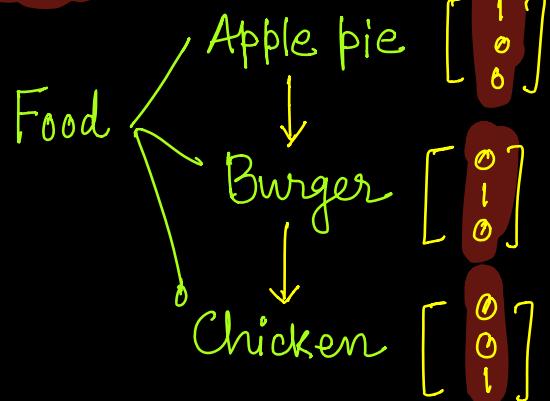
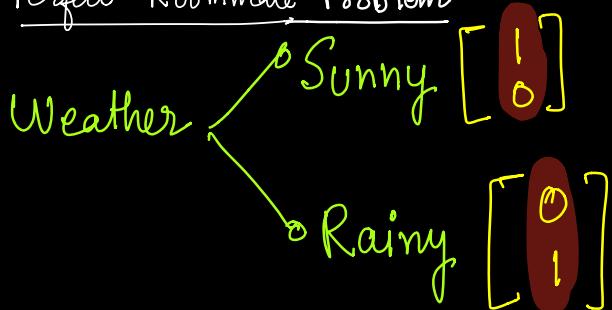


Part-A (RNN - Intuition) (Encoding)

Perfect Roommate Problem



Problem - 01 (Fixed rule)

⊗ Weather = Sunny \Rightarrow Apple pie

⊗ Weather = Rainy \Rightarrow Burger

O/P (food)



I/P (weather)

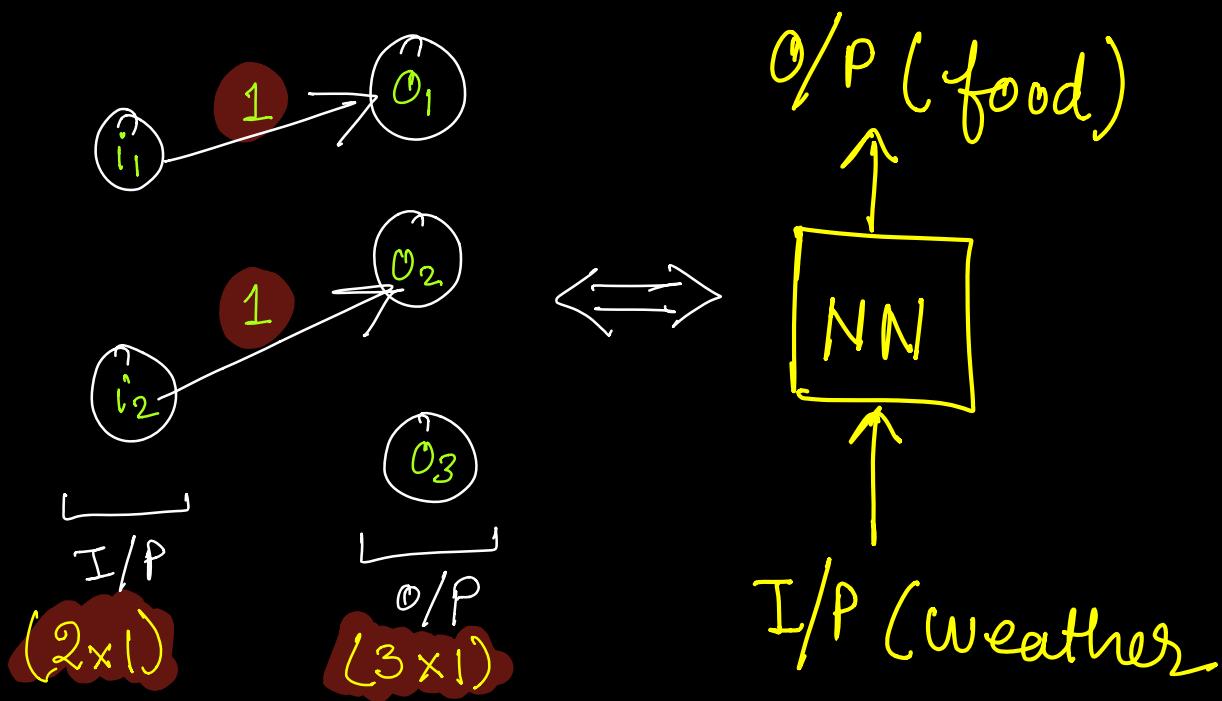
$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad R_1$$

Model Sunny Apple pie

i_1 i_2

$$O_1 \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad R_2$$

Model Rainy Burger



Problem - 02 (Round Robin Schedule)

→ Apple pie → Burger → chicken

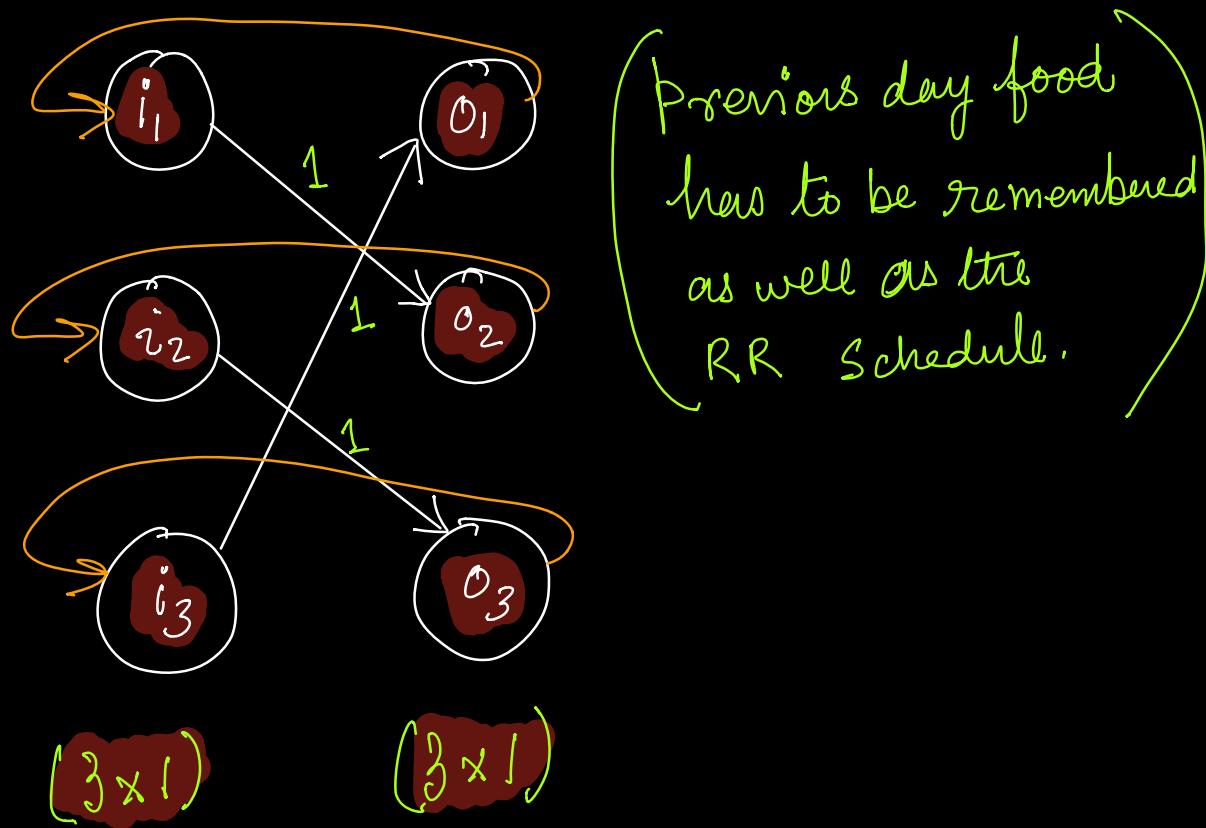
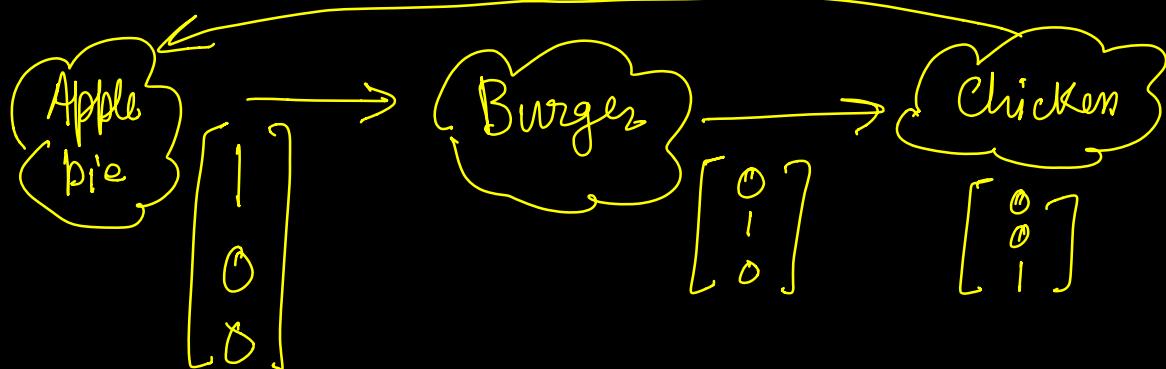
⊗ O/P has to be fed back as I/P

⊗ Network need to remember last food as cook accordingly.

