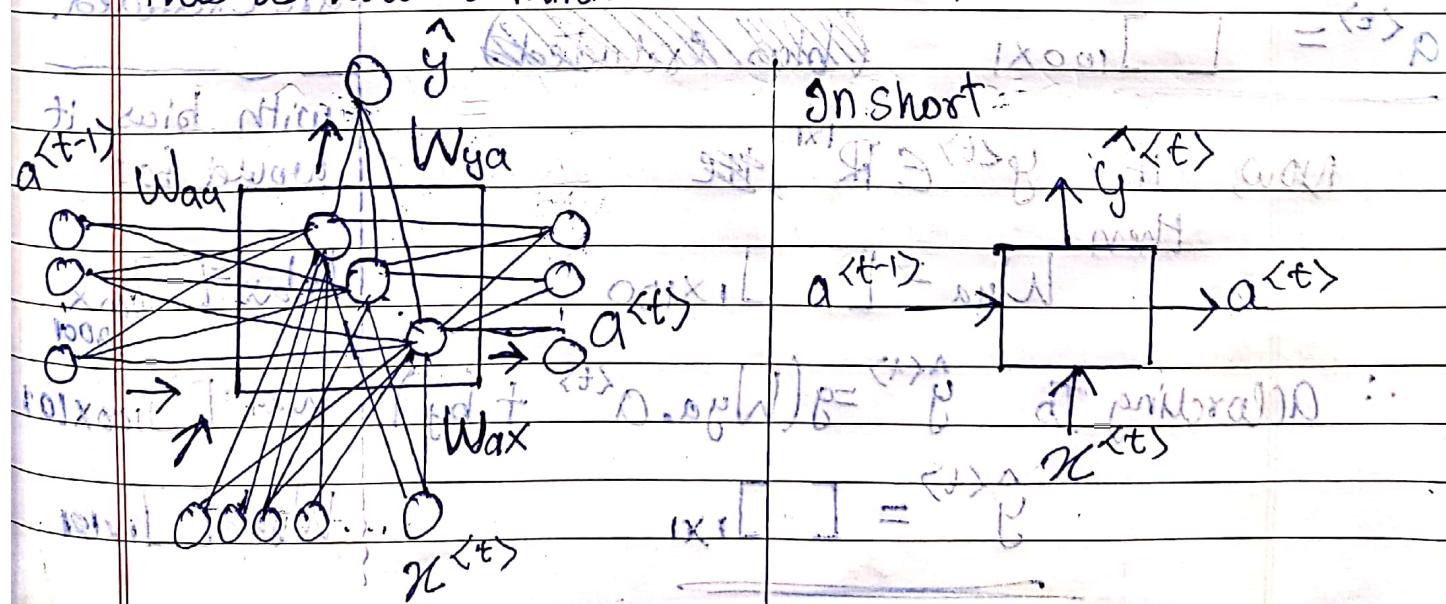
⇒ first, Standard NN. ~~using forward pass~~ standard



→ Now, RNN



→ This is how I think it looks like.



Say-Explanation

Date: _____ *enlive
Page: _____

Q. How input, weight and bias pass thru to calculate $a^{(t)}$ and outputs? (A) ~~Input~~ ~~Output~~ ~~Unit~~

⇒ for Example, Same as previous one.

* Consider, if $x^{(t)} \in \mathbb{R}^{10,000 \times 1}$, $a^{(t-1)} \in \mathbb{R}^{10,000 \times 1}$, $a^{(t)} \in \mathbb{R}^{100 \times 1}$
then,

$$W_{ax} = [\quad]_{100 \times 10,000} \text{ and }$$

$$W_{aa} = [\quad]_{100 \times 100}.$$

$$\text{as we know, } a^{(t)} = g(W_{aa} \cdot a^{(t-1)} + W_{ax} \cdot x^{(t)} + b_a)$$

$$a^{(t)} = g(W_a[a^{(t-1)}, x^{(t)}] + b_a) \in \mathbb{R}^{100 \times 1}$$

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

$$\therefore [W_{aa}, W_{ax}] \cdot [a^{(t-1)}, x^{(t)}]$$

$$= [\quad]_{100 \times 10,000} [\quad]_{10,000 \times 1}$$

$$a^{(t)} = [\quad]_{100 \times 1}$$

Now, if $y^{(t)} \in \mathbb{R}^{1 \times 1}$, then

$$W_{ya} = [\quad]_{1 \times 100} \quad W_{ax} = [\quad]_{100 \times 10,000}$$

$$\therefore \text{According to } y^{(t)} = g(W_{ya} \cdot a^{(t)} + b_y) \quad W_{aa} = [\quad]_{100 \times 100}$$

$$y^{(t)} = [\quad]_{1 \times 1} \quad W_{ya} = [\quad]_{1 \times 100}$$

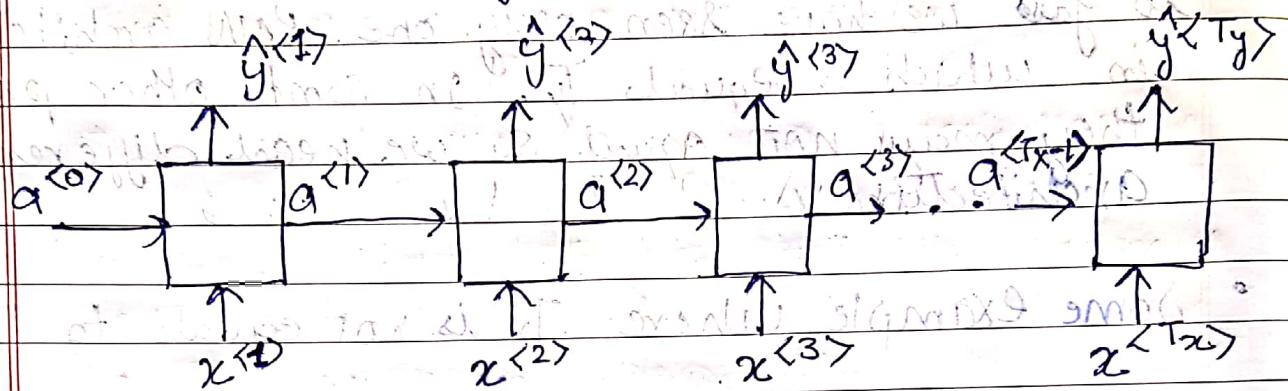
Different Types of RNNs

- So far we have seen only one RNN architecture in which T_x equals T_y . In some other problems, they may not equal so we need different architectures.
- Some examples where T_x is not equal to T_y ,

