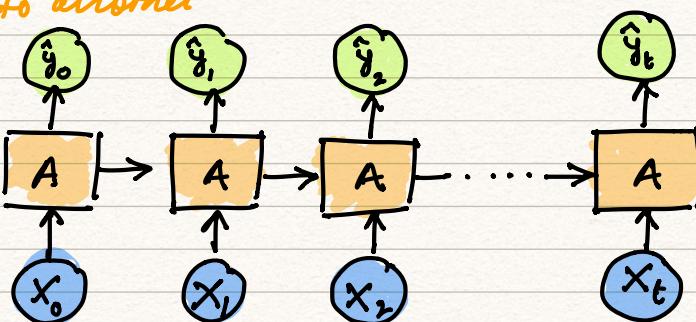
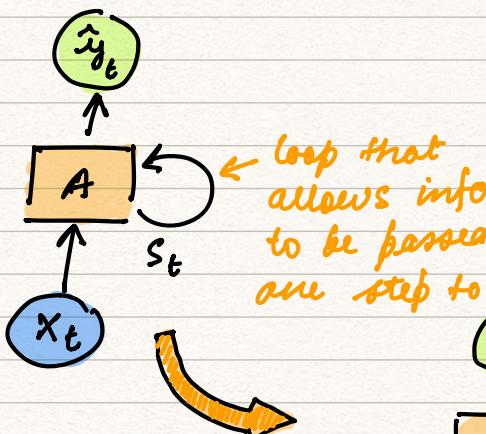
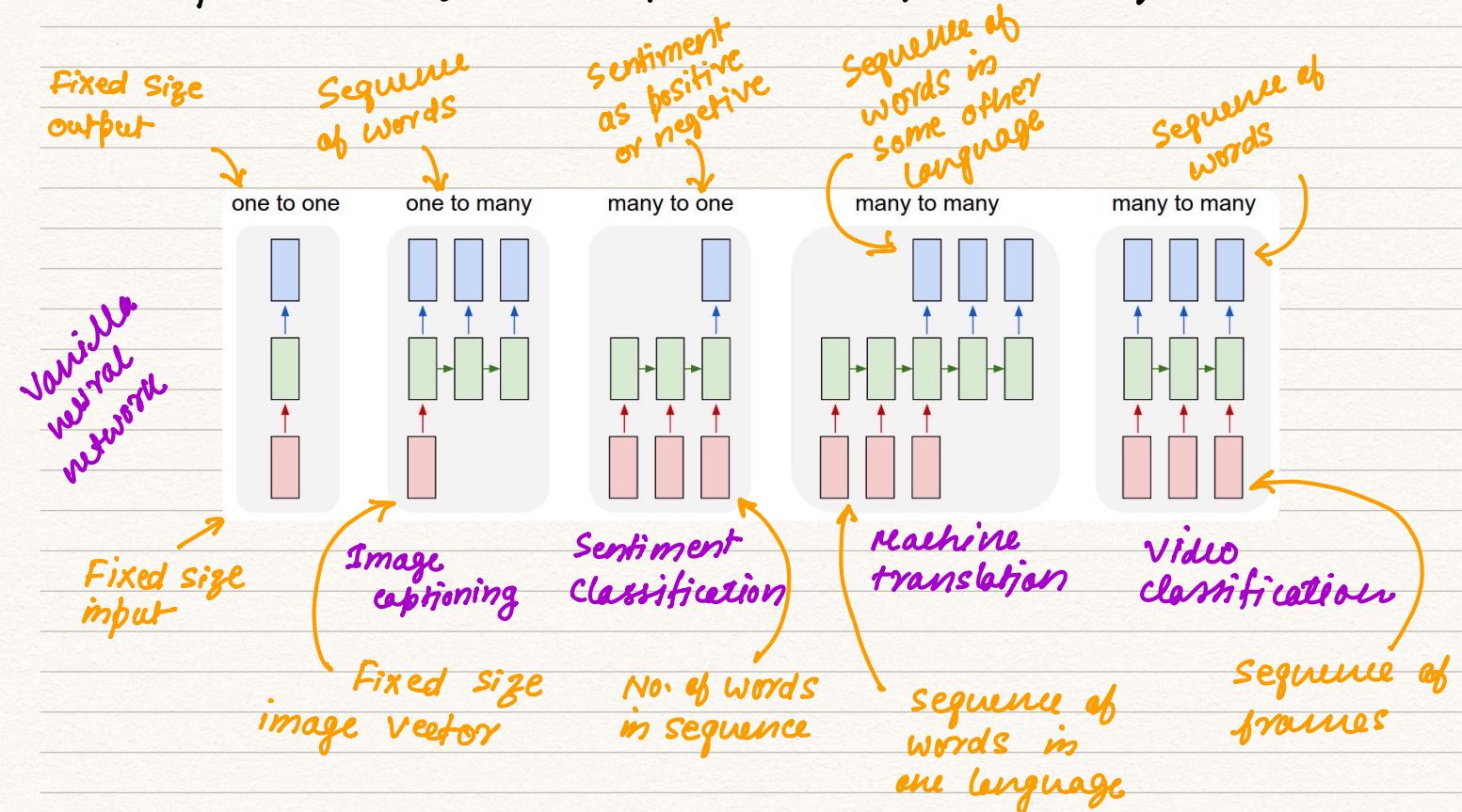


RNN → provides lot of flexibility in how you wire your neural net



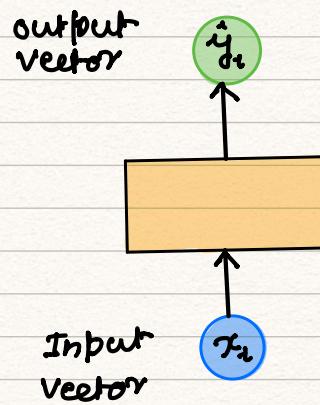
$$s_t = f_w(s_{t-1}, x_t)$$

↑ New State
↑ Recurrence function with parameters w

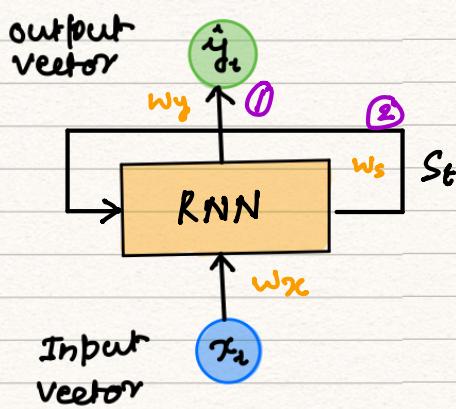
Same function and the same set of parameters are used at every time step.

