

Perception

* Neuron মাদ্রেকে control করতু সাবে স্মিলেট Neuron GI - Parameters.

* অন্যি Neuron কে সা provide করলো এই Control mechanism কে help করতু ক্ষেত্ৰফলকে রাখা হ'ব Hyperparameters

Example → Weight - Neuron কেখানকা weights generate কৰতু সাবে তাৰ
weights GI - parameters. এখন বিৰু কোই weights এৰ value
কমান্তে 100 টোৱা কৰিব ক্ষেত্ৰে cross কৰতু সাবে না, GI 100 টোৱা কৰিব
ক্ষেত্ৰে cross কৰতু সাবে না. GI neuron কে ক্ষেত্ৰে Hyper-
parameters কৰ - কৰতো,

∴ Hyperparameter neuron কে control কৰতু - আৰু - neuron মাদ্রেকে control
কৰতু তাৰা GI - Parameters.

* Transfer function → কোটি function কে - গোকোটি state GI transfer কৰা,
কোটি function GI - কৰিবলৈ কোই ক্ষেত্ৰকা কলা হ'ব
transfer function.

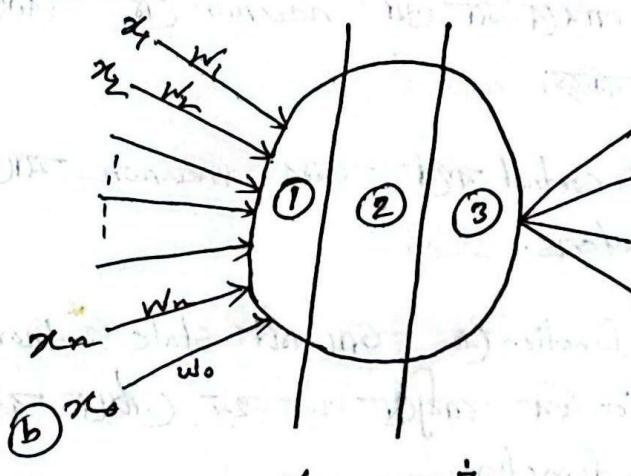
যিৰি, input and weights - কোটি external - ফিল্ট্ৰ, GI ক্ষেত্ৰ
process কৰতু neuron GI - ক্ষেত্ৰৰ Property information
pass কৰাৰ - procedure কৰালৈ - GI - transfer function.

* Compartment 1 GI weighted sum compute কৰে - ফিল্ট্ৰেকা transfer
কৰতু Compartment 2 কে নিয়ে - কোটি - কৰতু,

Weights and inputs - কুলো multiply কৰা - কৰতু, তখন এৰ impact
factor GI কৰা হ'ব, positive কোই positive কুলো multiply কৰালৈ
+ve কৰতু, positive কোই negative কুলো multiply কৰালৈ value কমান্তে
কৰতো - কৰতু, -ve and +ve multiply কৰালৈ +ve কৰতু যাৰ, -ve
লাইকুলার +ve and -ve কৰালৈ - কৰতু - কোটি - weighted sum কৰে
value reduction কৰতো যাবো, Negative input কোই negative weight
কুলো - কোই importance to - আৰু কামিকো কৰিবলৈ - কৰতু, GI কৰালৈ
কৰতো control কৰতু transfer function. - কোটি GI - অধিকৃত কৰিব

transfer করা হবে Activation function ।

* Activation function একটি security guard এর মতো, -গোস রাস্তা calculation এর condition check না করে Activation function information শুল্ক compartment- 3 টি প্রতে দিয়ে না, 3 no compartment 6 মাদ্যমিতে না দেয় -এর মধ্যে 3 no compartment হেতু এর outcome এর মধ্যে ছোট zero. -একটি Neuron activate হবে না, Neuron activate হল �outcome একটি Non-zero value -এর মধ্যে,



3 major components

(1) Transfer function / weighted sum function

$$\text{Weighted sum } V_j = \sum_{i=1}^n x_i w_i \quad i \rightarrow \text{num of inputs}$$

Value of jth neuron

(2) Activation function

$$\phi(V_j) = \begin{cases} 1; & \sum_{i=1}^n x_i w_i \geq \theta_j \\ 0; & \sum_{i=1}^n x_i w_i < \theta_j \end{cases}$$

If transfer function includes bias,

$$V_j = \sum_{i=1}^n x_i w_i + b_j$$

প্রার্থনা Perception
এর জন্য bias is inverse of threshold

If
Threshold = 0,
Bias এর বিপরীত
Threshold + 2k
consider এর
জন্য

$$b_j = x_0 w_0$$

bias * প্রার্থনা

$$\therefore V_j = \sum_{i=0}^n x_i w_i$$

(3) Output function

$$y_j = \phi(V_j)$$

Weighted sum

প্রার্থনা Threshold
এর জন্য bias

Activation function

- মাত্র, ক্ষেত্র control

জন্য neuron এর দ্রুতি

Data set \Rightarrow Excel sheet এ represent করা - এখন এই Data এর
columns শুল্ক মানে (প্রথম - অনেক শুল্ক category এর values
শুল্ক - বাধা-হ্রয়, শুল্ক) এবং neuron এর inputs x_i
এবং - গুরুত্ব importance - এবং যাই একই include করি,
- যেখানে একই একই বাদ দিয়ে দেওয়া

Excel sheet ৰেখ দিয়ে last column ৰেখাৰ output ৰেখাৰ বলে
ডেটা-ফার্ম মৈল Data set ৰেখাৰ বলে ৰেখাৰ supervised data set.