	x	y
→ Sentiment classification	"There is nothing but like in this movie!"	5 stars 1 to 5 rating
→ Machine translation (many-to-many but different length)	Voulez-Vous Chantar avec moi?	Do you want to sing with me?
→ Music generation	Null	Many words.
→ Name entity recognition	Yesterday, Harry potter met Ron	Yesterday, Harry potter met Ron

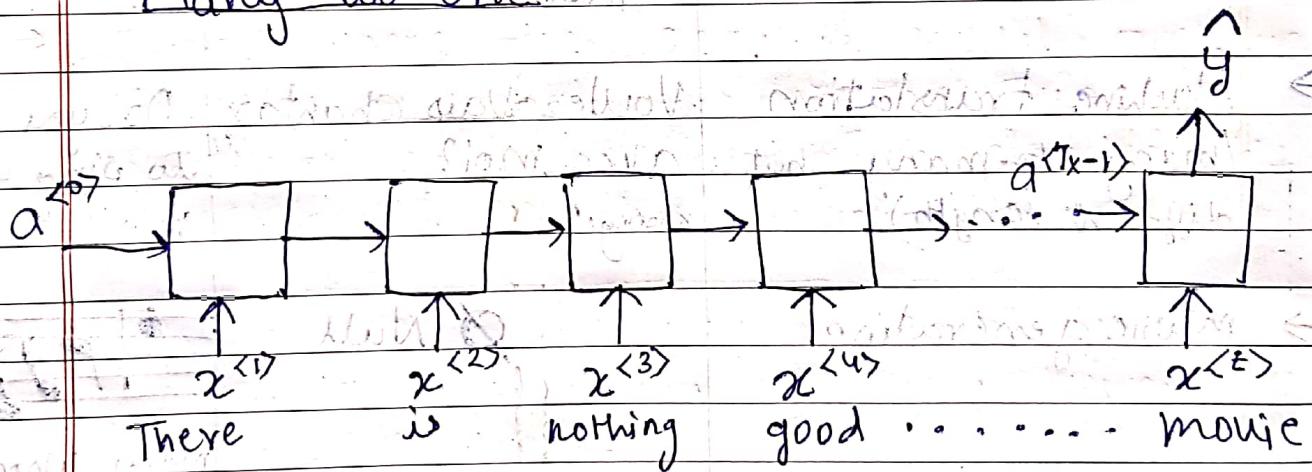
• Examples of RNN architectures:

→ when $T_x = T_y$, we have seen only this architecture so far. Many-to-many



→ In Sentiment analysis problem, X is a text while Y is an integer that ranges 1 to 5. The RNN architecture for that is:

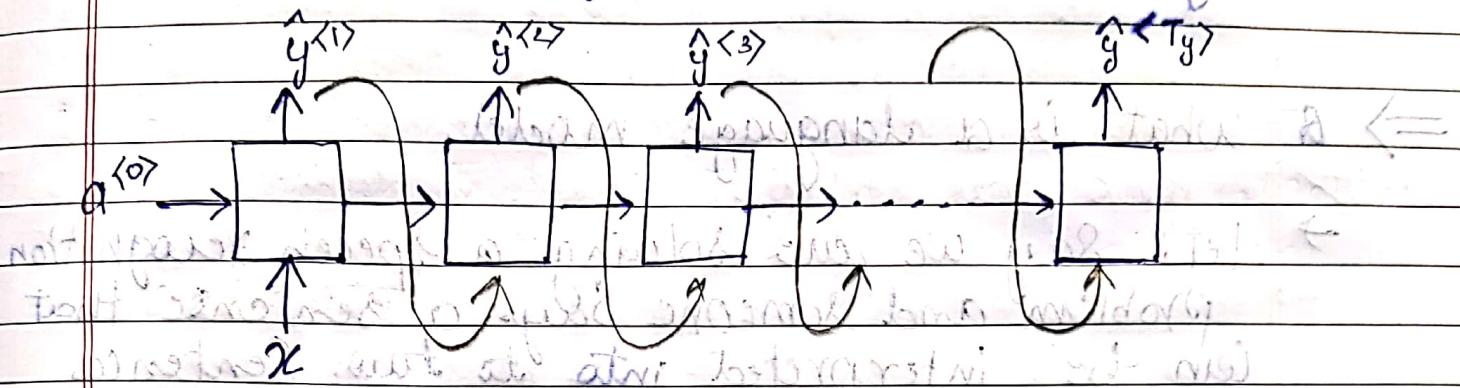
→ Many-to-One



→ classic, standard Neural Network known as one-to-one RNN



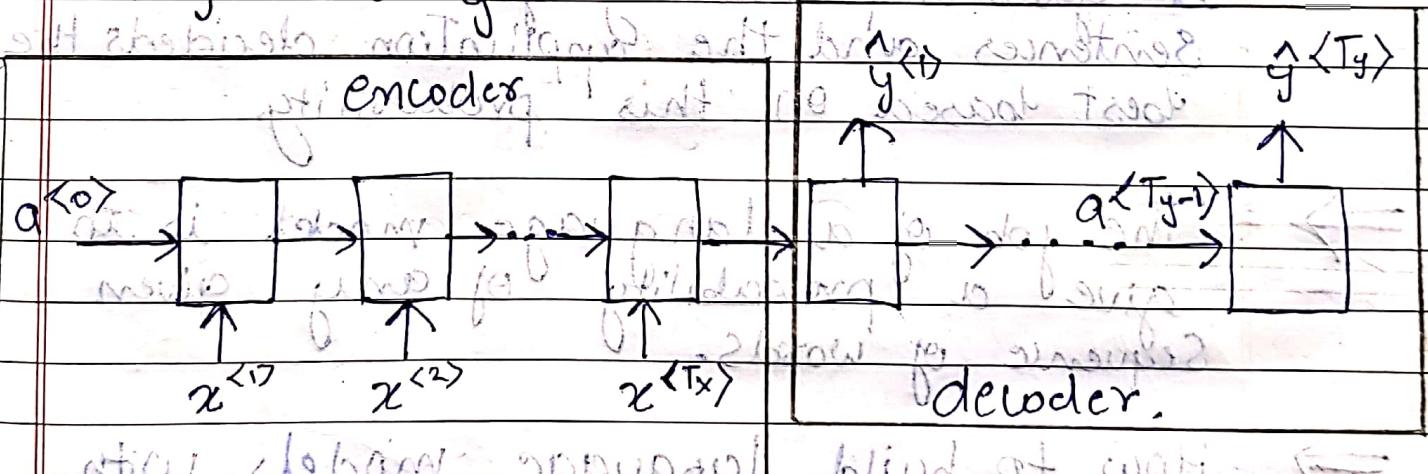
→ A one-to-many architecture in application would be music generation.



- Note that starting in the second layer we are feeding the generated output back to the ~~whole~~ Network.

→ ~~and soft softmax shown in the box~~

→ There are another interesting architecture in Many-to-many. Applications like machine translation inputs and outputs sequences have different lengths in most of the cases. ∵ alternative of Many-to-Many.



- There are an encoder and a decoder parts in this architecture. The encoder takes input ~~the~~ sequence into one matrix and feed it to the decoder to generate the output.

- Encoder and decoder have different weight matrices.



Language models and sequence generation.

⇒ Q What is a language model?

→ Let's say we are solving a speech recognition problem and someone says a sentence that can be interpreted into two sentences

• The apple and pear salad

• The apple and pears Salad

→ Pear and pear sounds exactly the same, so how would a speech recognition application choose from the two.

→ That's where the language model comes in. it gives a probability for the two sentences and the application decides the best based on this probability.

⇒ The job of a language model is to give a probability of any given sequence of words.

⇒ How to build language models with RNNs?

→ The first thing is to get a training set:

• A large corpus of target language text

• Structured training data

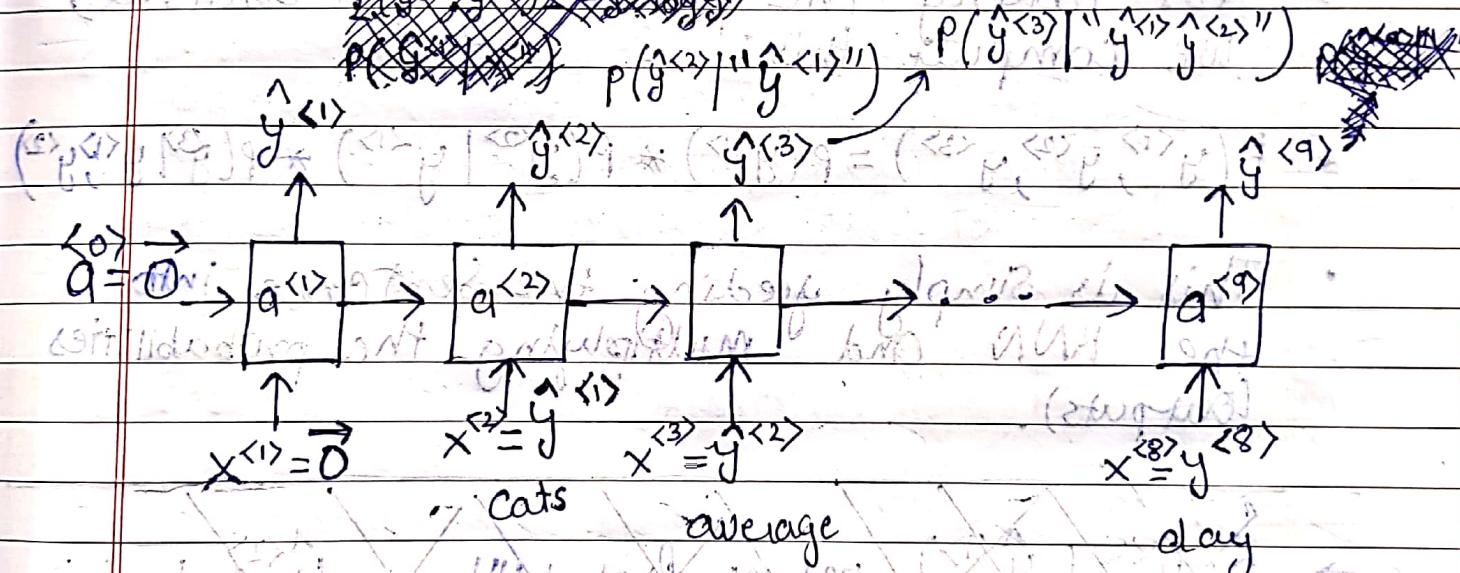
→ Then tokenize this training set by

register

getting the vocabulary and then one-hot each word.

→ Put an 'end of sentence' token $\langle \text{EOS} \rangle$ with the vocabulary and include it with each converted sentence. Also, use the token $\langle \text{UNK} \rangle$ for unknown words.