Allows us to use same function to every sequence of input irrespective of size of input/output sequence

- 4 Sequence modeling problem
1. Handle variable length sequence
 2. Track long term dependencies.
 3. Maintain information about order
 4. Share parameter across the sequence



In Vanilla feed forward neural network, we are going from input to output in one direction and it can not maintain information about sequence of data.



RNN computation involves:

- ① output calculation
- ② state update

2 weights

$$S_t = \phi (w_s \cdot s_{t-1} + w_x \cdot x_t)$$

activation fn

- 4 standard neural net operations:
 → Multiplication by weight matrices
 → Apply non-linearity

$$\hat{y}_t = w_y \cdot s_t$$

output vector

$$\text{Total loss, } L = (l_1 + l_2 + l_3 + \dots + l_t)$$

Major limitation of RNN

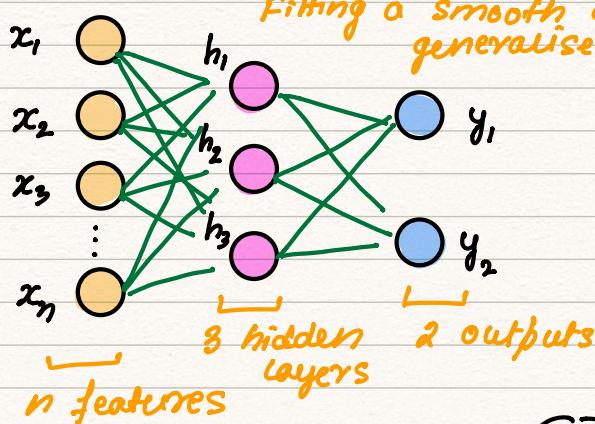
unable to capture relationships that spans more than 8 or 10 steps back

↳ called "vanishing gradient problem" ↳ LSTM solves this problem
 ↳ contribution of information decays geometrically over time

Feed forward Refresher

↳ It is non-linear function approximation

Fitting a smooth curve over the given data such that it generalises well for even the new set of data



As there are only 1 hidden layer, steps:
 ① Finding \bar{h} → from input x and weight w_1 ,
 ② finding \bar{y} → from \bar{h} and weights w_2

Feed forward.

Both these steps involve
linear combination

↳ use matrix multiplication

$w_1 \quad w_2$

w_{ij}

→ weight that
connects x_i to h_j

$[i \rightarrow 1 \text{ to } n]$
 $[j \rightarrow 1 \text{ to } 3]$

$x = [x_1, x_2, x_3, \dots, x_n]$

$$w_1 = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ \vdots & \vdots & \vdots \\ w_{n1} & w_{n2} & w_{n3} \end{bmatrix}$$

← 3 columns →

n rows

h'

$$\therefore [h'_1, h'_2, h'_3] = x \cdot w_1$$

\downarrow

1×3

$1 \times n$

$= \frac{n \times 3}{\text{Same}}$

$$h'_1 = x_1 \cdot w_{11} + x_2 \cdot w_{12} + x_3 \cdot w_{13} + \dots + x_n \cdot w_{n1}$$

$$h = \phi(h')$$

activation function

↳ Allow the network to represent non-linear relationships between its inputs and its outputs

Most world data is non-linear

$$\frac{[y_1, y_2]}{1 \times 2} = \frac{[h_1, h_2, h_3]}{1 \times 3} \cdot \frac{\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}}{3 \times 2}$$

2 cols →

$$(y) = h \cdot w_2$$

Output ← Activation fn is optional here.

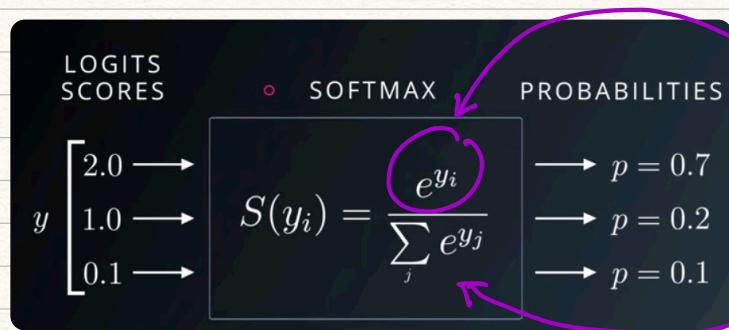
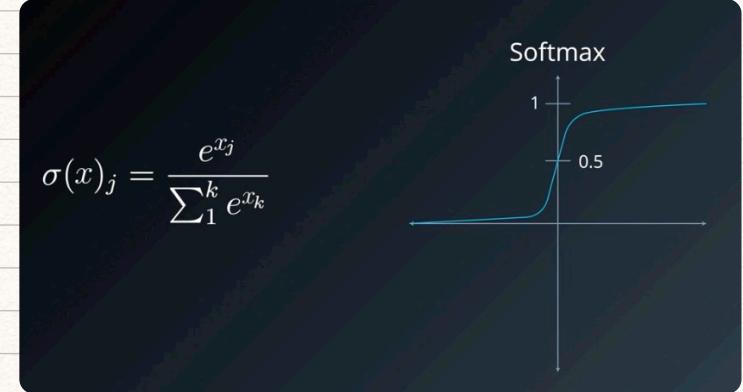
↳ you can use softmax function if you want probability score

Softmax function

A activation function that turns numbers (logits) into probabilities that sum to one



outputs a vector that represents the probability distributions of a list of potential outcomes.



Converts all outputs to 0 or positive

divide it by sum to normalise

↳ logit layers are used as last layer neurons layer for neural net for classification task

↳ produces raw prediction values as real numbers ranging from $-\infty$ to $+\infty$

Loss Function → Cross Entropy Loss is the key loss function used with Softmax function

Note that sigmoid function can be seen as special case of Softmax where logit values are $[0, \infty]$ or for binary classification.

$$\therefore S(x_i) = \frac{e^0}{(e^0 + e^x)} = \frac{1}{1+e^x}$$



cat or dog?