* Neuron কুলোৱা Row অনুসৰে Learn কৰিবত = সকল - ফার্ম,

* Normally column কুলোৱা পঠিতা কৰা হয়, আবার rows কুলোৱা
ৰেখাৰ value কুলোৱাৰneuron কৰিবত = ফার্ম

* One epoch → কোনোৰ মধ্যে training example হৈলে, একই training
examples ৰেখাৰ starting থক্ষে আৰু বাবু ending পৰ্যন্ত full ^{কোনোৰ} revision
ডেটাগুৰু বলে হৈলে ৰেখাৰ one epoch.

* এক epochথেকে next epoch ৰেখাৰ shuffling কৰি বিতৰি হয়,
বাবুণ কোনোৰ revision ডেটাগুৰু style কৰা- আৰু অন্ধেৰ বাবুণ
revision ডেটাগুৰু style same হৈলো, এই process কুলোৱা neuron
ৰেখাৰ হৈলো, Training ৰেখাৰ কোনোৰ বিশেষজ্ঞ test ৰেখাৰ
evaluate কৰিবত হৈলো, Training ৰেখাৰ কোনোৰ বিশেষজ্ঞ test ৰেখাৰ
একই epoch ৰেখাৰ process ৰেখাৰ চলে

* Training proper way তে সাক্ষী বিশেষজ্ঞ কোনোৰ কোনোৰ
validation কৰি আৰু

* Good dataset সাক্ষী ডেটা হৈলে (যোগোক ও নাই জোক)
কৰিবত হৈলে → Training, validation and Testing.

Training বিশেষজ্ঞ validation ensure কৰি training ৰেখাৰ properly proceed
কৰিবত বিশেষজ্ঞ

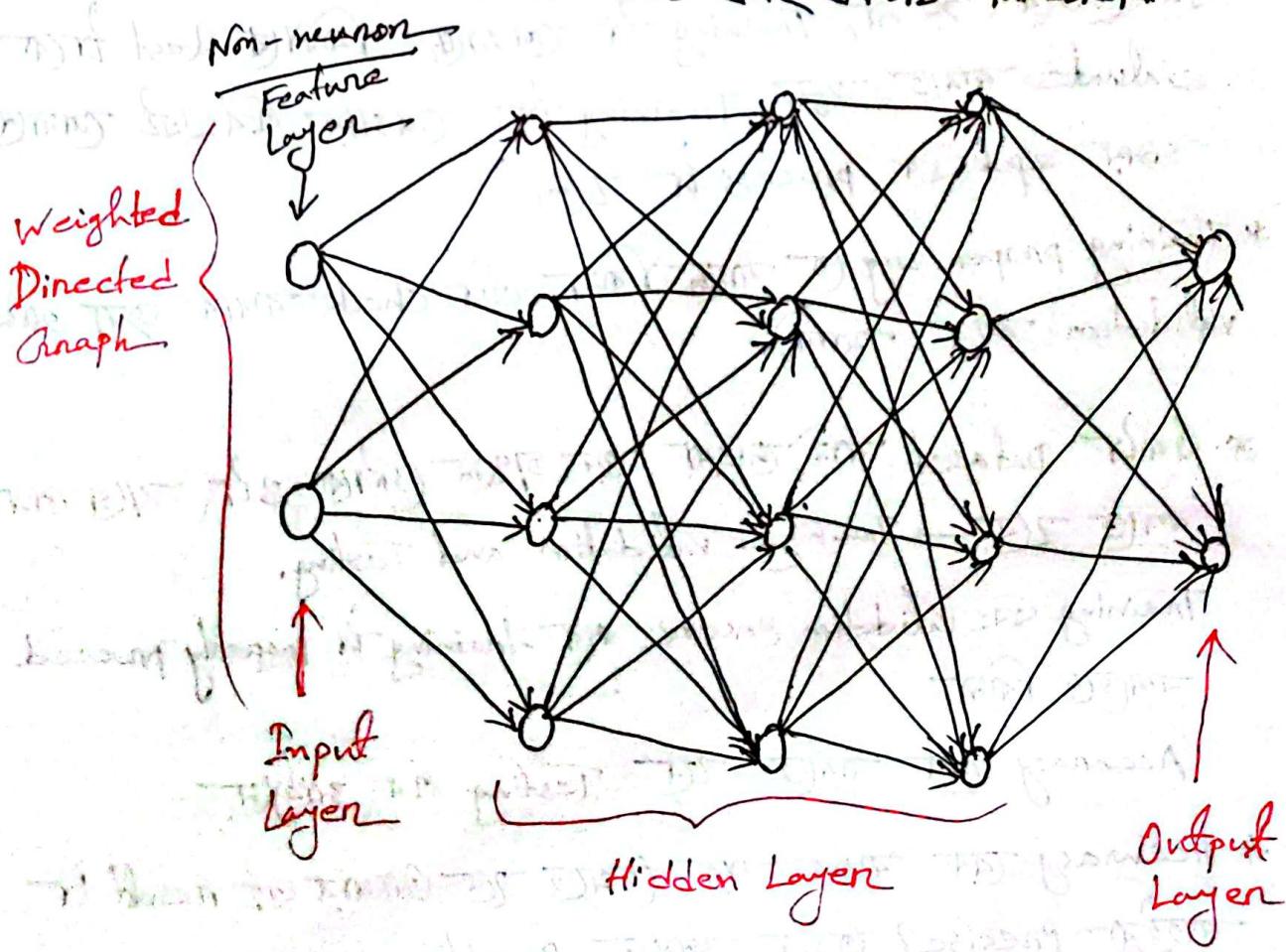
Accuracy কৰি কৰিবত হৈলে Testing ৰেখাৰ জায়গে