⇒ Given the sentence, $y = \langle \text{UNK} \rangle \text{ cuts average is hours of sleep a day.} \langle \text{EOS} \rangle$



→ Loss function is defined by cross-entropy loss:

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

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

Note $L(\hat{y}, y) = -\hat{y} \cdot \log(\hat{y})$
 if $\hat{y}=0 = -0 \cdot \log(0)$
 $= 0$

Now talk about next word prediction sentence

→ i is of all elements in the Corpus, this fall for all timesteps.

→ To use this model:

→ for predicting the chance of next word, we feed the sentence to the RNN and then get the final $y^{<t>}$ hot vector and sort it by maximum probability.

→ for taking the probability of a sentence, we compute this:

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

This is simply feeding the sentence into the RNN and multiplying the probabilities (outputs).

⇒ $P(y^{<2>} | y^{<1>})$ means probability of $y^{<2>}$ to be "average" due to $y^{<1>} (\text{Cats})$ as input.

⇒ $P(y^{<3>} | y^{<1>} | y^{<2>})$ means probability of $y^{<3>}$ to be "is" due to $y^{<1>}$ and $y^{<2>} (\text{Cat Average})$.

⇒ Likewise goes for all till $P(y^{<n>} | y^{<1>} | \dots | y^{<n-1>})$

⇒ Probability of sentence will be:

$$P(y^{<1>} | y^{<2>} | \dots | y^{<n>}) = P(y^{<1>}) * P(y^{<2>} | y^{<1>}) * \dots * P(y^{<n>} | y^{<1>} | \dots | y^{<n-1>})$$

As mentioned :

$y = \text{"Cats average 15 hours of sleep a day." (EOS) }$

$y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad y^{(4)} \quad y^{(5)} \quad y^{(6)} \quad y^{(7)} \quad y^{(8)} \quad y^{(9)}$

\Rightarrow here, $P(\cdot)$ means probability,

at first time step, $P(y^{(1)})$ means what is the probability of $y^{(1)}$ to be $y^{(1)}$ i.e "Cats", as $y^{(1)}$ will be output vector coming thru softmax, which means exactly ~~expresses~~ what will be the probability of "Cats" in output vector when input is feed to Network. $x^{(1)} = \vec{0}$, $c_1^{(0)} = \vec{0}$

Same way, at second time step, $P(y^{(2)} | y^{(1)})$ means what is probability of $y^{(2)}$ to be $y^{(2)}$ i.e "average" when $y^{(1)}$ was feed into network as $x^{(2)} = y^{(1)}$.

Same way, at third time step, $P(y^{(3)} | y^{(1)}, y^{(2)})$ means what is probability of $y^{(3)}$ to be $y^{(3)}$ i.e "15" when $y^{(2)}$ was feed into network as $x^{(3)} = y^{(2)}$.

So. on, till the end $P(y^{(1)} | y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(8)})$

\Rightarrow To calculate probability of Sentence.

$$P(y^{(1)}, y^{(2)}, y^{(3)}, y^{(4)}, y^{(5)}, y^{(6)}, y^{(7)}, y^{(8)}, y^{(9)}) =$$

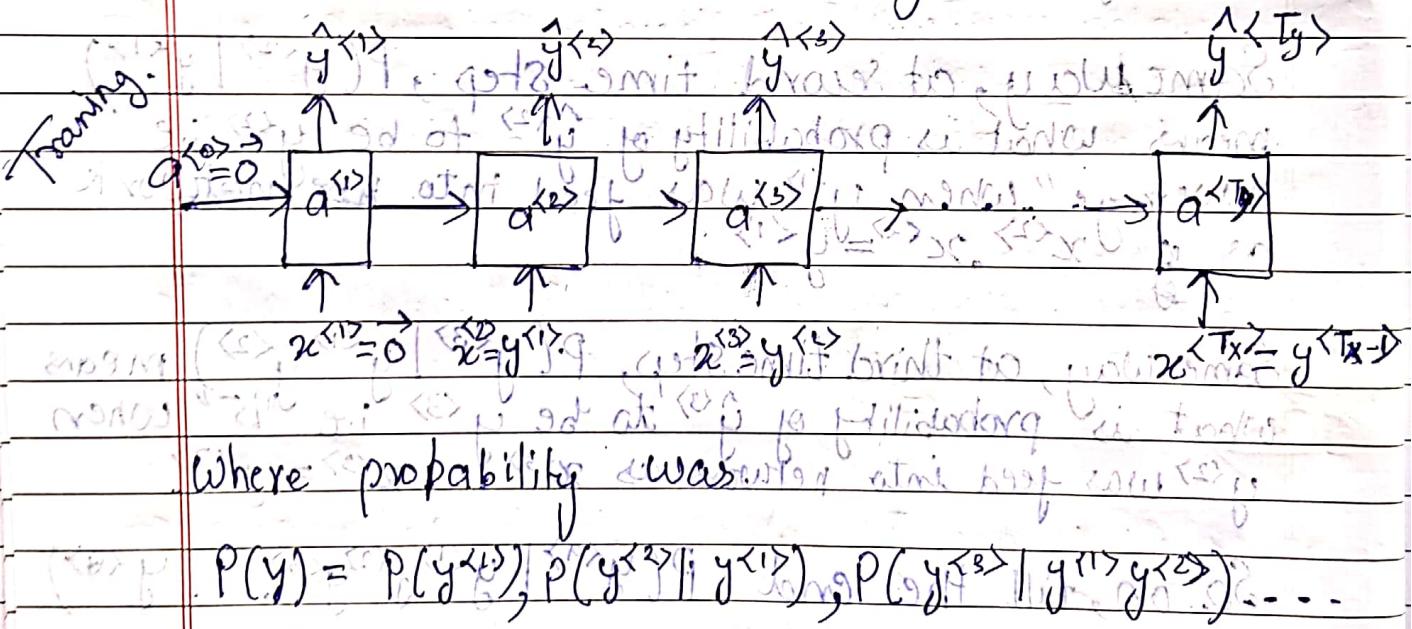
$$P(y^{(1)}) * P(y^{(2)} | y^{(1)}) * P(y^{(3)} | y^{(1)}, y^{(2)}) * \dots * \\ * P(y^{(9)} | y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(8)})$$

★ Sampling novel sequences.