Permutation matrix

$$\begin{matrix}
 o_1 & i_1 & i_2 & i_3 \\
 o_2 & & & \\
 o_3 & & &
 \end{matrix}
 \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}
 \begin{bmatrix} a \\ b \\ c \end{bmatrix} = \begin{bmatrix} c \\ a \\ b \end{bmatrix}$$

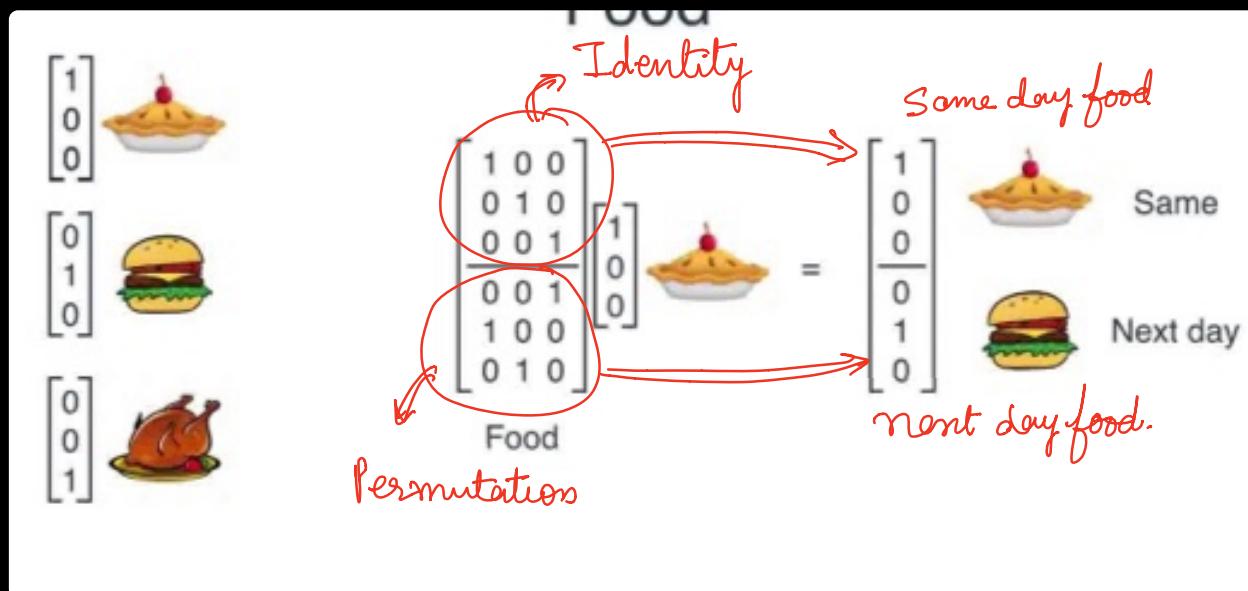
Above permutation matrix can be used as a model to realize the RR scheduling which takes a food as an i/p and gives next day food.



Problem-03 : ⚡ If Sunny, do not cook and just repeat the last day's left over food.

⚡ If Rainy, Cook the next day's food according to the RR scheduling.

Now along with the previous day food (for RR) one should have to check the weather. Depending upon the weather choose yesterday's food or today's food. Along with yesterday's food today's weather is also an i/p



$\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ 	$\begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}$ Weather $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ 	$=$ $\begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
--	---	---

Sunny selects same food.

$\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ 	$\begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}$ Weather $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$ 	$=$ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}$
--	---	---

Rainy selects next day food

If this is the i/P

 	$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ O/P after food Matrix	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}$ O/P after weather matrix	$=$ $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}$ Expected O/P
------	---	--	---

+
 Same
 Next day
 Additions

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

Food



$$+ \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

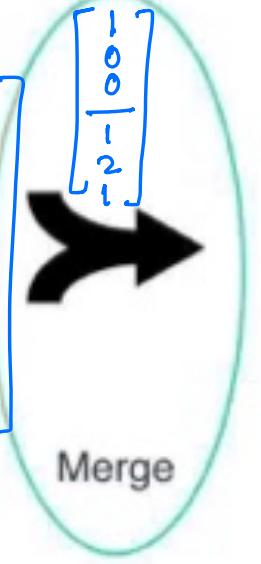
Add

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

Weather

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 2 \end{bmatrix}$$

Merge



$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 2 \\ 1 \end{bmatrix} \xrightarrow{\text{(non-linear)}} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Max operation

Merge

$$\begin{bmatrix} 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0+0 \\ 0+1 \\ 0+0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

Adding top & bottom.