* Accuracy কৰি কৰিবত পৰ্যন্ত কৰি কোনোৰ কোনোৰ result ৰেখাৰ
কাতুলোৱা precise? কোনোৰ কোনোৰ precision.

- * বাত্তুলাতে pit falls করাইলি যাই ক্ষেত্রে গোড়া recall
- * প্রয়োগ-আধুনিক statistical measurement score (Accuracy, precision এবং Recall এবং Precision বনানো করা), ক্ষেত্রে কলাম F1 - Measurement.
- * প্রয়োগ-গোড়া Explainable AI. - আমি প্রয়োগ করে বলতেছি - আমরা এই model সেখানকে সহজে বলতে চাই যে কিন্তু আমি বলতে, model কৈমানে কিভাবে explain করবে, এবং particular accuracy ও রীতিটা কেন হচ্ছে?
- * প্রয়োগ-গোড়া Perception perform করবে

Parameters → Bias, Weights, Threshold

একটি Neuron G একটি Bias এবং একটি Threshold



How many parameters are there?

① This is a fully connected neural network

Each neuron of the next layer is connected with every other neuron of its previous layer: No connection in own layer's neurons

If fully connected

Input layer 2 (not considering as neurons)

H_1 4 neurons

H_2 4 neurons

H_3 4 neurons

Output layer 2 neurons

1 bias and 1 Threshold
at each neuron

$$\text{So, } 4+4+4+2 \rightarrow 14 \text{ biases}$$
$$14 \text{ Thresholds}$$

$$\text{So far parameters} = 14+14=28$$

What about weights?

$$(4 \times 2) + (4 \times 4) + (4 \times 4) + (2 \times 4)$$

H_1 input H_2 H_1 H_3 H_2 output H_3

$$= 48$$

$$\therefore \text{Total parameters} = 48 + 28 = 76$$

→ 76 for variables ~~array~~
→ store ~~array~~ 2x2
Matrix operations is
much faster than
normal variable
operations

* यही बात 2x2 Neural network four layer,

— 6x2 (or) 5 or layer तो 2x2 Input layer

— 2x2 तो 4 or layer (3 neuron 2x2 तो,

5 layers of neural networks, including inputs.

Row \Rightarrow dimension 1 \rightarrow Row vector

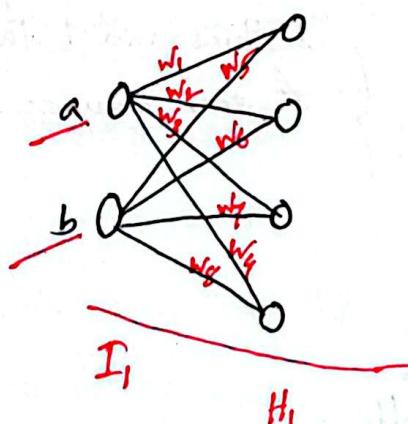
Column \Rightarrow dimension 1 \rightarrow Column vector

Matrix representation \Rightarrow main target computation faster than

3rd Matrix \Rightarrow neuron control \Rightarrow weight matrix,

Threshold matrix and
Bias matrix

Matrix Representation (Fully connected NN)



$$O_1 = \begin{matrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix}$$

Z_i^L \leftarrow i^{th} neuron
 Z_i^L \leftarrow weighted sum
 L^{th} layer

H_1 \leftarrow layer
 $H_{11} H_{12} H_{13} H_{14}$ \leftarrow Num of neuron

$$\begin{matrix} a & \begin{bmatrix} w_1 & w_2 & w_3 & w_4 \\ w_5 & w_6 & w_7 & w_8 \end{bmatrix} \\ b & \end{matrix}$$

Transpose

$$\begin{bmatrix} w_1 & w_5 \\ w_2 & w_6 \\ w_3 & w_7 \\ w_4 & w_8 \end{bmatrix} \quad (4 \times 2)$$