After a sequence model is trained on a language model, to check what the model has learned you can apply it to sample a novel sequence.

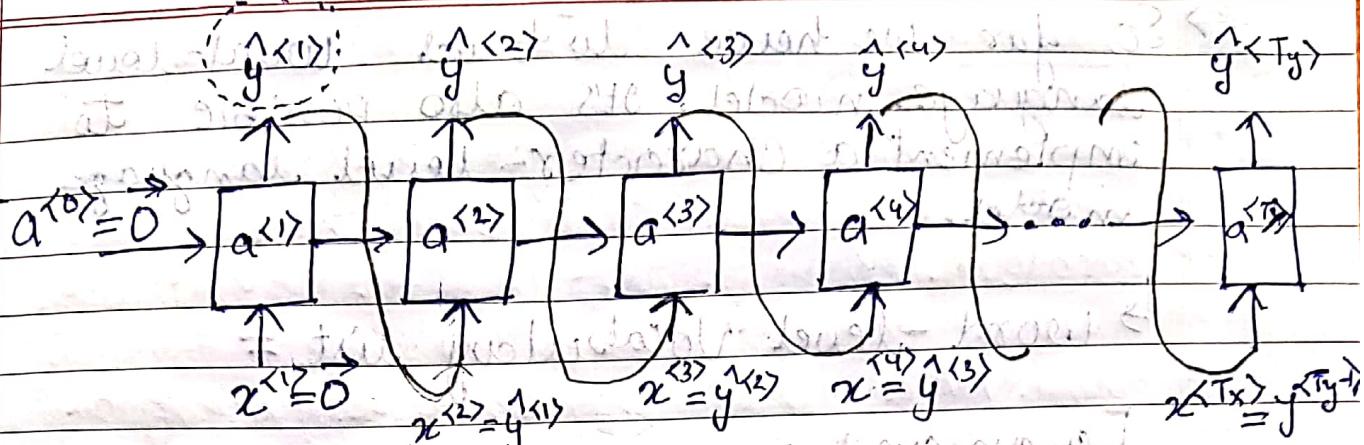
Let's see the steps of how we can sample a novel sequence from a trained sequence language model:

First of all, the network was trained using this structure



Now let's see how Novel Sampling is done on a trained language model.

Sampling



- We first pass $a^{<0>} = \vec{0}$, and $x^{<1>} = \vec{0}$
- Then we choose a prediction randomly from distribution obtained by $\hat{y}^{<1>}$. For example it could be anything, let's say "The".
- ↳ In numpy this can be implemented using: `numpy.random.choice(..)`
- ↳ This is the end where you get a random beginning of the sentence each time you sample run a novel sequence.
- We pass the last prediction word with the calculated $a^{<1>}$ next time step as input.
- We keep doing 3 & 4 steps for a fixed length or until we get the <EOS> token.
- You can reject <UNK> or <OOV> tokens if you mind finding it in your output.

→ So far we have to build word-level language model. It's also possible to implement a character-level language model.

↳ Word-level Vocabulary list =

[a, an, and, can, ..., z, the, ,]

↳ character-level Vocabulary list =

[a, b, c, d, ..., z, A, B, C, D, ..., Z, 1, 2, 3, ..., 9, 0]

→ In character-level language models, the vocabulary will contain [a-zA-Z0-9], punctuation, special characters and possibly tokens.

→ In character-level language model, we give input as: "CCITS Average15m. <EOS>
y¹y²y³y⁴y⁵ ... y^{T_y}"

→ character-level language model has some pros and cons compared to the word-level language model.

↳ Pros:

→ There will be no token; it can create any word in the training data.

→ Cons:

- The main disadvantage is that you end-up with much longer sequences.
- Character-level language models are not as good as word-level language models at capturing long-range dependencies between how the earlier parts of the sentence also affect the later part of the sentence.
- Also, more computationally expensive and harder to train.
- The trend Andrew Ng Sir has seen in NLP is that for the most part, a word-level language model is still used, but as computers get faster there are more and more applications where people are, at least in some special cases, starting to look at more character-level models. Also, they are used in specialized applications where you might need to deal with unknown words or other vocabulary words a lot. Or they are also used in more specialized application where you have a more specialized vocabulary.

* Vanishing gradients with RNNs

One of the problems with a naive RNNs that they run into Vanishing gradient problem.

- An RNN after processing a sequence of data with the size of 10000 time steps, has 10,000 depth layers which is very hard to optimize.