Python implementation

$$x = [x_1, x_2, x_3, \dots, x_n]$$

$$\text{exp} = \text{np.exp}(x)$$

$$S = \frac{\text{exp}}{\text{np.sum(exp)}}$$

Note that for large values of x , the respective exp function value will be nan

↓
use normalization

Standard

$$\frac{x - \text{np.mean}(x)}{\text{np.std}(x)}$$

min-max

$$\frac{x - \text{np.min}(x)}{\text{np.max}(x) - \text{np.min}(x)}$$

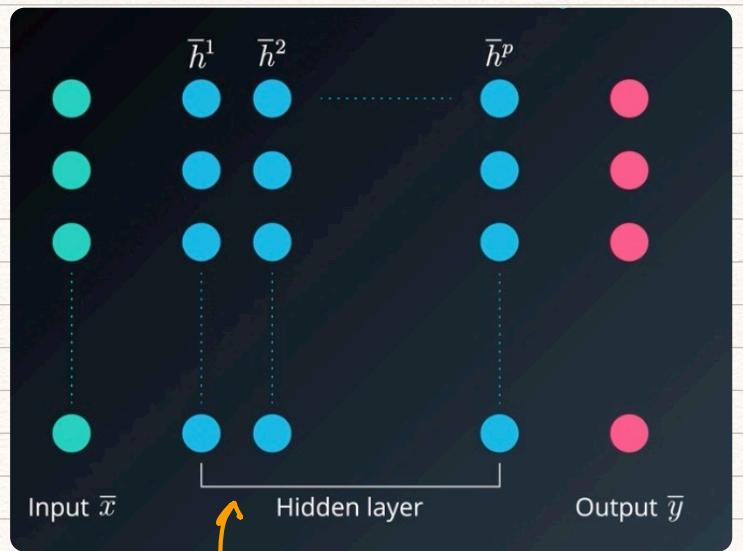
Feed Forward ... cont

For a good approximation of y there will be more hidden layers

lets say there are ϕ hidden layers

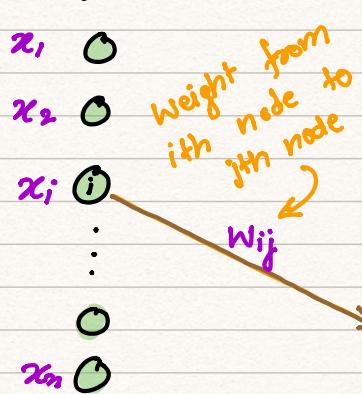
$$h^p = \phi(x \cdot w_p)$$

output of \overline{x} input (vector) activation functions
 a hidden layer (vector) weights at layer p (matrices)



No. of neurons can be different in each layers

Input \overline{x} Hidden layer
Layer k Layer $k+1$



$$h_1 = \phi \left(\sum_{i=1}^m x_i \cdot w_{i1} \right)$$

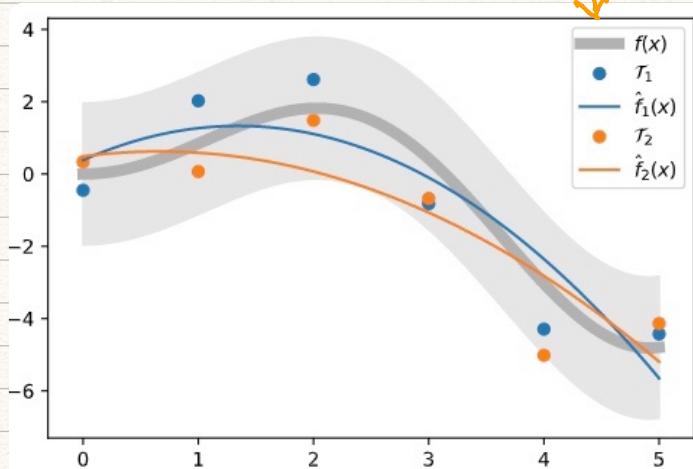
$$h_2 = \phi \left(\sum_{i=1}^m x_i \cdot w_{i2} \right)$$

$$h_m = \phi \left(\sum_{i=1}^m x_i \cdot w_{im} \right)$$

Note that bias is just another input

Error

T_1 & T_2 are 2 different Experiments



lets say, we want to find an approximation to $f(x)$

We can never observe this function directly.

We observe some training set T (T_1 & T_2)

We use this to arrive at approximation $\hat{f}_T(x)$

In general, approximation won't be a perfect fit, because of 2 source of error:

- ① Systematic error (bias) \rightarrow comes from the choice of model
- ② Random error (variance) \rightarrow comes up from the randomness inherent in the training set

Overall, we want to estimate $f(x)$ by observing \mathcal{T}
 we come up with estimator $\hat{f}_{\mathcal{T}}(x)$ and expect $[\hat{f}_{\mathcal{T}}(x) - f(x)]$ to be
 low

MSE (Mean Squared Error) \rightarrow common approach to measure average error in regression.

$$\hat{f}_{\mathcal{T}}(x) = \hat{y} \leftarrow \text{predicted value}$$

$$f(x) = y \leftarrow \text{Actual value}$$

$$\text{MSE} = \frac{\mathbb{E}_{\mathcal{T}} (\hat{y} - y)^2}{\substack{\leftarrow \text{sum of error squared} \\ \leftarrow \text{average over training set}}}$$

$$\text{Bias}(x) = \mathbb{E}_{\mathcal{T}} (\hat{y} - y)$$

$$\text{Let's consider, average prediction} = \mu(x) = \mathbb{E}_{\mathcal{T}} (\hat{y})$$

$$\text{Bias}(x) = \mu(x) - f(x)$$

\hookrightarrow here we have distributed expectation over the two terms (because expectations are linear).