Weighted sum
depends on a & b

$$\begin{matrix} \text{Layer} \\ \text{Weighted sum} \\ \text{Neuron no} \end{matrix} \quad \begin{bmatrix} Z_1^1 \\ Z_2^1 \\ Z_3^1 \\ Z_4^1 \end{bmatrix} = \begin{bmatrix} w_1 a + w_5 b \\ w_2 a + w_6 b \\ w_3 a + w_7 b \\ w_4 a + w_8 b \end{bmatrix}$$

$$\begin{bmatrix} Z_1^1 \\ Z_2^1 \\ Z_3^1 \\ Z_4^1 \end{bmatrix} = \begin{bmatrix} w_1 & w_5 \\ w_2 & w_6 \\ w_3 & w_7 \\ w_4 & w_8 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} + \begin{bmatrix} \text{bias} \\ \vdots \\ \vdots \end{bmatrix}$$

$$\Rightarrow Z_i^1 = W_L^T \begin{bmatrix} a \\ b \end{bmatrix}$$

for i^{th} neuron
for i^{th} bias
 (4×1)

$$F(z_i^1) = f\left(\begin{bmatrix} z_1^1 \\ z_2^1 \\ z_3^1 \\ z_4^1 \end{bmatrix}\right) = \begin{bmatrix} f(z_1^1) \\ f(z_2^1) \\ f(z_3^1) \\ f(z_4^1) \end{bmatrix} = \begin{bmatrix} o_1 \\ o_2 \\ o_3 \\ o_4 \end{bmatrix}$$

Activation function

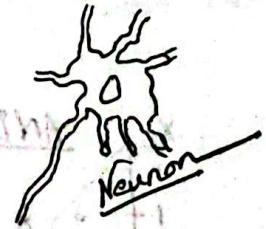
input \rightarrow Matrix size $1 \times 4 \rightarrow (2 \times 1)$

Hidden layer 1 \rightarrow output size (4×1) \rightarrow 2×1

Hidden layer 2 \rightarrow 2×1 input.

$$\begin{bmatrix} a \\ b \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \end{bmatrix}$$

(2×1) (4×2) (4×4) (3×4) (3×1)
Output (final)



pretrained Model \rightarrow Weights মূলে আগে প্রেসে ঢোকা আছে. এবং neuron
 কে শর্কি train করানো আছে, last layer এর
 output layer কি check করুন Accuracy বা
 validation কে করা হচ্ছে,

Exploring Boolean Function Through Perceptron

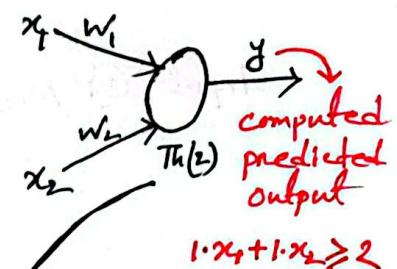
AND

Data set $\{x\}$ - Range (0-1)

Training data set $\{x\}$ & $\{y\}$
 result $y = w_1 x_1 + w_2 x_2$ (with $w_1=1, w_2=1$)
 desired result $y \geq 2$ (with $w_1=1, w_2=1$)

Threshold

x_1	x_2	$x_1 \cdot x_2$	Weighted sum $v_j = w_1 x_1 + w_2 x_2$	$v_j \geq 2$	Error
0	0	0	0	0	$d-y$
0	1	0	1	0	0
1	0	0	1	0	0
1	1	1	2	1	0

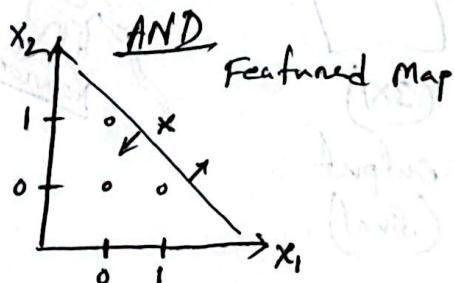


Let, $w_1=1, w_2=1, \theta=2$

If $\theta=1$

Threshold θ \neq weighted sum equal

OR gate (≥ 2) \neq OR gate (≥ 1)
 OR gate (≥ 1) nature show OR gate



Data point দুটি প্রক্রিয়া

Single line দ্বারা divide

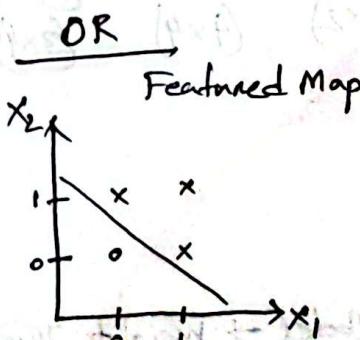
বর্ণনা করা

OR Gate Dataset (≥ 1)

linearly separable dataset.

AND - linearly separable dataset দ্বারা single perceptron

That করি করা সহজে predict করা মুশক



Linearly separable dataset.