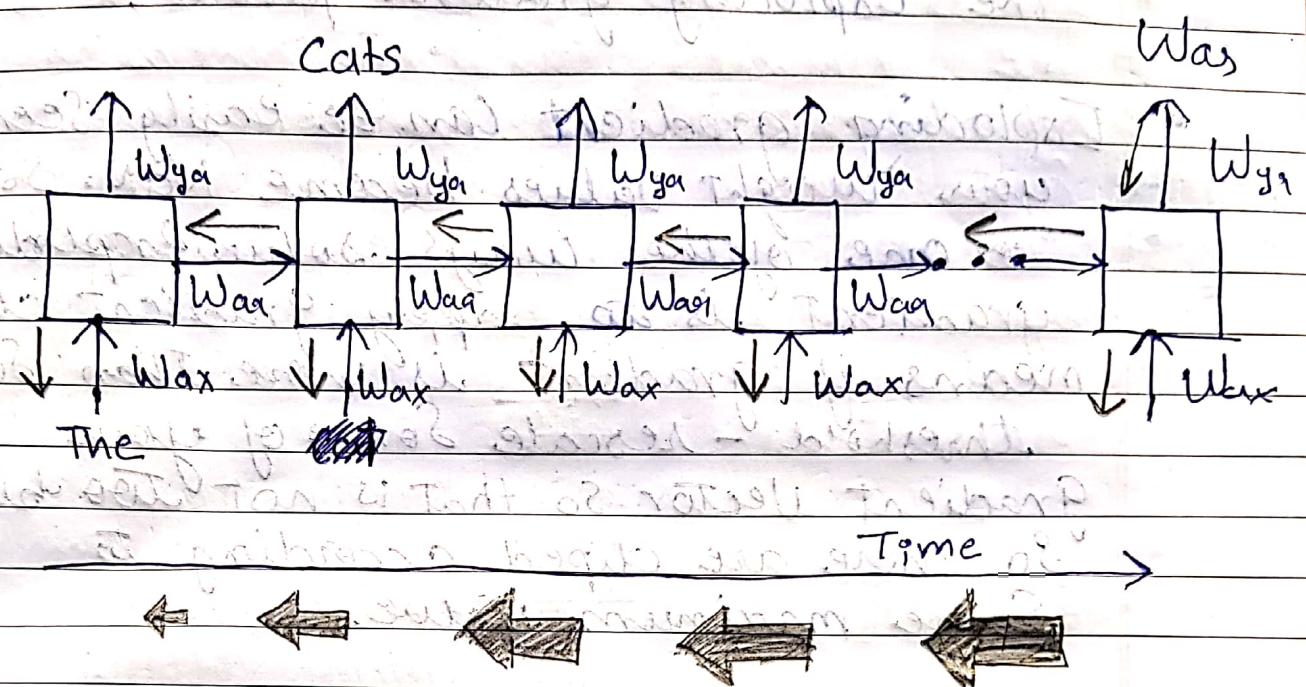
- Let's take an example. Suppose we are working with language modeling problems and there are two sequences that model tries to learn:

- "The cat which already ate . . . , was full"
- "The cats which already ate . . . , were full."
- Dots represent many words in between.

- What we need to learn here that "Was" came with "Cat" and that "Were" came with "Cats". The naive RNN is not very good at capturing very long-term dependencies like this.

- As we have discussed in Deep neural networks, deeper networks are getting into Vanishing gradient problem!

That also happens with RNNs with a long sequence size.

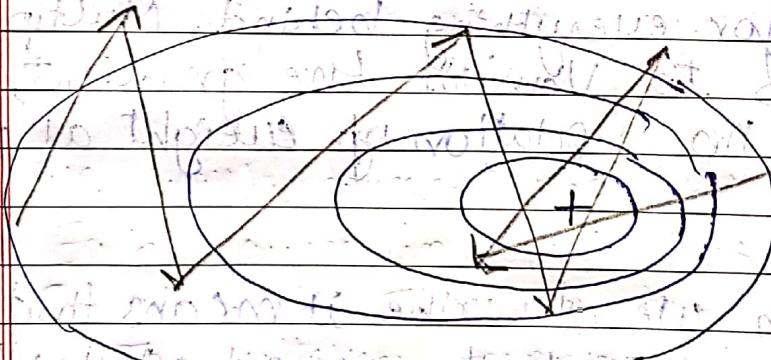


- For computing the word "Was", we need to compute the gradient for everything behind. Multiplying fraction tends to Vanish the gradient, which means no updation of weight at early time steps.
- In the problem we describe it means that its hard for the network to memorize "Was" word all over back to "cats". So in this case, the network won't identify the singular/plural words so that it gives it the right grammar form of verb Was/Were.
- The conclusion is the RNN are not good in long-term dependencies because some of the weight may not be update properly.

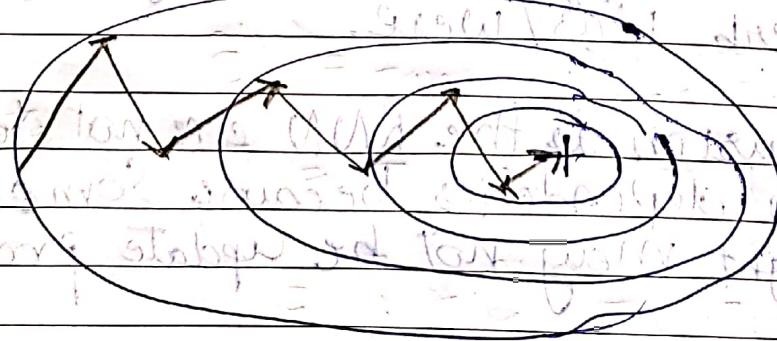
Vanishing gradient problem tends to be the bigger problem with RNNs than the exploding gradients problem.

Exploding gradient can be easily seen when your weight values become NaN. So one of the ways solve exploding gradient is to apply gradient clipping means if gradient is more than some threshold - rescale some of your gradient vector so that is not too big. So there are clipped according to some maximum value.

Without gradient clipping.