$$\text{variance/Var}(x) = \mathbb{E}_{\mathcal{T}} (\hat{y} - \mu(x))^2$$

The Bias-Variance decomposition

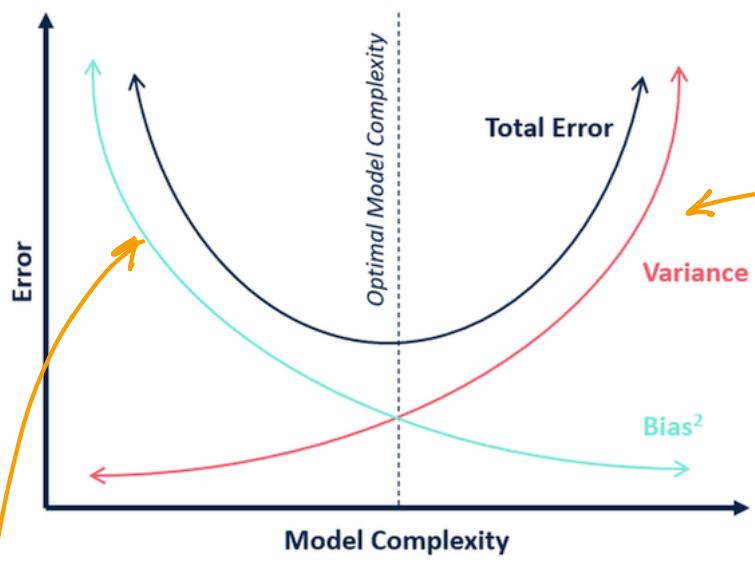
$$\begin{aligned} \text{MSE}(x) &= \mathbb{E}_{\mathcal{T}} (\hat{y} - y)^2 \\ &= \mathbb{E}_{\mathcal{T}} (\hat{y}^2 - 2\hat{y}y + y^2) \\ &= \mathbb{E}_{\mathcal{T}} (\hat{y}^2) - 2\mathbb{E}_{\mathcal{T}} \hat{y} \cdot y + \mathbb{E}_{\mathcal{T}} y^2 \\ &\quad \leftarrow \mu(x) \\ &= \mathbb{E}_{\mathcal{T}} (\hat{y})^2 - 2\mu(x)f(x) + f(x)^2 \\ &= \mathbb{E}_{\mathcal{T}} (\hat{y})^2 - \mu(x)^2 + \mu(x)^2 - 2\mu(x)f(x) + f(x)^2 \\ &= \underbrace{\mathbb{E}_{\mathcal{T}} (\hat{y})^2 - \mu(x)^2}_{\text{variance}} + \underbrace{(\mu(x) - f(x))^2}_{\text{Bias}} \\ &= \text{Var}(x) + \text{Bias}(x)^2 \end{aligned}$$

As model complexity increases, more of MSE can be attributed to variance

\hookrightarrow using regularizations, we artificially increase bias to reduce variance.

So, Expected prediction error or EPE(\hat{y}) = $\text{Var}(\hat{y}) + \text{Bias}(\hat{y})^2 + \sigma^2$

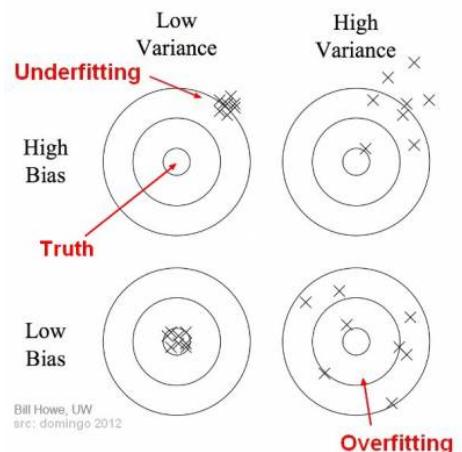
irreducible
error (Noise)



Model with high bias pays very little attention to the training data and oversimplifies the model.

↳ Leads to high error on both training and test data

Model with high variance pays a lot of attention to training data and does not generalize on new data



Backpropagation

Goal → find weights that minimize error

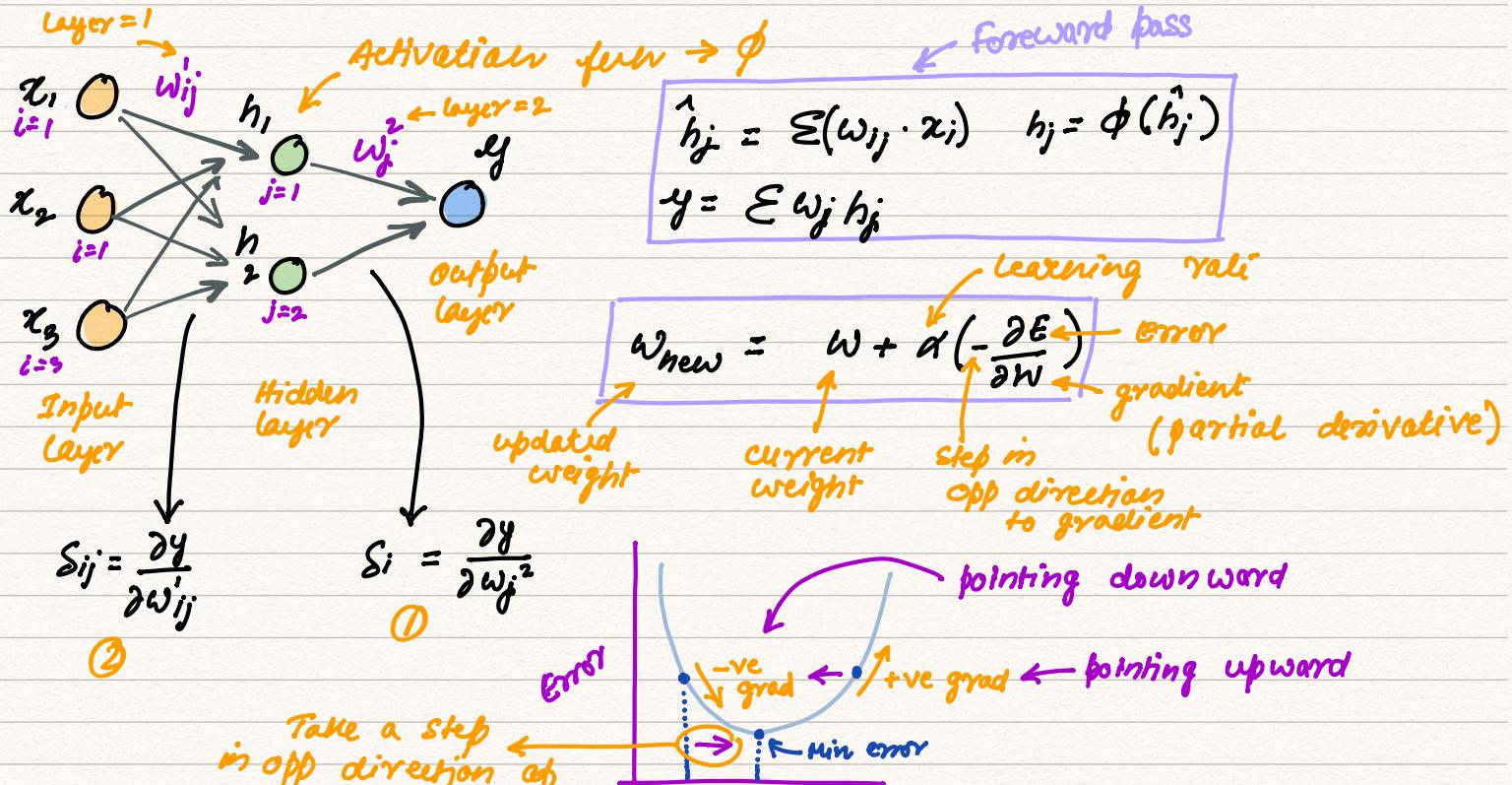
Backpropagation algorithm looks for the min of the error function in weight space using the method of gradient descent