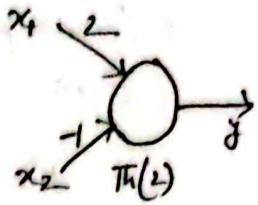
Single perceptron $\theta=1$ predict কর

সহজ

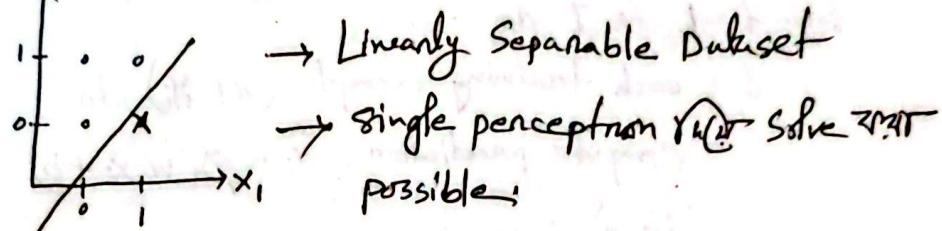
* Featured Map automatically তৈরি হলে কোথাও এল হল Deep learning.

* $x_1 \bar{x}_2$

x_1	x_2	\bar{x}_2	$x_1 \cdot \bar{x}_2$	Weighted sum $y_j = 2 \cdot x_1 + (-1) \cdot x_2$	$y_j \geq 2$	$d-y$ error
0	0	1	0	0	0	0
0	1	0	0	-1	0	0
1	0	1	1	2	1	0
1	1	0	0	1	0	0

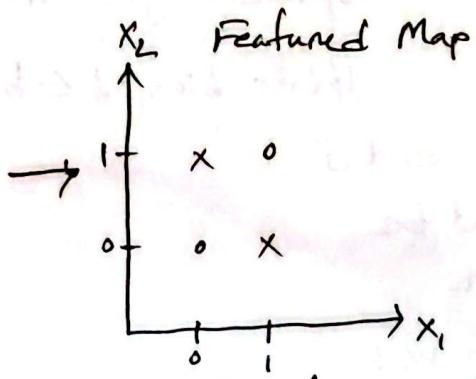


↓
Feature Map



* Exclusive OR

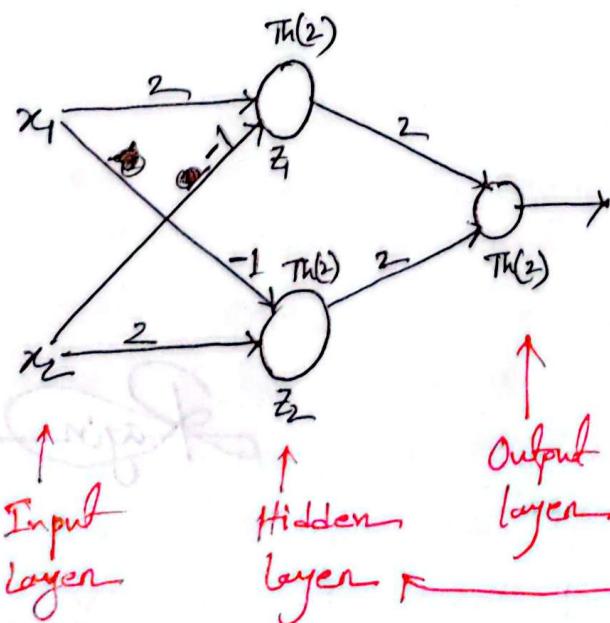
x_1	x_2	$f(x_1, x_2)$
0	0	0
0	1	1
1	0	1
1	1	0



$$x_1 \oplus x_2 = x_1 \bar{x}_2 + x_2 \bar{x}_1$$

$$\quad z_1 \quad \quad z_2$$

Not linearly separable dataset
Single perceptron $\text{Th}(1)$ can't map this



← Multilayer
Perception

Input
layer

Hidden
layer

multiple line द्वारा, (एटे)
multiple line एट थ्रॉउट
boundary region द्वारा बनाया

* Error column এর মাধ্যে Learning করা -

* Error produce - তথ্য পরিসরের Error গুরুত্বের ক্ষেত্রে Process
করা হলো - এর Perception Learning algorithm