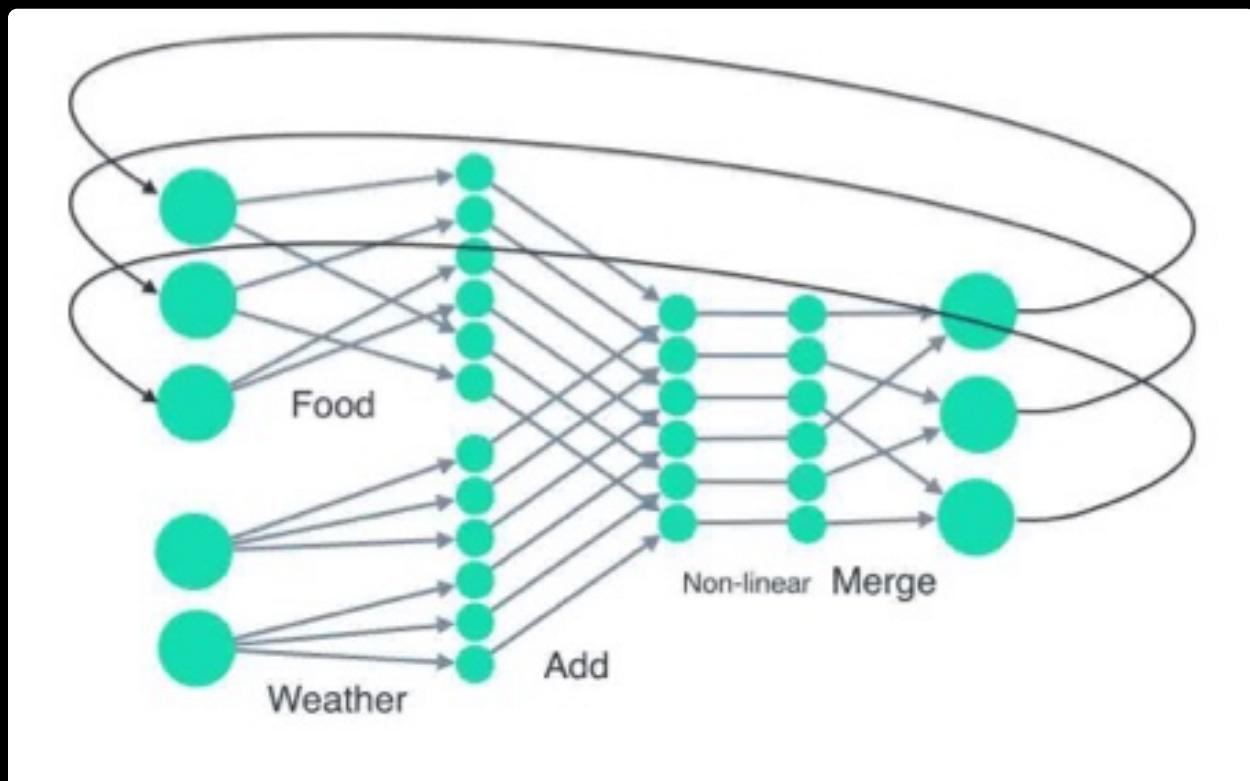
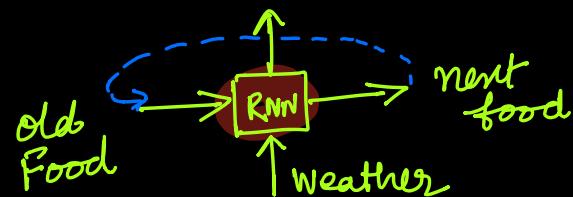
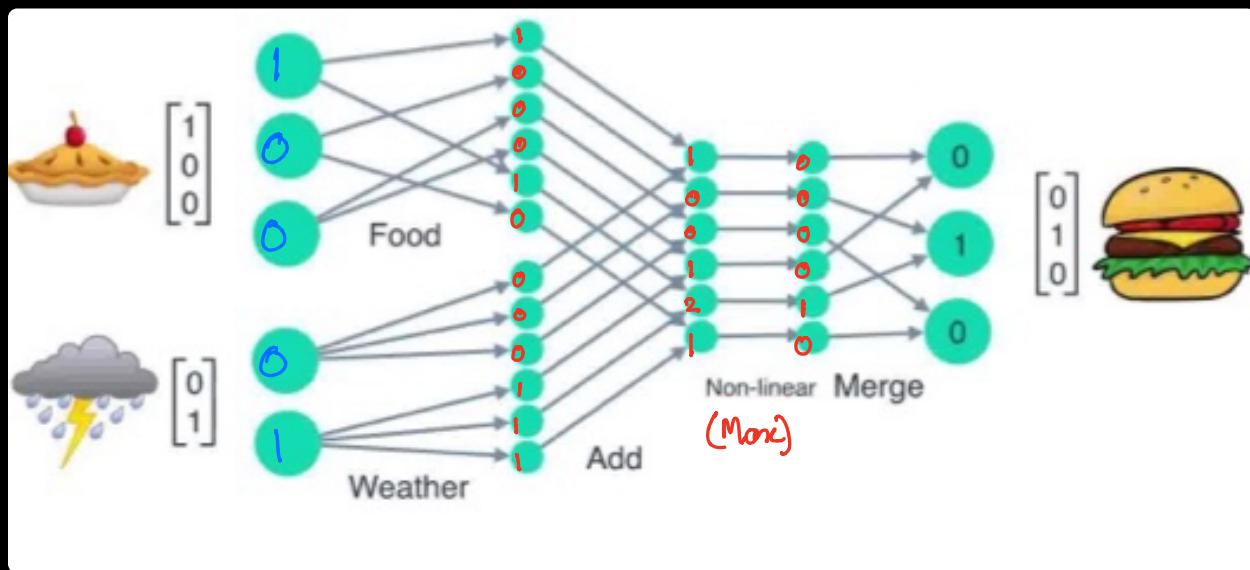
$$\begin{bmatrix} 1 & 0 & 0 & | & 1 & 0 & 0 \\ 0 & 1 & 0 & | & 0 & 1 & 0 \\ 0 & 0 & 1 & | & 0 & 0 & 1 \end{bmatrix}$$

Merge

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0+0 \\ 0+1 \\ 0+0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$





PART-B: Mathematical formulation of RNN

Name recognition in the text modules

useful for indexing
people in news article

(A) Given any sentence, find / label out the names. (T_x^{ent} labeling)

Say:

$X : \begin{matrix} X <1> & X <2> & \dots & X <t> & \dots & X <T_x> \\ \text{(1st word)} & \text{(2nd word)} & & & & \end{matrix}$
 played Shoib very well.

$Y : \begin{matrix} 1 & 0 & 1 & 0 & 0 \\ Y <1> & Y <2> & \dots & Y <t> & Y <T_y> \end{matrix}$

Here $T_x = T_y$ (Total length of (X) and (Y) Seq)
 (may or may not be same)

$X^i < t > \Rightarrow i^{\text{th}} \text{ i/p training example.}$

T_x^i and $T_y^i \Rightarrow$ Corresponding lengths.

(B) Word representation: (1 Hot encoding)

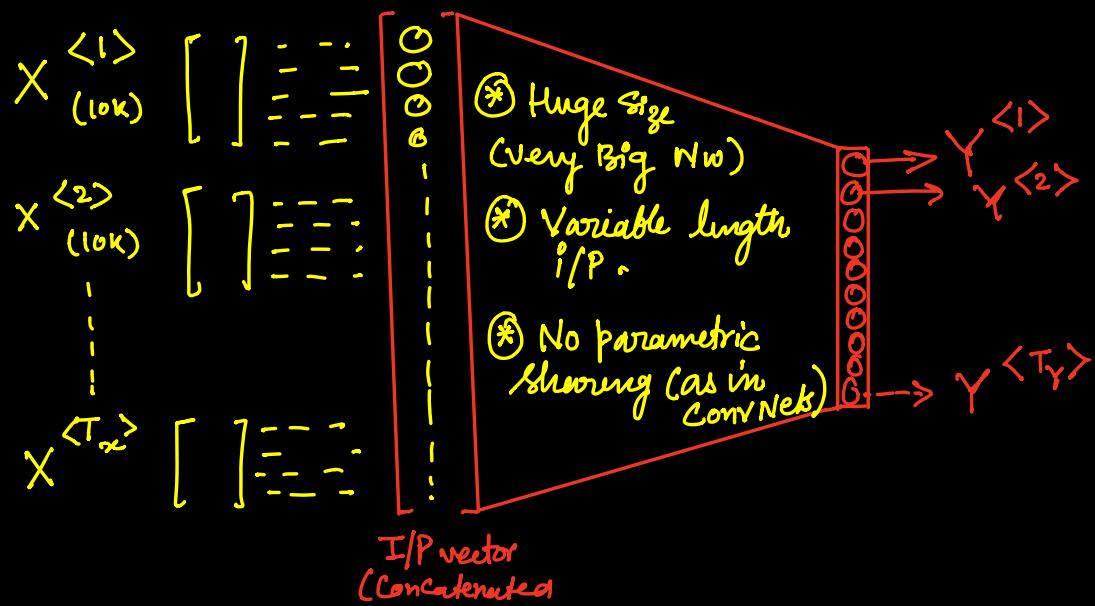
$\left[\begin{matrix} \text{word}_1 \\ \text{word}_2 \\ \text{word}_3 \\ \vdots \\ \text{word}_{10k} \end{matrix} \right]$

* make a dictionary of unique words in your training data.
 (say $10k$)