For gradient calculation, we must guarantee the continuity and differentiability of the error function.

$$w_{ij} = w_{ij} + \Delta w_{ij}$$

gradient

$$\begin{aligned} \Delta w_{ij} &= -\alpha \frac{\partial E}{\partial w_{ij}} = -\alpha \frac{\partial \frac{(\hat{y} - y)^2}{2}}{\partial w_{ij}} \Rightarrow \boxed{\Delta w_{ij} = \alpha (\hat{y} - y) \frac{\partial \hat{y}}{\partial w_{ij}}} \\ E &= \frac{(\hat{y} - y)^2}{2} \\ &= -\alpha \frac{\partial \frac{(\hat{y} - y)^2}{2}}{\partial w_{ij}} = -\alpha \frac{\partial \frac{1}{2} (\hat{y} - y)^2}{\partial w_{ij}} \end{aligned}$$



① $\frac{\partial y}{\partial \omega_i} = \frac{\partial (\omega_i h_i)}{\partial \omega_i}$ Linear summation of terms

$$= h_i$$

$$\frac{\partial y}{\partial h} = \frac{\partial (\omega_j h_j)}{\partial h} = \omega_j^2$$

② $\frac{\partial y}{\partial \omega_{ij}}$

Chain rule

$$\frac{\partial y}{\partial \omega} = \frac{\partial y}{\partial h} \times \frac{\partial h}{\partial \omega}$$

$$= \omega_j^2 \times \phi'_j \times x_i$$

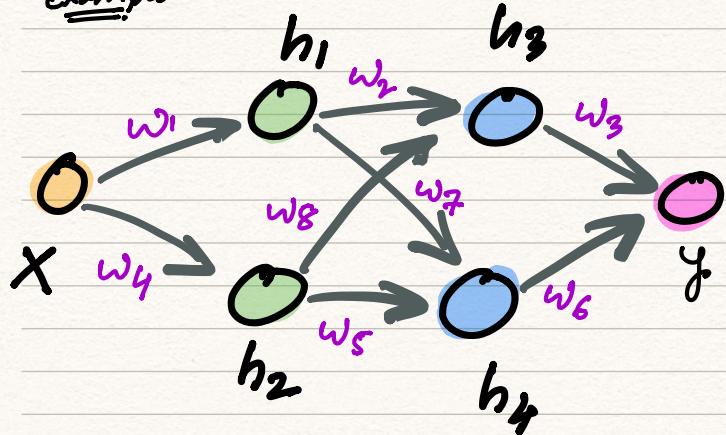
Partial derivative of activation func

$$\frac{\partial h}{\partial \omega} = \frac{\partial (\phi_j(x_i \omega_{ij}))}{\partial \omega_{ij}}$$

$$= \partial \phi_j / \partial \omega_{ij} \times \frac{\partial (x_i \omega_{ij})}{\partial \omega_{ij}}$$

$$= \phi'_j \times x_i$$

Example



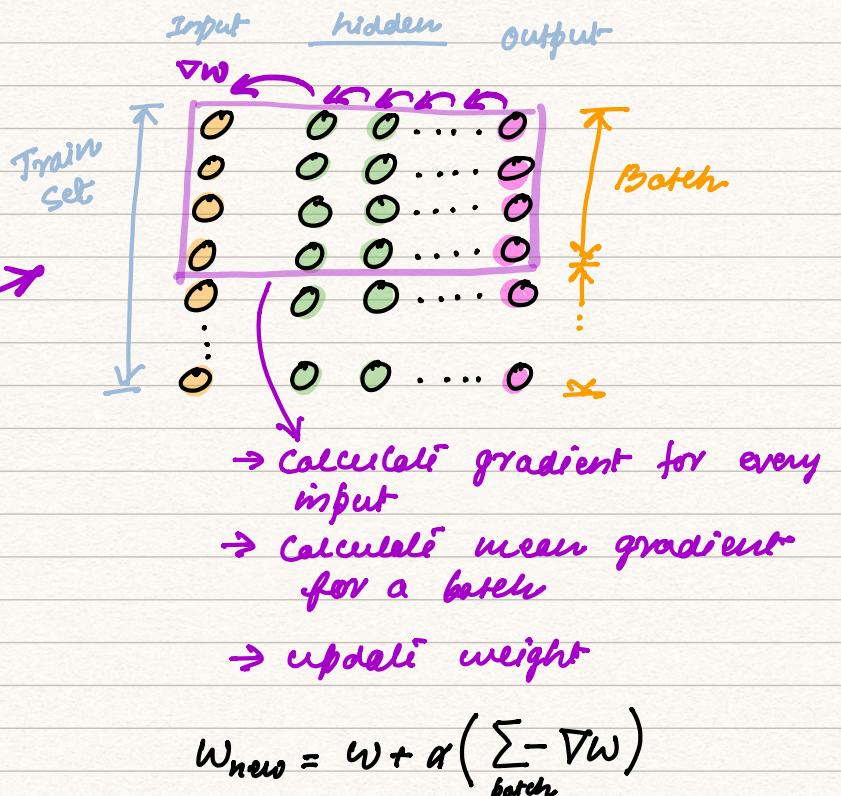
What is update rule for weight matrix $\underline{\underline{\omega}}$

$$\frac{\partial y}{\partial \omega_1} = \left[\frac{\partial y}{\partial h_3} \times \frac{\partial h_3}{\partial h_1} \times \frac{\partial h_1}{\partial \omega_1} \right. \\ \left. + \frac{\partial y}{\partial h_4} \times \frac{\partial h_4}{\partial h_1} \times \frac{\partial h_1}{\partial \omega_1} \right]$$

Gradient Descent

Stochastic
gradient (SG)

mini
Batch
Gradient Descent



Input	hidden	output
0	0 0 ... 0	
0	0 0 ... 0	
0	0 0 ... 0	
0	0 0 ... 0	
0	0 0 ... 0	
⋮	⋮	⋮
0	0 0 ... 0	