Perception Learning Algorithm → slide

Weigh vector w and bias b

initialize $w \leftarrow 0, b \leftarrow 0$

for $t=1$ to T do max num of iteration

 for each training example (x_i, y_i) do

 Compute prediction: $\hat{y}_i = \phi(w \cdot x_i + b)$

 if $\hat{y}_i \neq y_i$ then

 Update weights: $w \leftarrow w + \eta \cdot \hat{y}_i \cdot x_i$

 Update bias: $b \leftarrow b + \eta \cdot \hat{y}_i$ error boolean

 end if

end for

end for

return w, b

hyperparameter
→ Learning rate (η)
→ Activation function

error
boolean

Rajin

Multilayer Perception

Sigmoidal Function

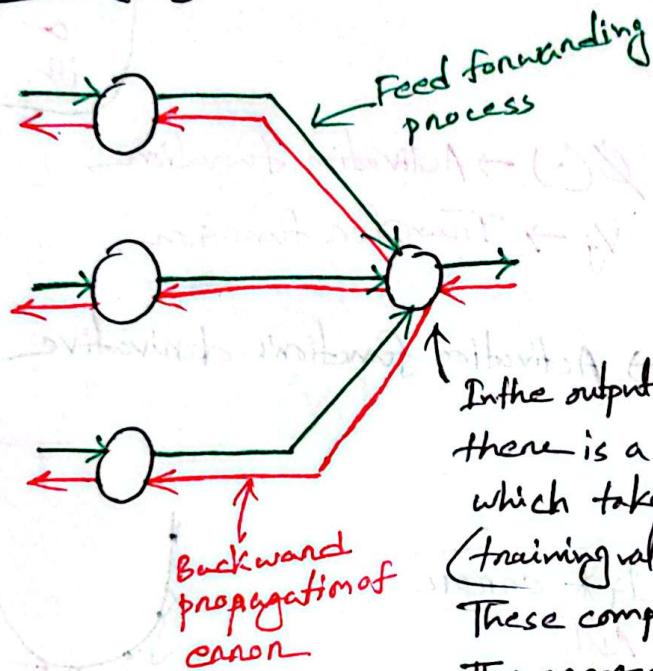
- * Differentiable
- * Support non-linearity

$$v_j = \sum_{i=0, \dots, m} w_{ji} \delta_i$$

$$\varphi(v_j) = \frac{1}{1 + e^{-av_j}}$$

$$\varphi'(v_j) = a\varphi(v_j)[1 - \varphi(v_j)]$$

Learning Algorithm



In the output layer, there is a loss function which takes two inputs - (training value, computed output value). These computes error.

The errors are propagated with red lines

Backpropagation algorithm

Let NN get weights adjusted by trying to minimize loss function (at value minimize)

Loss function \rightarrow The avg. square error

- * Error signal of output neuron j n-th training example

Total energy at time n : $e_j(n) = d_j(n) - \hat{y}_j(n)$

Avg. Square error: $E(n) = \frac{1}{2} \sum_{j \in C} e_j^2(n)$

Measure of learning performance: $E_{AV} = \frac{1}{N} \sum_{n=1}^N E(n)$

Our Goal: Adjust weights of NN to minimize E_{AV}

Minimize the avg. square error which is our loss function, so we need to adjust our NN's weights in such a way so, our loss function can be minimized.

$e_j \rightarrow$ Error at output of neuron j .

$d_j \rightarrow$ Output of neuron j .

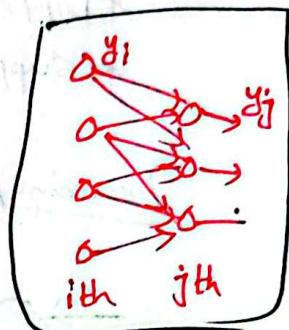
$v_j \rightarrow \sum_{i=0, \dots, m} w_{ji} y_i$ Induced local field of neuron j .

$$e_j = d_j - y_j$$

y_j \rightarrow jth neuron

$\phi(\cdot)$ \rightarrow Activation function

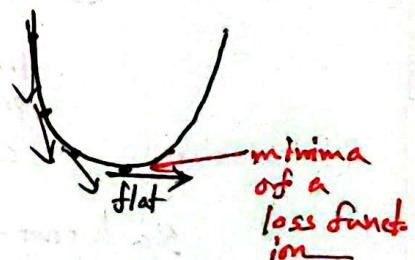
$v_j \rightarrow$ Transfer function



$$\frac{dy_j}{dv_j} = \phi'(v_j) \rightarrow \text{Activation function's derivative}$$