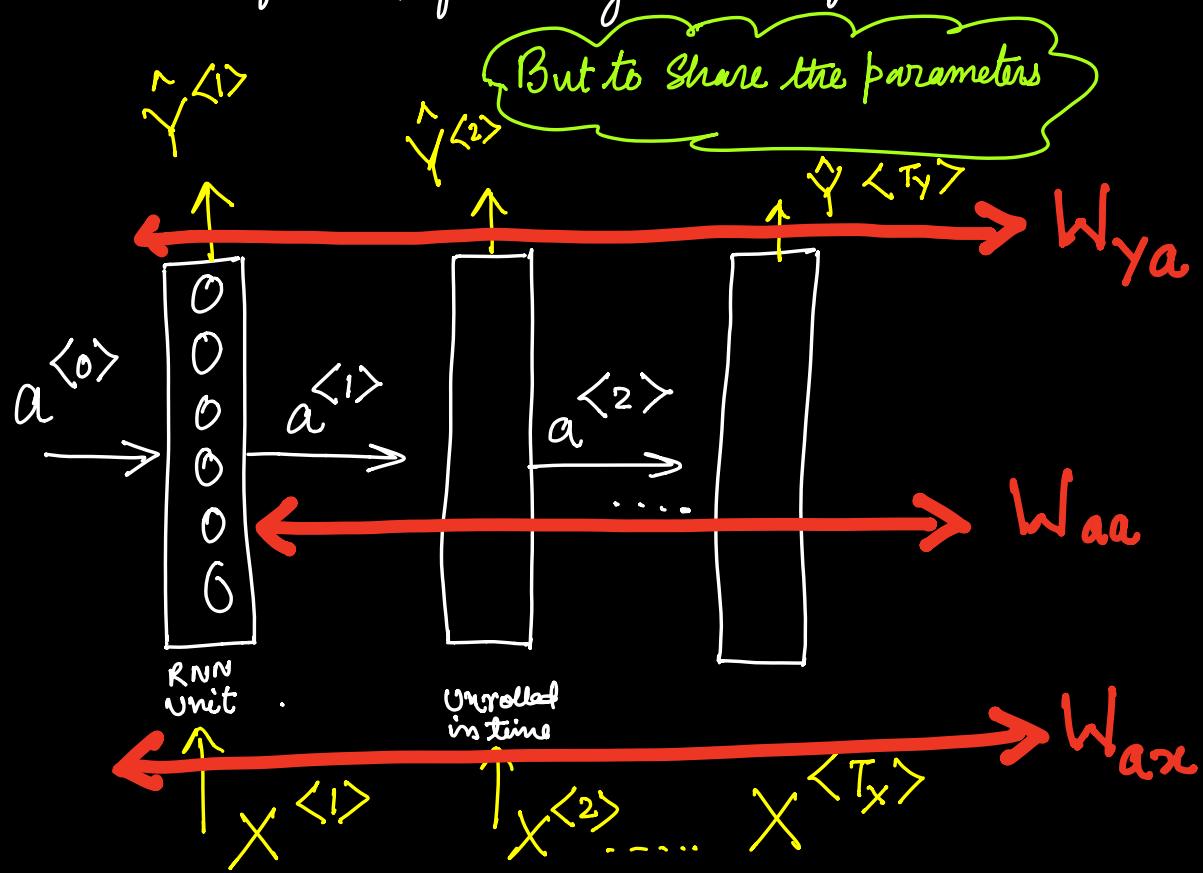
* Each word can be represented as a $10k$ dimensional 0/1 (binary vector) \Rightarrow 1-Hot encoding.

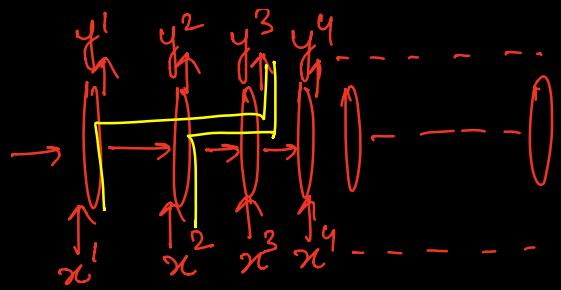
* Our goal to get a label (Y) for given (X).

Learning a mapping from $\underline{\underline{(X)}} \Rightarrow \underline{\underline{(Y)}}$



(*) We require our network to read the sentence word by word from left to right (Seeing data as a Time Seq)



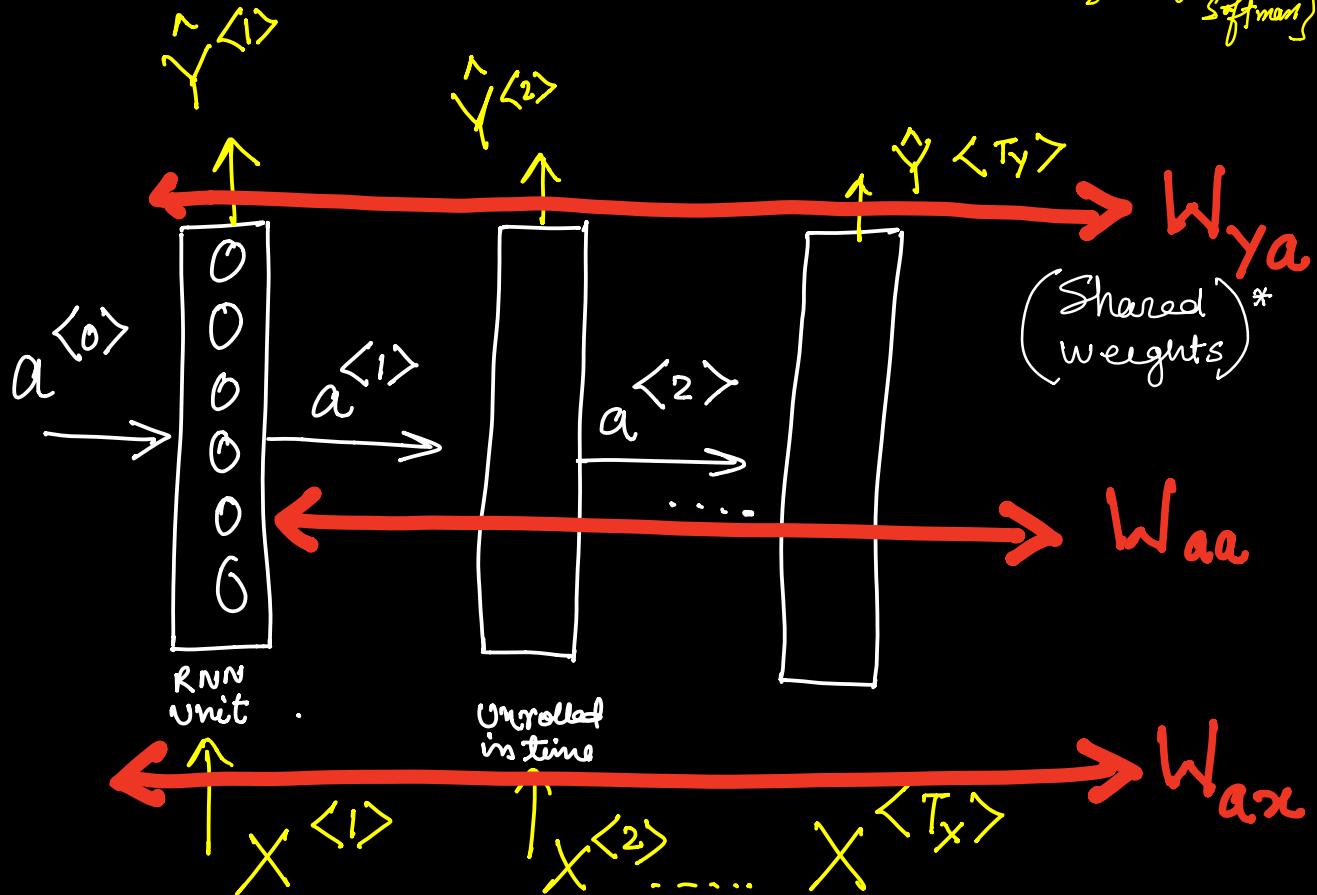


In order to predict $\hat{y}^{(3)}$ from $x^{(3)}$
 $x^{(1)}$ & $x^{(2)}$ they are
also utilized.

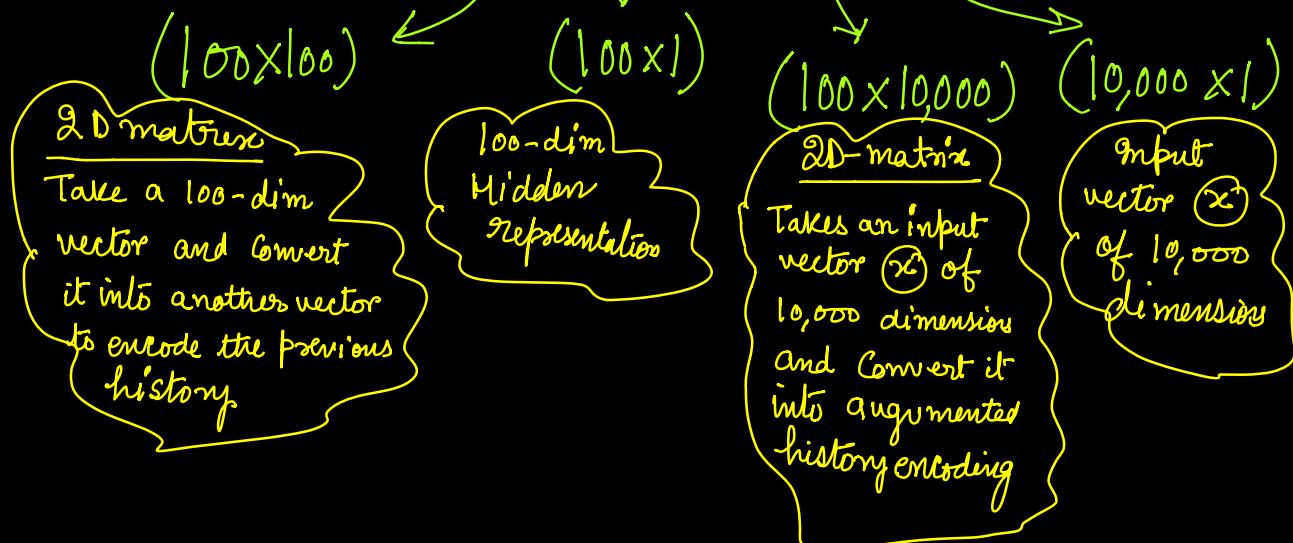
$$\textcircled{1} \quad a^{(0)} = \vec{0} \quad (\text{initialized})$$

$$\textcircled{2} \quad a^{(t)} = g_1 (w_{aa} a^{(t-1)} + w_{ax} x^{(t)} + b_a)$$

$$\textcircled{3} \quad \hat{y}^{(t)} = g_2 (w_{ya} a^{(t)} + b_y) \quad \begin{cases} g_1 \equiv \tanh / \text{ReLU} \\ g_2 \equiv \text{Sigmoid} / \text{Softmax} \end{cases}$$