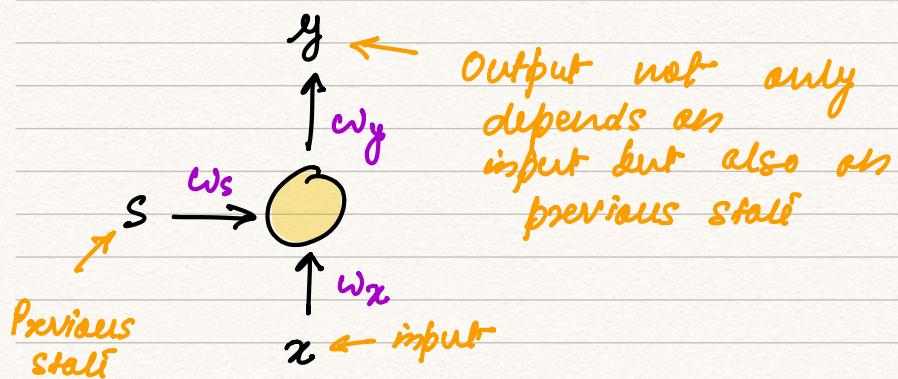
- calculate gradient for every input
- update weight

$$w_{\text{new}} = w + \alpha (-\nabla w)$$

Back to ... RNN (Recurrent Neural Network).

↳ Takes care of temporal dependences

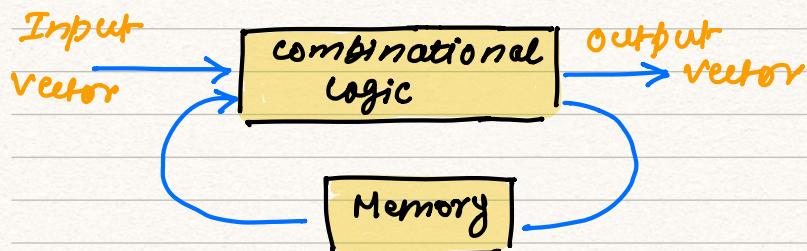
↳ dependences over time



usage

- ① Sentiment analysis
- ② Speech recognition
- ③ Time series prediction
- ④ NLP
- ⑤ Gesture recognition

State Machine

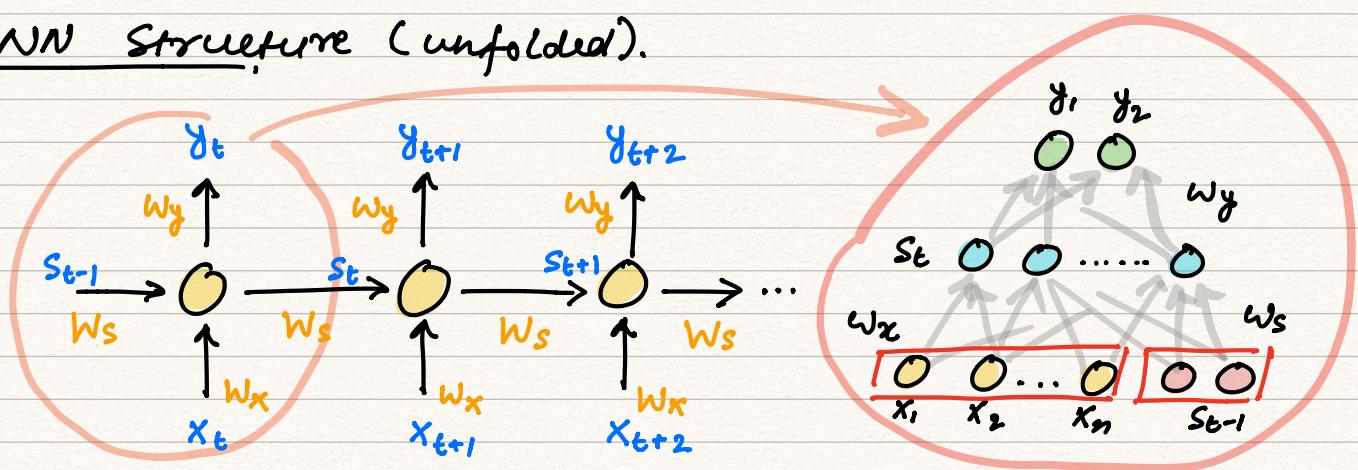


2 key feature of a RNN →

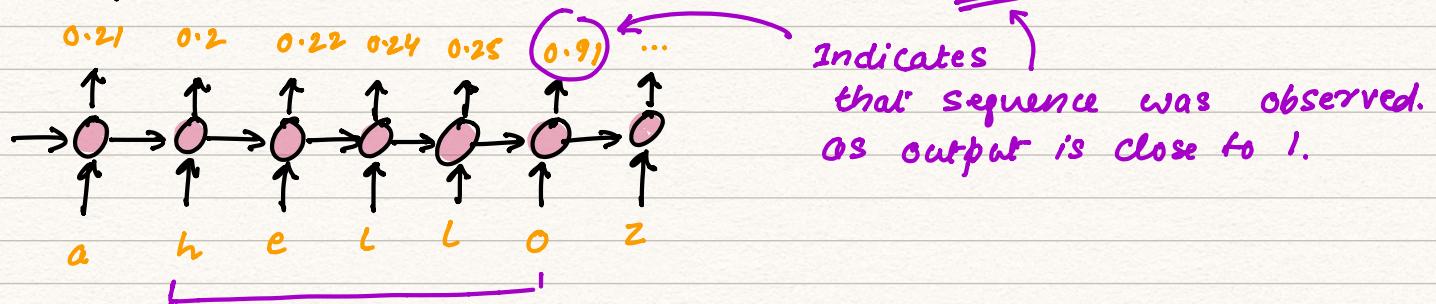
- ① Sequences as input in training phase
- ② memory elements — output of hidden layer neurons, which will serve as additional input to the network during next training step

Elman Network → the basic 3 layer neural network with feedback that serve as memory inputs

RNN Structure (unfolded).



Example \rightarrow we train model on word 'hello'

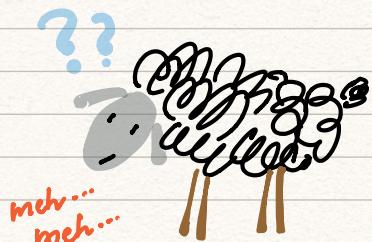


Note that we can train RNN on different lengths of words.
RNN can deal with varying sequence lengths.