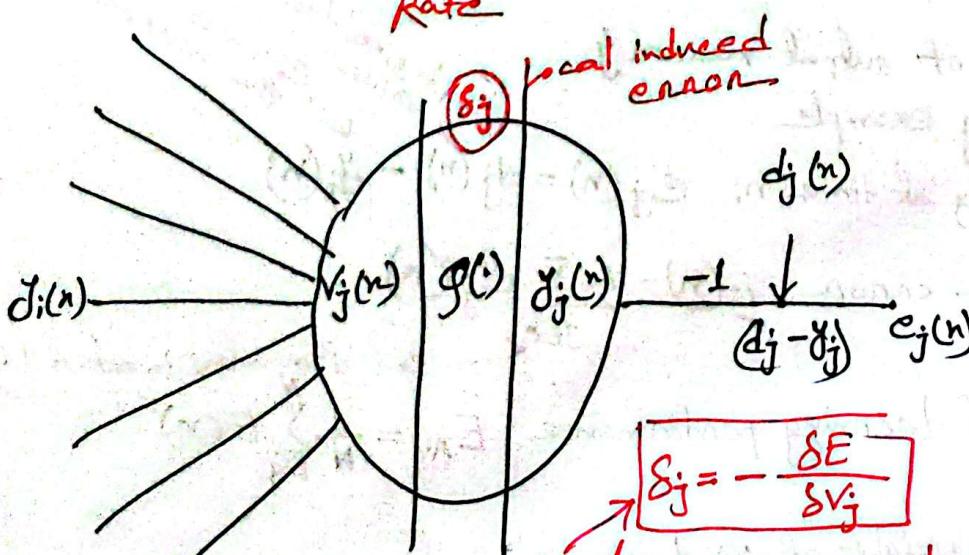
Hedbs rule

$$w_{new} = w_{old} + LR * I_j * \frac{\delta E}{\delta w_j}$$



$$\Delta w_{ji} = -\eta \frac{\delta E}{\delta w_{ji}}$$

Learning Rate



Local induced error computed inside the neuron itselfs

$$V_j = \sum w_{ji} y_i$$

$$\phi(V_j) = \underline{\quad}$$

$$y_j = \phi(V_j)$$

$$e_j = d_j - y_j$$

$$E = \frac{1}{2} \sum e_j^2$$

$$= \frac{1}{2} (e_1^2 + e_2^2 + \dots + e_j^2 + \dots + e_n^2)$$

$$\frac{\delta E}{\delta e_j} = 0 + 0 + 0 + \frac{1}{2} \cdot 2 e_j = e_j$$

$$e_j = d_j - \delta_j$$

$$= 0 - 1 = -1$$

$$\frac{\delta \delta_j}{\delta v_j} = \phi'(v_j)$$

$$\begin{aligned} \frac{\delta E}{\delta v_j} &= \frac{\delta E}{\delta e_j} \cdot \frac{\delta e_j}{\delta \delta_j} \cdot \frac{\delta \delta_j}{\delta v_j} \\ &= e_j \cdot (-1) \cdot \phi'(v_j) \end{aligned}$$

$$\therefore \delta_j = - \frac{\delta E}{\delta v_j}$$

$$\begin{aligned} &= - \{e_j \cdot (-1) \cdot \phi'(v_j)\} \\ &= e_j \phi'(v_j) \end{aligned}$$

Local induced error of the output layer

$$\delta_j = e_j \phi'(v_j)$$

$$\delta_j = - \frac{\delta E}{\delta v_j}$$

$$v_j = \sum_{i=0, \dots, n} w_{ji} y_i$$

$$\frac{\delta v_j}{\delta w_{ji}} = y_i$$

$$\Delta w_{ji} = -\eta \frac{\delta E}{\delta w_{ji}}$$

$$= -\eta \cdot \frac{\delta E}{\delta v_j} \frac{\delta v_j}{\delta w_{ji}}$$

$$= \eta \cdot \delta_j \cdot y_i$$

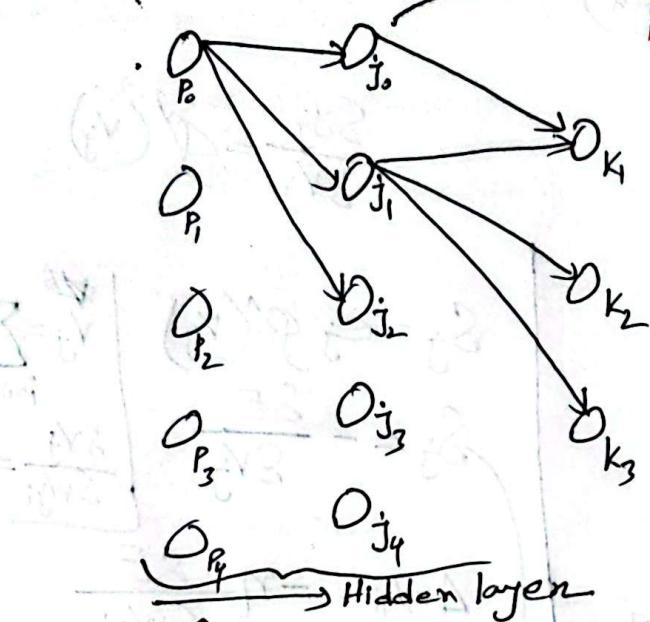
$$\therefore \Delta w_{ji} = \eta * \delta_j * y_i$$

Case: 1 $j \rightarrow$ output neuron

$$e_j = d_j - \delta_j$$

$$\delta_j = (d_j - \delta_j) \phi'(v_j)$$

Case 2: $J \rightarrow$ Hidden neuron

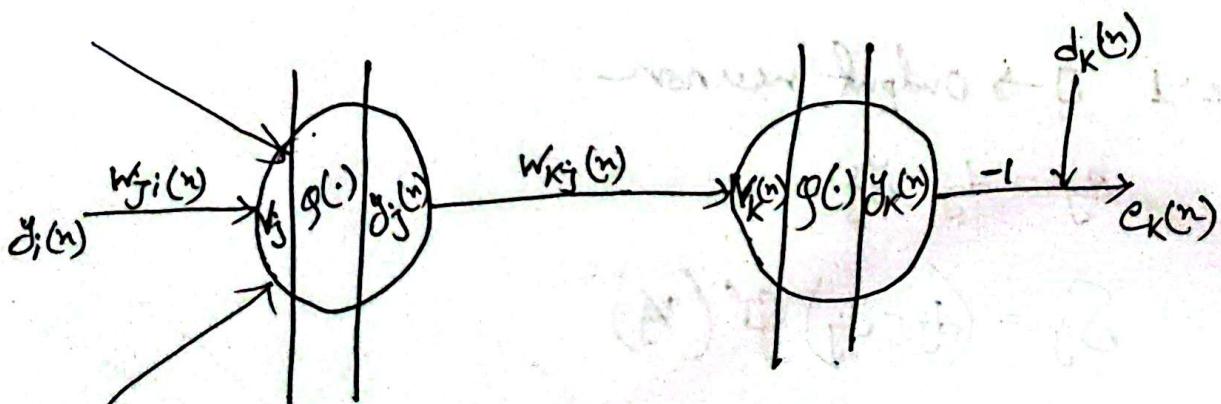


The local gradient of for neuron j is recursively determined in terms of the local gradients of all neurons to which neuron j is directly connected.