With gradient clipping



* Extra: (Q4) Fixing gradient problems

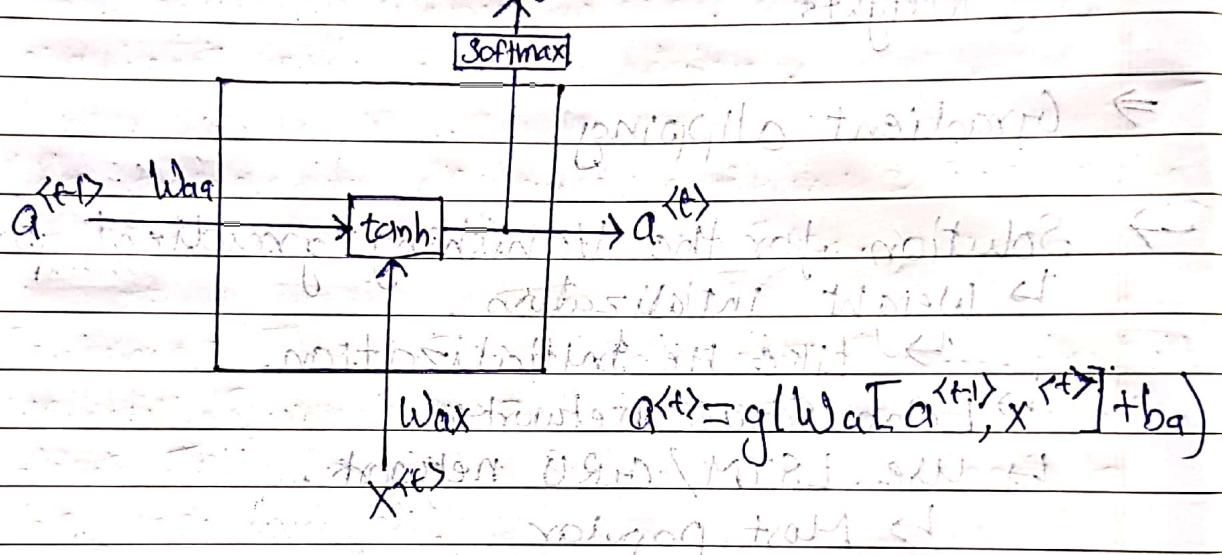
- Solution for the Exploding gradient problem:
 - ⇒ Truncated backpropagation.
 - ↳ Not to update all the weights in the way back.
 - ↳ Not optimal. You won't update all the weights.
- ⇒ Gradient clipping.
- Solution for the Vanishing gradient problem.
 - ↳ Weight initialization.
 - ↳ Like He Initialization.
 - ↳ Echo State networks.
 - ↳ Use LSTM / GRU network.
 - ↳ Most popular



Grated Recurrent Unit (GRU)

- GRU is an RNN type that can help solve the vanishing gradient problems and can remember the long-term dependencies.

- The basic RNN can be visualized to be like this:



- We will represent the GRU with similar drawings
- Each layer in GRUs has a new variable C which is the memory cell. it can tell to whether memorize something or not.
- In GRU, $C^{t+1} = a^{t+1}$.
- Equations

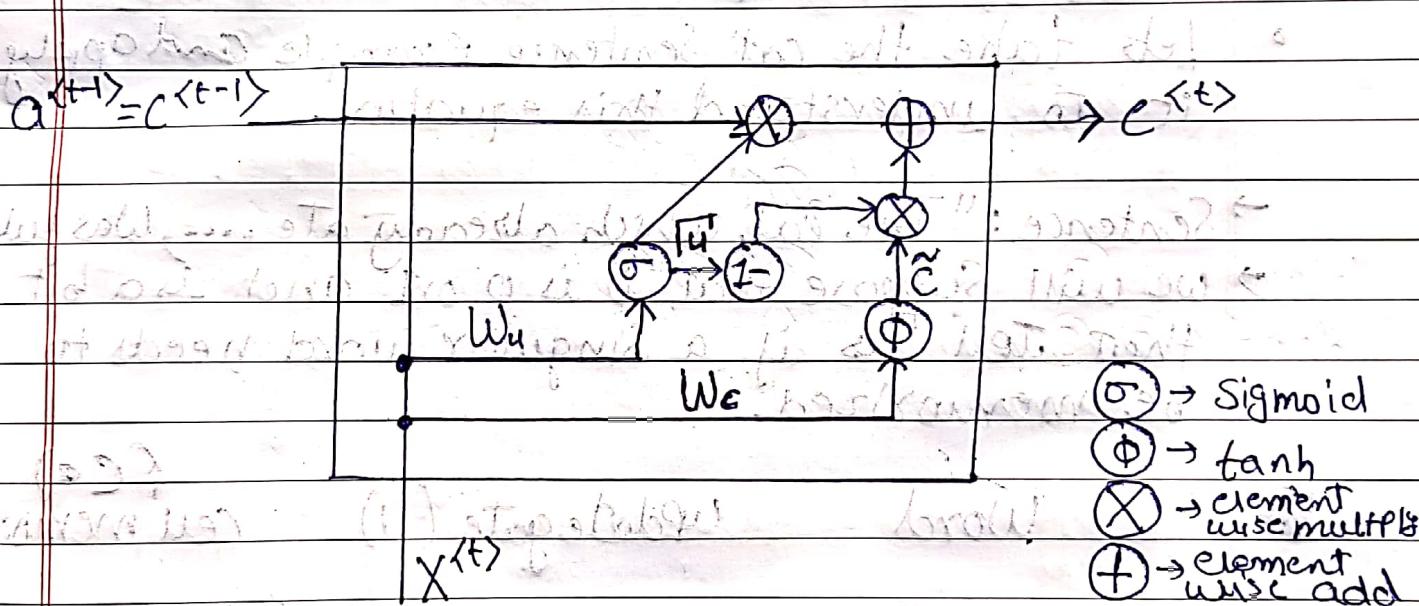
$$\tilde{C}^{t+1} = \tanh(W_c [C^{t-1}, x^{t+1}] + b_c)$$

update gate $\rightarrow \Gamma_u = \sigma(W_u [C^{t-1}, x^{t+1}] + b_u)$

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

- ↳ The update gate is between 0 and 1.
- ↳ To understand GRUs, imagine that the update gate is either 0 or 1 most of the time.

- So we update the memory cell based on update cell and the previous cell.
- Drawing of the GRU (Simplified).



• Because the update gate U_t is usually a small number like 0.00001, GRUs doesn't suffer the vanishing gradient problem.