Backpropagation through Time (BPTT)

$$S_t = \tanh(x_t \cdot w_x + S_{t-1} \cdot w_s)$$

$\overbrace{x_t}$ activation fn
 $\overbrace{w_x}$ Input
 $\overbrace{S_{t-1}}$ State at t
 $\overbrace{w_s}$ State at (t-1)



$$y_t = S_t \cdot w_y ; \text{ In case you are using softmax}$$

$$y_t = S(S_t \cdot w_y)$$

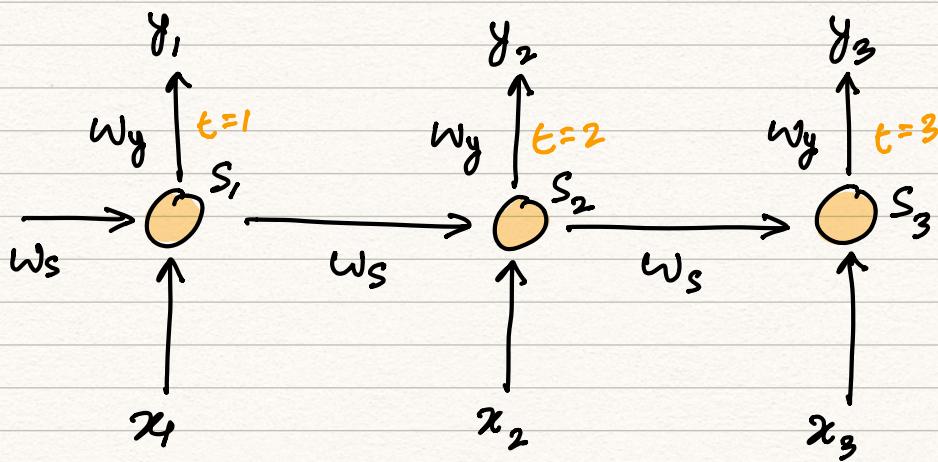
$$E_t = (\hat{y}_t - y_t)^2$$

$\overbrace{\hat{y}_t}$ prediction
 $\overbrace{y_t}$ softmax
 $\overbrace{E_t}$ Error

$$E_1 = (\hat{y}_1 - y_1)^2$$

$$E_2 = (\hat{y}_2 - y_2)^2$$

$$E_3 = (\hat{y}_3 - y_3)^2$$

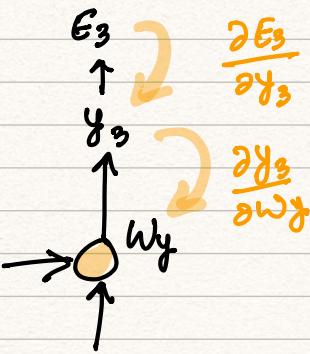


Let's adjust weight at $E = 3$.

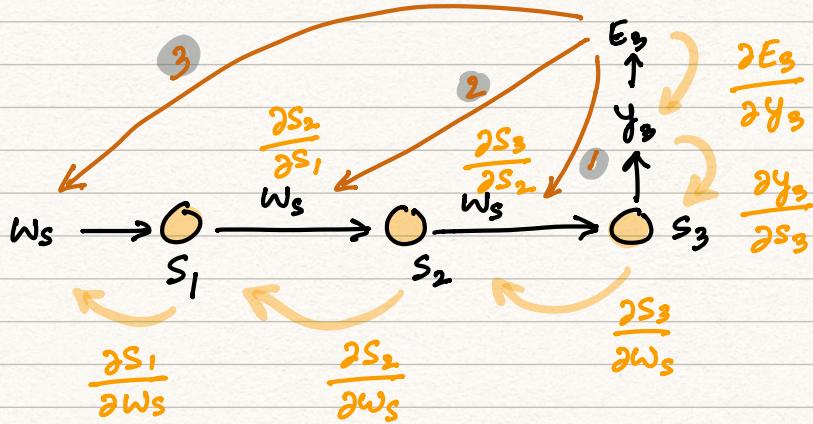
① Adjust weight matrix w_y

$$E_3 = (y_{ij} - y_j)^2$$

$$\frac{\partial E_3}{\partial w_y} = \frac{\partial E_3}{\partial y_3} \times \frac{\partial y_3}{\partial w_y}$$



② Adjust weight matrix w_s



$$\frac{\partial E_3}{\partial w_s} = 1 + 2 + 3$$

$$1) \frac{\partial E_3}{\partial w_s} = \frac{\partial E_3}{\partial y_3} \times \frac{\partial y_3}{\partial s_3} \times \frac{\partial s_3}{\partial w_s}$$

$$2) \frac{\partial E_3}{\partial w_s} = \frac{\partial E_3}{\partial y_3} \times \frac{\partial y_3}{\partial s_3} \times \frac{\partial s_3}{\partial s_2} \times \frac{\partial s_2}{\partial w_s}$$

$$3) \frac{\partial E_3}{\partial w_s} = \frac{\partial E_3}{\partial y_3} \times \frac{\partial y_3}{\partial s_3} \times \frac{\partial s_3}{\partial s_2} \times \frac{\partial s_2}{\partial w_s}$$

$$\therefore \frac{\partial E_3}{\partial w_s} = \frac{\partial E_3}{\partial y_3} \times \frac{\partial y_3}{\partial s_3} \times \frac{\partial s_3}{\partial w_s}$$

+

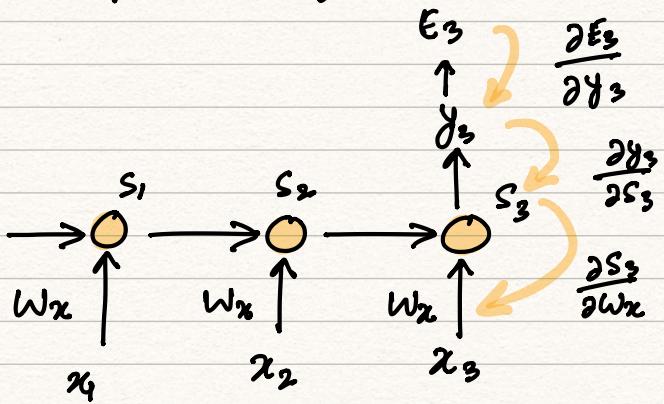
$$\frac{\partial E_3}{\partial y_3} \times \frac{\partial y_3}{\partial s_3} \times \frac{\partial s_3}{\partial s_2} \times \frac{\partial s_2}{\partial w_s}$$

$$\frac{\partial E_3}{\partial y_3} \times \frac{\partial y_3}{\partial s_3} \times \frac{\partial s_3}{\partial s_2} \times \frac{\partial s_2}{\partial s_1} \times \frac{\partial s_1}{\partial w_s}$$

General Equation

$$\frac{\partial E_n}{\partial w_s} = \sum_{i=1}^n \frac{\partial E_n}{\partial y_m} \times \frac{\partial y_m}{\partial s_i} \times \frac{\partial s_i}{\partial w_s}$$

③ Adjust weight matrix w_x



Gradient = accumulative gradient of each of the previous steps.