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



Matrix Augmentation :

$$\textcircled{a} \quad \left[\begin{array}{c|c} W_{aa} & W_{ax} \end{array} \right] = W_a \quad \begin{matrix} (\text{Weight Matrix}) \\ 100 \times 100 \\ \text{Horizontal Concat} \end{matrix}$$

$$\textcircled{b} \quad \left[\begin{array}{c} a^{(t-1)} \\ \hline x^{(t)} \end{array} \right] = \left[\begin{array}{c} a^{(t-1)} \\ x^{(t)} \end{array} \right] \quad \begin{matrix} 100 \times 1 \\ \text{Vertical Concat} \end{matrix}$$

$$\therefore W_a [a^{(t-1)}, x^{(t)}] \Rightarrow [100 \times 1]$$

$$[I|I] \left[\begin{array}{c} a \\ b \end{array} \right] = [a+b]$$

Therefore our equations :

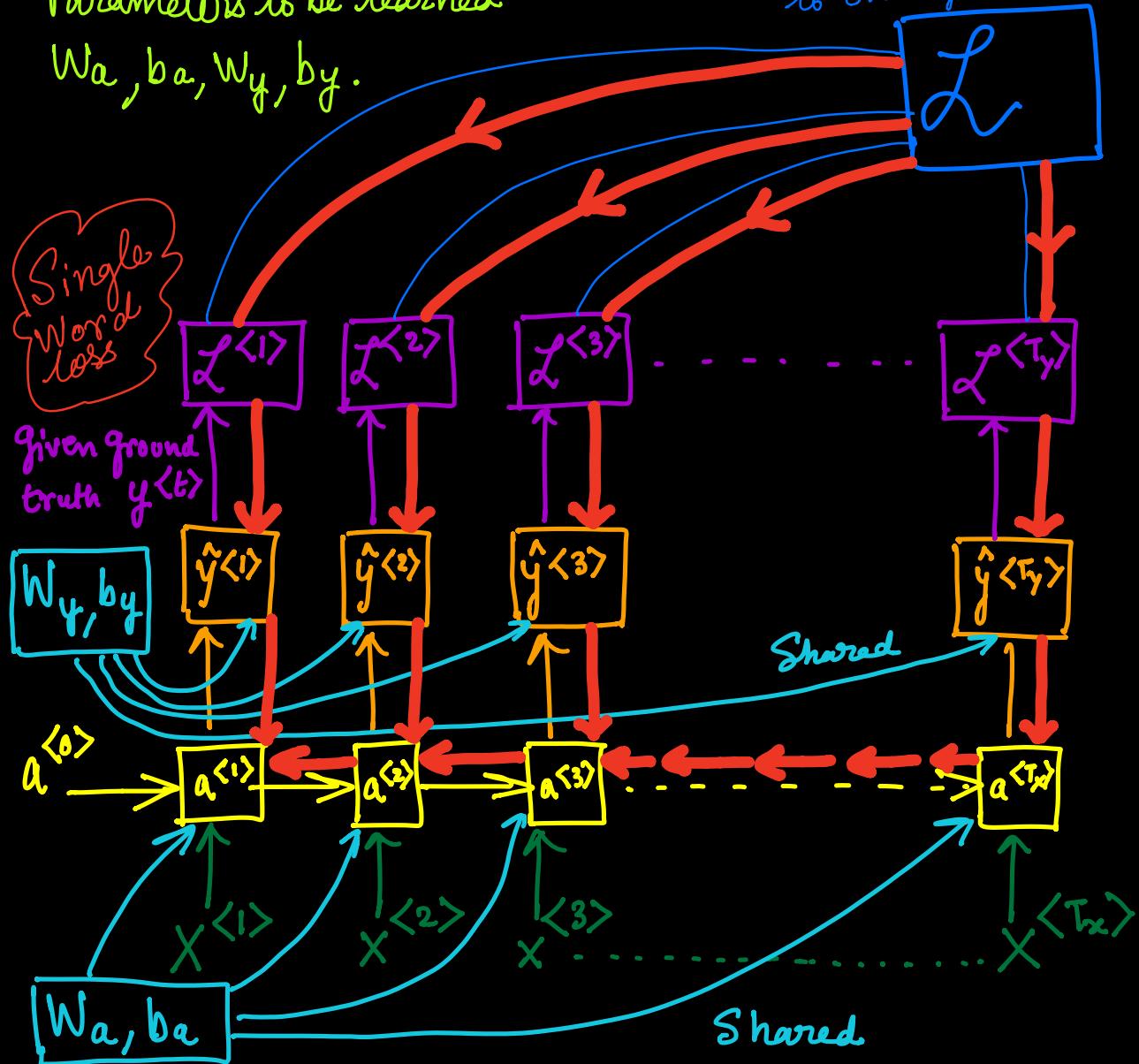
$$a^{(t)} = g_1 \left(W_a [a^{(t-1)}, x^{(t)}] + b_a \right)$$

$$\hat{y}^{(t)} = g_2 (W_y a^{(t)} + b_y)$$

Parameters to be learned

W_a, b_a, W_y, b_y .

Total loss associated
to training sample.

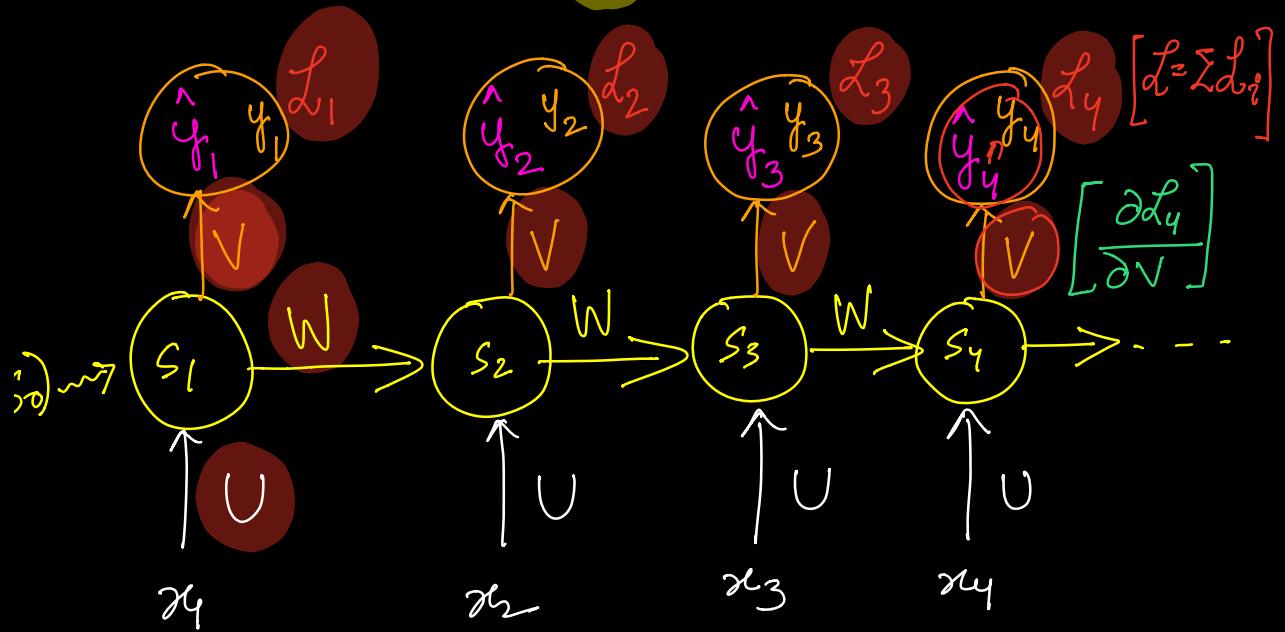


① Word level loss

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

② Sentence level loss

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



Part C : Vanishing Gradient (RNN backprop)

$$\mathcal{L} = \mathcal{L}_1 + \mathcal{L}_2 + \mathcal{L}_3 + \mathcal{L}_4 \quad \text{where } \left[\mathcal{L}_i = -\log(y_i^{t_e}) \right]$$