↳ In the equation this makes $C^{(t)} = C^{(t-1)}$ in a lot of cases cases.

Note : Update gate \tilde{U} can also be written as U

*enliv

Date : _____

Page: _____

Shapes

- $a^{(t)}$ Shape is (No. of hidden Neurons)
- $C^{(t)}$ is the Same as $a^{(t)}$
- $C^{(t)}$ is the Same as $a^{(t)}$
- $U^{(t)}$ is also same dimension as of $a^{(t)}$

- The multiplications in the equation are element wise multiplication.
- Let's take the cat Sentence Example and apply it to understand this equation.
 - Sentence : "The cat, which already ate, was full"
 - we will suppose that U is 0 or 1 and is a bit that tell us if a Singular word needs to be memorized.

Word	Update.gate (U)	cell memory ($C^{(t)}$)
The	0	Val
Cat	1	new-value
which	0	New-Value
already	0	New-Value
....	0	New-Value
Was	1	newer-value
full	0	

⇒ What has been described so far is the Simplified GRU unit. Let's describe the full one:

- The full Gated GRU contains a new gate that is used with to calculate the candidate C . The gate tells you how relevant is C^{t-1} to C^t .

→ Equations :

$$\Gamma_u = \sigma(W_u[C^{t-1}, x^t] + b_u]$$

$$r = \sigma(W_r[C^{t-1}, x^t] + b_r]$$

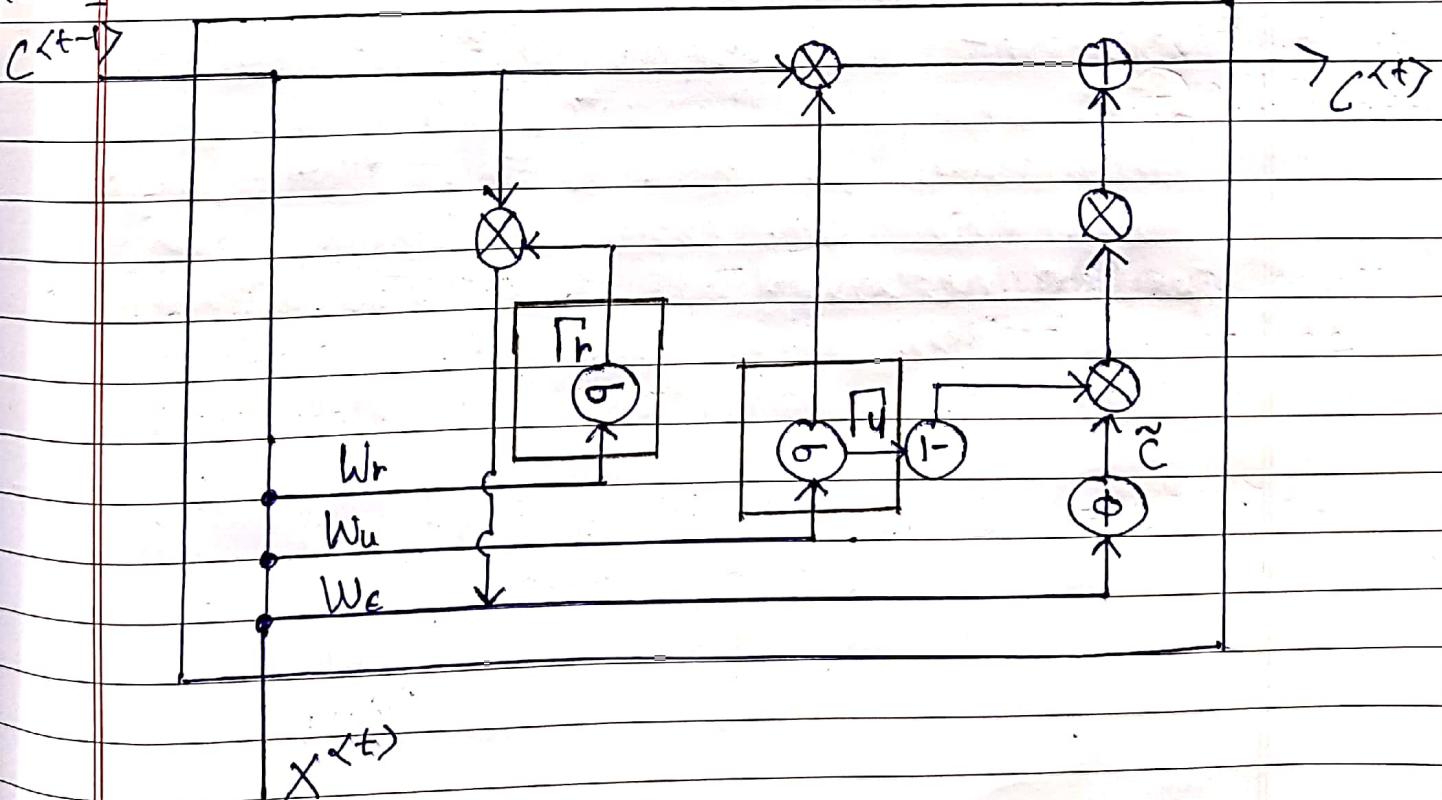
$$\tilde{C}^t = \tanh(W_c[r * C^{t-1}, x^t] + b_c)$$

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

→ Diagram

⑥ → Sigmoid, ⑦ → tanh, ⑧ → element wise multiplication, ⑨ → addition

$$a^{t-1} =$$



Shapes are the same. To understand this we can look at the diagram below.

So why we use these architectures why don't we change them, how we know they will work, why not add another gate, why not use the simpler GRU instead of the full GRU well researchers has experimented over years all the various types of these architectures with many different version and also addressing the vanishing gradient problem. They have found that full GRUs are one of the best RNN architecture to be used for many different problems. you can make your design but putting in mind that GRUs and LSTMs are standards.