Stali S_3

Stali S_2

$$\frac{\partial E_3}{\partial w_x} = \frac{\partial E_3}{\partial y_3} \times \frac{\partial y_3}{\partial s_3} \times \frac{\partial s_3}{\partial w_x} + \frac{\partial E_3}{\partial y_3} \times \frac{\partial y_3}{\partial s_3} \times \frac{\partial s_3}{\partial s_2} \times \frac{\partial s_2}{\partial w_x}$$

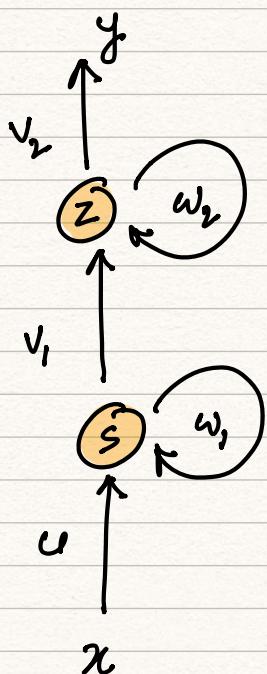
$$+ \frac{\partial E_3}{\partial y_3} \times \frac{\partial y_3}{\partial s_3} \times \frac{\partial s_3}{\partial s_2} \times \frac{\partial s_2}{\partial s_1} \times \frac{\partial s_1}{\partial w_x}$$

Stali S_1

General Equation

$$\frac{\partial E_m}{\partial w_x} = \sum_{i=1}^m \frac{\partial E_m}{\partial y_m} \times \frac{\partial y_m}{\partial s_i} \times \frac{\partial s_i}{\partial w_x}$$

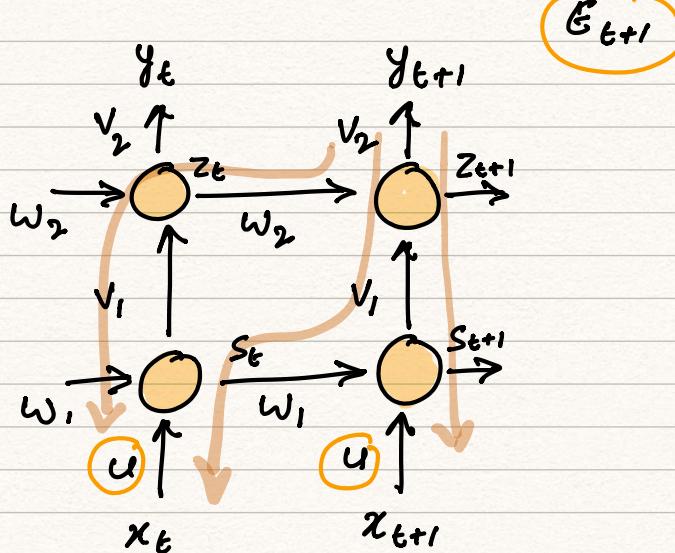
Example



Assume that error is E ! What is update rule of weight matrix w at time $t+1$?

Solution

the diagram is unfolded —

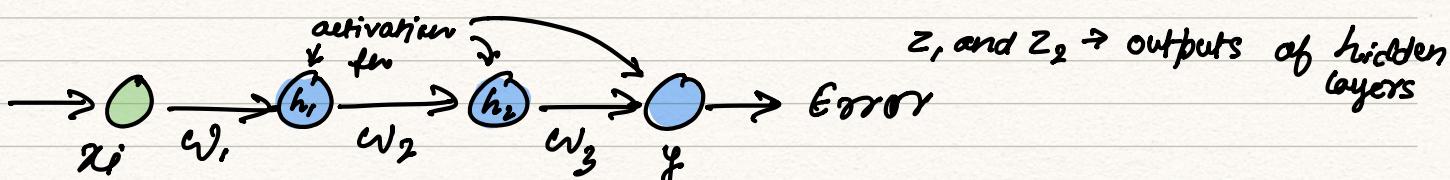


$$\begin{aligned}\frac{\partial E_{t+1}}{\partial u} = & \frac{\partial E_{t+1}}{\partial y_{t+1}} \times \frac{\partial y_{t+1}}{\partial z_{t+1}} \times \frac{\partial z_{t+1}}{\partial s_{t+1}} \times \frac{\partial s_{t+1}}{\partial u} \\ & + \frac{\partial E_{t+1}}{\partial y_{t+1}} \times \frac{\partial y_{t+1}}{\partial z_{t+1}} \times \frac{\partial z_{t+1}}{\partial s_{t+1}} \times \frac{\partial s_{t+1}}{\partial s_t} \times \frac{\partial s_t}{\partial u} \\ & + \frac{\partial E_{t+1}}{\partial y_{t+1}} \times \frac{\partial y_{t+1}}{\partial z_{t+1}} \times \frac{\partial z_{t+1}}{\partial s_t} \times \frac{\partial s_t}{\partial u}\end{aligned}$$

Vanishing gradient problem → If we backpropagate more than 8-10 time steps, then gradient will become too small (contribution of information decays exponentially over time).

↳ In summary, temporal dependency spanning more than 8-10 timestep will be disregarded.

↓
LSTM addresses this issue.



Assume, we are using sigmoid activation fn.

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial y} \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$\hat{z}_2 = h_2 \times w_3 \quad | \quad z_2 = \sigma(\hat{z}_2)$

$\frac{\partial(\sigma(\hat{z}_2))}{\partial \hat{z}_2} \times w_3 = [\sigma' \cdot w_3]$

$[\sigma' \times w_2]$

weights → chosen using

Gaussian dist with mean=0 & Std=1 ∴ all terms here are ≤ 1 hence prod will be very small and will with increasing layers

Exploding Gradients → the gradient grows uncontrollably.