Predicted probability
of true character at ($t=i$)

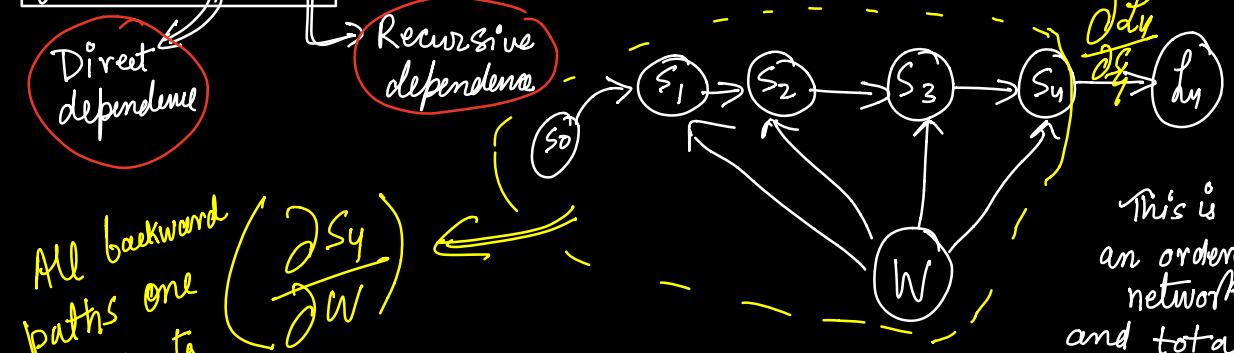
$$\left(\frac{\partial \mathcal{L}}{\partial V/W/U} \right) = \sum_{k=1}^4 \frac{\mathcal{L}_k}{\partial V/W/U}$$

as $\frac{\partial \mathcal{L}_k}{\partial V} = \frac{\partial f(g(v), \hat{y}_k)}{\partial V} = \frac{\partial f(g(v), \hat{y}_3)}{\partial V}$ [direct computation]

See Computation Graph for (W)

$$\frac{\partial \mathcal{L}_4}{\partial W} = \frac{\partial \mathcal{L}_4}{\partial S_4} * \frac{\partial S_4}{\partial W} \quad (\text{But, } S_4 \text{ is recursively dependent on } (S_3 \dots S_1))$$

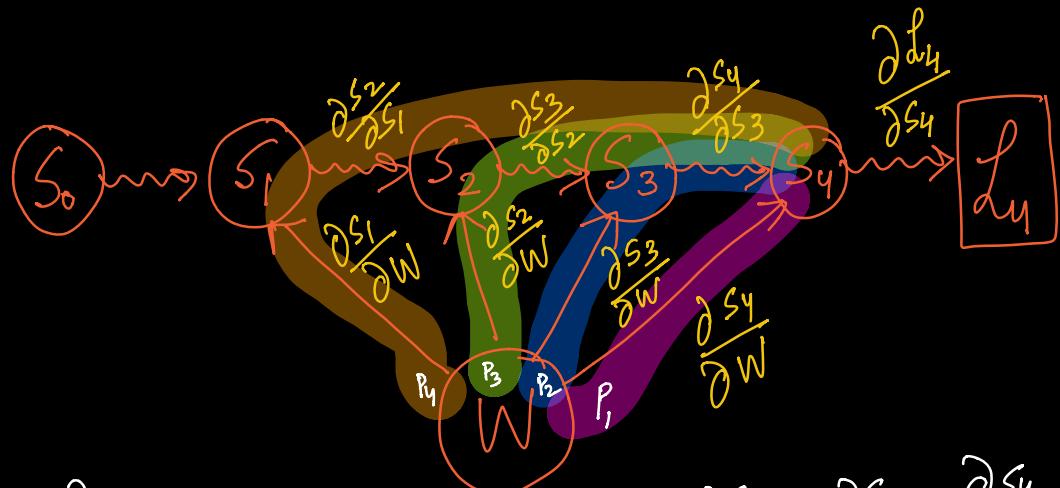
simply treating $(S_3 \dots S_1)$ const and differentiate about W .



All backward paths one can take to reach (W) from (S_4)

To Compute
 $\left[\frac{\partial S_4}{\partial W} \right] \downarrow$

This is an ordered network and total derivative need to be computed after adding all the possible paths.



$$\therefore \frac{\partial S_4}{\partial W} = \text{Adding all Paths derivatives} = \frac{\partial S_4}{\partial S_4} + \frac{\partial S_4}{\partial S_3} * \frac{\partial S_3}{\partial W} + \underbrace{\frac{\partial S_4}{\partial S_3} * \frac{\partial S_3}{\partial S_2} * \frac{\partial S_2}{\partial W}}_{\frac{\partial S_4}{\partial S_2}} + \underbrace{\frac{\partial S_4}{\partial S_3} * \frac{\partial S_3}{\partial S_2} * \frac{\partial S_2}{\partial S_1} * \frac{\partial S_1}{\partial W}}_{\frac{\partial S_4}{\partial S_1}}$$

$$\frac{\partial S_4}{\partial W} = \sum_{K=1}^4 \frac{\partial S_4}{\partial S_K} * \frac{\partial S_K}{\partial W}$$

$$\therefore \frac{\partial L_y}{\partial W} = \frac{\partial L_y}{\partial S_4} * \frac{\partial S_4}{\partial W} = \frac{\partial L_y}{\partial S_4} * \left[\sum_{K=1}^4 \frac{\partial S_4}{\partial S_K} * \frac{\partial S_K}{\partial W} \right]$$

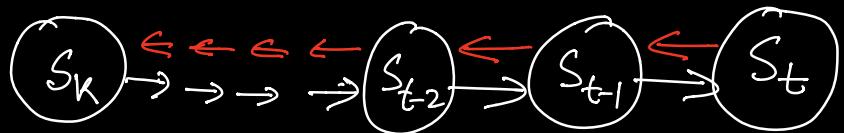
$$\frac{\partial L}{\partial W} = \frac{\partial L_1}{\partial W} + \frac{\partial L_2}{\partial W} + \frac{\partial L_3}{\partial W} + \frac{\partial L_4}{\partial W}$$

Bounded

$S_i = f(W, S_{i-1})$
 $\therefore \text{related to } [W]$

This is called BPTT

This algorithm back propagation in times need to compute a chained differentiation which requires recursion and repetitive multiplication. $\left(\frac{\partial S_t}{\partial S_k} \right)$



$$\frac{\partial S_t}{\partial S_k} = \frac{\partial S_t}{\partial S_{t-1}} * \frac{\partial S_{t-1}}{\partial S_{t-2}} * \dots * \frac{\partial S_{k+1}}{\partial S_k} = \prod_{j=k}^{t-1} \frac{\partial S_{j+1}}{\partial S_j}$$

Since these (S_i 's) are activated o/p hence are bounded therefore their derivatives are also bounded

$$(\text{Sigmoid})' \leq \frac{1}{4}$$

$$(\tanh)' \leq 1$$

$$\therefore \frac{\partial \mathcal{L}}{\partial W} \text{ (may explode/Vanish)}$$

① Word level loss

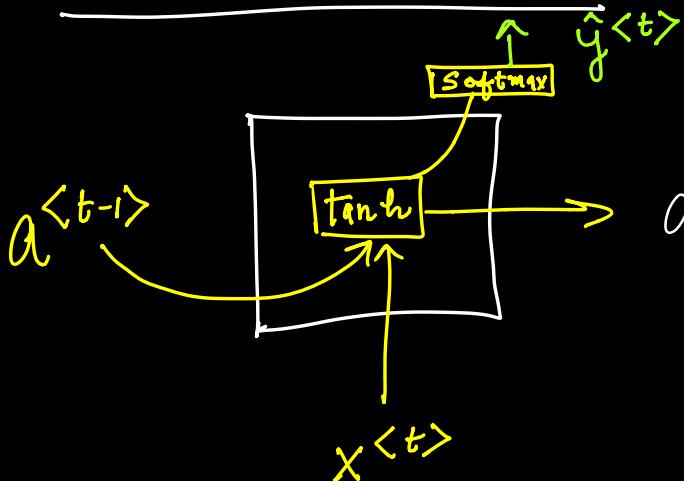
$$\mathcal{L}^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1-y^{(t)}) \log (1-\hat{y}^{(t)})$$

② Sentence level loss

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

Major Problem RNN faces is Vanishing Gradients and as the backpropagation is in time it fails to update well with respect to some old time.

Vanilla RNN Block diagram



This causes RNN to forget about long term events and their effect.

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

PART-D: Gated Recurrent Unit (GRU)

* Virat is a good player. He is the Captain
dependency

* Three persons are going with us and all of them are long terms dependency.

Problem is that how to remember that Subject is
 { Singular / Plural } // { Male / female }
 was / were he / she

\Rightarrow We consider our hidden state memory cells to understand such dependencies and tract/keep on memorizing them till the point they are required or something more important need to be remembered.

$$\Gamma_u = 1 \quad \Gamma_u = 0 \quad \text{---} \underset{\text{Keep on running}}{\text{---}} \quad \Gamma_u = 0 \quad \Gamma_u = 0 \Rightarrow (\text{When to toggle})$$

Virat is a good player. He is the Captain

dependency

$C^{t=1}$ $C^{t=1} - \dots - C^{t=1}$

(Some bit value got set)

Now you may update

Remember he is MALE (what to keep)

We need to remember Subjects Gender or it's Singular or plural.

$\begin{cases} \text{male} \\ \text{female} \end{cases}$

* Here (C^t) \Rightarrow It can be seen as the cell value at time t .

(Γ_u) Γ_u \Rightarrow It can see as the gate [allows/block] Some information

1	→	✓
0	→	
1	→	✓
0	→	
0	→	
0	→	



Element wise multiplication will help to produce the effect of gating.

Our I/p's are at t

$$(\Gamma_u^t) (C^t)$$

$$\boxed{C^{<t-1>}} \text{ and } \boxed{X^{<t>}}$$

New Cell State needs to be computed

a) Candidate update value

$$\tilde{C}^{<t>} = \tanh \left(W_c [C^{<t-1>} ; X^{<t>}] + b_c \right)$$

Neural network parameterized by (W_c, b_c) with input as $C^{<t-1>}$ and $X^{<t>}$ [what to update]

b) Update gate Γ_u (when to update)

$$\Gamma_u = \sigma \left(W_u [C^{<t-1>} ; X^{<t>}] + b_u \right)$$

Neural network parameterized by (W_u, b_u) with input as $C^{<t-1>}$ and $X^{<t>}$

This network tries to learn whether we need to update the current cell state or not.

Γ_u allows flexible cell updation as it is activated using Sigmoid fn which is always in b/w $\{0, 1\}$

c) Cell update :

Simplified GRU

$$C^{<t>} = \underbrace{\Gamma_u \odot C^{<t-1>}}_{\text{What needs to be updated}} + \underbrace{(1 - \Gamma_u) C^{<t-1>}}_{\text{How much to retain}}$$

Γ_u is 1 \Rightarrow signifies that to update the cell content by $C^{<t>}$ else keep the same.

But update governed only by one variable Γ_u

$C^{<t>}, C^{<t>}, \Gamma_u$ all are of same dimensions
say (100×1) then cell update
equation is element wise multiplication

(Γ_u) understands few key cell states bits behaviours (what they are learning)
and just try to gate their o/p effectively to as reduce the loss.

Full GRU:

$$\textcircled{1} \quad C^{<t>} = \tanh(W_c [\Gamma_r \odot C^{<t-1>} \times^{<t>}] + b_c)$$

$$\textcircled{2} \quad \Gamma_u = \sigma(W_u [C^{<t-1>} \times^{<t>}] + b_u)$$

$$\textcircled{3} \quad \Gamma_r = \sigma(W_r [C^{<t-1>} \times^{<t>}] + b_r)$$

Relevance gate, another neural network parameterized by (W_r, b_r) with input as $(C^{<t-1>} \times^{<t>})$

$$\textcircled{4} \quad C^{<t>} = \Gamma_u \odot C^{<t>} + (1 - \Gamma_u) \odot C^{<t-1>}$$

after several updates handle longer range dependencies, vanishing gradient and better convergence. $[a^{<t>} = C^{<t>}]$

How relevant any bit of $C^{<t-1>}$ can be
Another neural network for gating

PART-E: GRU to LSTM

GRU

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

LSTM

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

In LSTM (no) relevance gate

update gate (GRU)

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

$C^{<t>}$ and $C^{<t-1>}$ are weighted by r_u .

Signifies which bits requires update remaining bits are not changed/modifed.

Update gate (LSTM)

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

$C^{<t>}$
will only be weighted

This signifies which bits need to be update.

Relevance gate (GRU) absent in LSTM

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

Forget gate (LSTM) ~~absent in LSTM GRU~~

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

LSTM will learn a forget gate separately. Basically how we need to forget from c^{t-1} (previous state)

Output gate (LSTM)

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

~~absent in GRU~~
 $a^t = c^t$

c^t gated by Γ_o is the final o/p

Updation (GRU)

$$c^t = \Gamma_u \odot \tilde{c}^t + (1 - \Gamma_u) \odot c^{t-1}$$

Updation (LSTM)

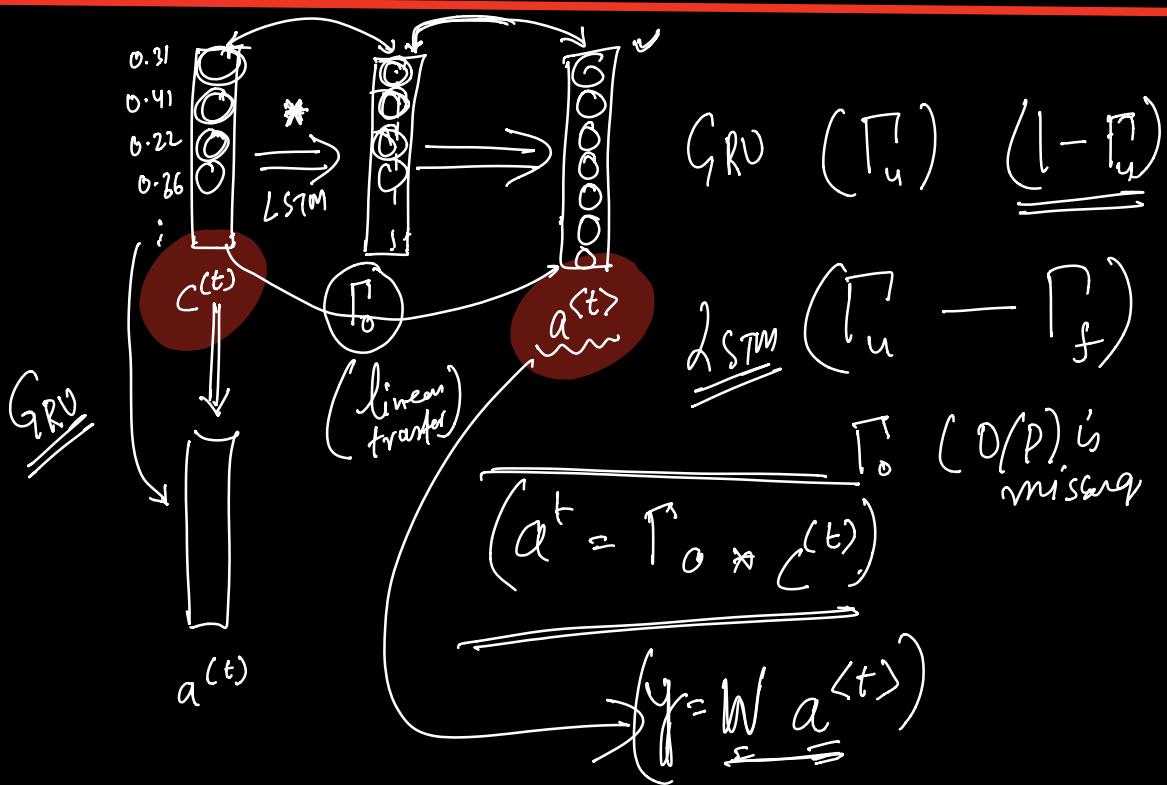
$$c^t = \Gamma_u \odot \tilde{c}^t + \Gamma_f \odot c^{t-1}$$

final o/p GRU

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

final o/p LSTM

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



Sequence Prediction Problems

PART F : Review

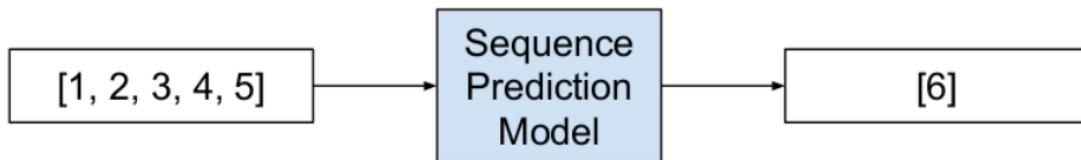
- Sequence Prediction.
- Sequence Classification.
- Sequence Generation.
- Sequence-to-Sequence Prediction.

Sequence Prediction.

Sequence prediction involves predicting the next value for a given input sequence. For example:

Input Sequence: 1, 2, 3, 4, 5
Output Sequence: 6

Listing 1.1: Example of a sequence prediction problem.



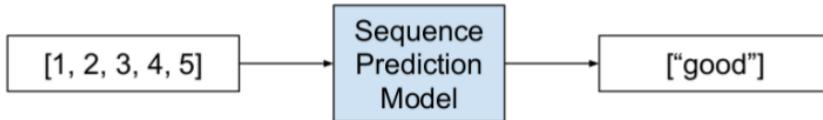
- Weather Forecasting. Given a sequence of observations about the weather over time, predict the expected weather tomorrow.
- Stock Market Prediction. Given a sequence of movements of a security over time, predict the next movement of the security.
- Product Recommendation. Given a sequence of past purchases for a customer, predict the next purchase for a customer.

Sequence Classification

Sequence classification involves predicting a class label for a given input sequence. For example:

```
Input Sequence: 1, 2, 3, 4, 5  
Output Sequence: "good"
```

Listing 1.2: Example of a sequence classification problem.



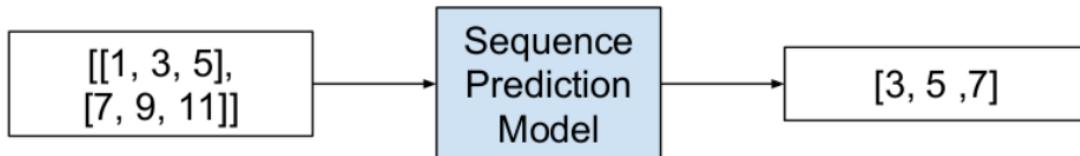
- DNA Sequence Classification. Given a DNA sequence of A, C, G, and T values, predict whether the sequence is for a coding or non-coding region.
- Anomaly Detection. Given a sequence of observations, predict whether the sequence is anomalous or not.
- Sentiment Analysis. Given a sequence of text such as a review or a tweet, predict whether the sentiment of the text is positive or negative.

Sequence Generation.

Sequence generation involves generating a new output sequence that has the same general characteristics as other sequences in the corpus. For example:

```
Input Sequence: [1, 3, 5], [7, 9, 11]  
Output Sequence: [3, 5 ,7]
```

Listing 1.3: Example of a sequence generation problem.



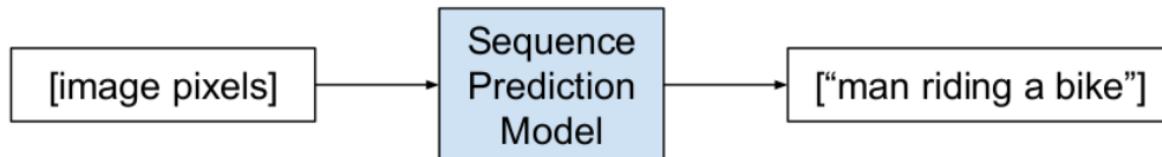
- Text Generation. Given a corpus of text, such as the works of Shakespeare, generate new sentences or paragraphs of text that read they could have been drawn from the corpus.
- Handwriting Prediction. Given a corpus of handwriting examples, generate handwriting for new phrases that has the properties of handwriting in the corpus.
- Image Caption Generation. Given an image as input, generate a sequence of words that describe an image.

Sequence-to-Sequence Prediction

For example:

```
Input Sequence: [image pixels]
Output Sequence: ["man riding a bike"]
```

Listing 1.4: Example of a sequence generation problem.



- Multi-Step Time Series Forecasting. Given a time series of observations, predict a sequence of observations for a range of future time steps.
- Text Summarization. Given a document of text, predict a shorter sequence of text that describes the salient parts of the source document.
- Program Execution. Given the textual description program or mathematical equation predict the sequence of characters that describes the correct output.

Limitations of Multilayer Perceptrons

- Stateless. MLPs learn a fixed function approximation. Any outputs that are conditional on the context of the input sequence must be generalized and frozen into the network weights.
- Unaware of Temporal Structure. Time steps are modeled as input features, meaning that network has no explicit handling or understanding of the temporal structure or order between observations.
- Messy Scaling. For problems that require modeling multiple parallel input sequences, the number of input features increases as a factor of the size of the sliding window without any explicit separation of time steps of series.
- Fixed Sized Inputs. The size of the sliding window is fixed and must be imposed on all inputs to the network.
- Fixed Sized Outputs. The size of the output is also fixed and any outputs that do not conform must be forced.

Promise of Recurrent Neural Networks

- Recurrent Neural Networks, or RNNs for short, are a special type of neural network designed for sequence problems.
- Given a standard feedforward MLP network, an RNN can be thought of as the addition of loops to the architecture.
- For example, in a given layer, each neuron may pass its signal laterally (sideways) in addition to forward to the next layer.
- The output of the network may feedback as an input to the network with the next input vector. And so on.
- The recurrent connections add state or memory to the network and allow it to learn and harness the ordered nature of observations within input sequences.

Introduction to LSTMs

- The Long Short-Term Memory, or LSTM, network is a type of Recurrent Neural Network.
- Long Short-Term Memory (LSTM) is able to solve many time series tasks unsolvable by feedforward networks using fixed size time windows.
- The addition of sequence is a new dimension to the function being approximated.
- Instead of mapping inputs to outputs alone, the network is capable of learning a mapping function for the inputs over time to an output.
- The internal memory can mean outputs are conditional on the recent context in the input sequence, not what has just been presented as input to the network.
- In a sense, this capability unlocks sequence prediction for neural networks.

RNN vs LSTM

- Like RNNs, the LSTMs have recurrent connections so that the state from previous activations of the neuron from the previous time step is used as context for formulating an output.
- But unlike other RNNs, the LSTM has a unique formulation that allows it to avoid the problems that prevent the training and scaling of other RNNs.
- The key technical historical challenge faced with RNNs is how to train them effectively.
- Experiments show how difficult this was where the weight update procedure resulted in weight changes that quickly became so small as to have no effect (vanishing gradients) or so large as to result in very large changes or even overflow (exploding gradients).
- LSTMs overcome this challenge by design.

LSTM architecture

- The computational unit of the LSTM network is called the memory cell, memory block, or just cell for short.
- The term neuron as the computational unit is so ingrained when describing MLPs that it too is often used to refer to the LSTM memory cell.
- LSTM cells are comprised of weights and gates.

LSTM Weights

A memory cell has weight parameters for the input, output, as well as an internal state that is built up through exposure to input time steps.

- Input Weights. Used to weight input for the current time step.
- Output Weights. Used to weight the output from the last time step.
- Internal State. Internal state used in the calculation of the output for this time step.

LSTM Gates

The key to the memory cell are the gates. These too are weighted functions that further govern the information flow in the cell. There are three gates:

- Forget Gate: Decides what information to discard from the cell.
- Input Gate: Decides which values from the input to update the memory state.
- Output Gate: Decides what to output based on input and the memory of the cell.

The forget gate and input gate are used in the updating of the internal state. The output gate is a final limiter on what the cell actually outputs.