$$\delta_{j_1} = f_1(\delta_{k_1}, \delta_{k_2}, \delta_{k_3})$$

$$\delta_{P_0} = f_2(\delta_{j_0}, \delta_{j_1}, \delta_{j_2})$$

$$= f_2(f(\delta_{k_1}), f_1(\delta_{k_1}, \delta_{k_2}, \delta_{k_3}), f \dots)$$



$$\begin{aligned}\delta_j &= -\frac{\delta E}{\delta v_j} = -\frac{\delta E}{\delta \delta_j} \cdot \frac{\delta \delta_j}{\delta v_j} \\ &= -\frac{\delta E}{\delta \delta_j} \cdot \varphi'(v_j) \left[\frac{\delta \delta_j}{\delta v_j} = \varphi'(v_j) \right]\end{aligned}$$

$$\delta_k = e_k \varphi'(v_k)$$

$$\varphi(v_k) = \delta_k$$

$$e_k = d_k - \delta_k$$

$$\frac{\delta E}{\delta y_j} = ?$$

$$E(y) = \frac{1}{2} \sum_{k \in C} e_k^T(u)$$

$$\frac{\delta E}{\delta y_j} = \sum_{k \in C} e_k \cdot \frac{\delta e_k}{\delta y_j}$$

$$\begin{aligned} s_j &= -\frac{\delta E}{\delta v_j} \\ &= -\frac{\delta E}{\delta y_j} \cdot \varphi'(v_j) \end{aligned}$$

$$= \sum_{k \in C} s_k w_{kj} \cdot \varphi'(v_j)$$

$$= -\sum_{k \in C} e_k \cdot \left(\frac{\delta e_k}{\delta y_j} \right)$$

$$= -\sum_{k \in C} e_k \cdot \left(\frac{-\delta e_k}{\delta v_k} \right) \frac{\delta v_k}{\delta y_j}$$

$$= -\sum_{k \in C} e_k \underbrace{\varphi''(v_k)}_{\text{Local induced error of output layer}} w_{kj}$$

$$\begin{aligned} v_k &= \sum_{j=0, \dots, m} w_{kj} y_j \\ \frac{\delta v_k}{\delta y_j} &= w_{kj} \end{aligned}$$

$$\therefore \frac{\delta E}{\delta y_j} = -\sum_{k \in C} e_k \varphi'(v_k) w_{kj}$$

$$= -\sum_{k \in C} s_k w_{kj}$$

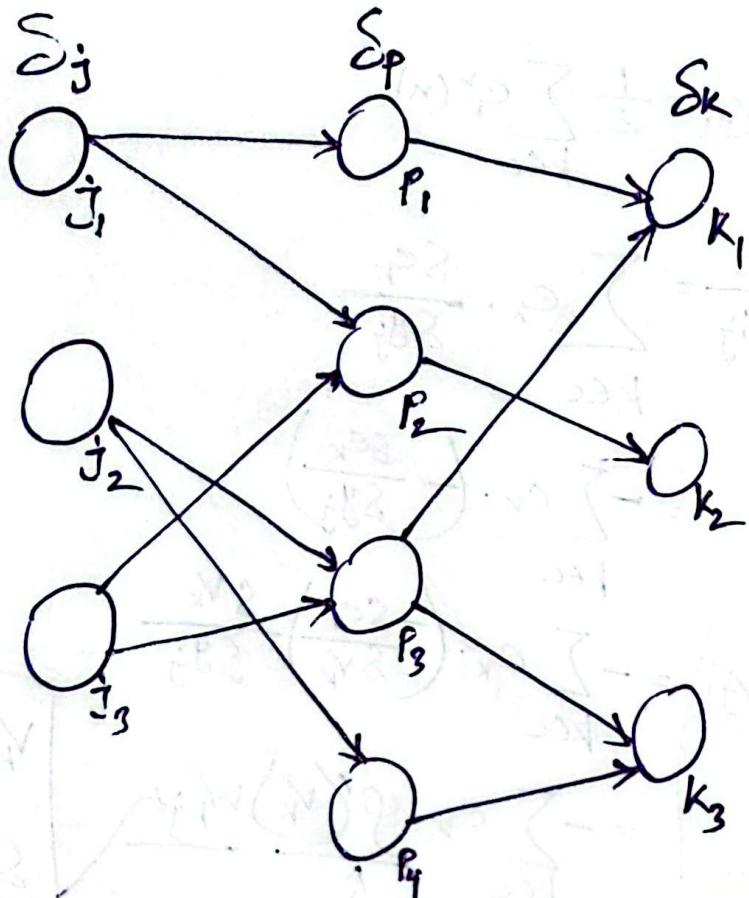
Local induced error of output layer

$$\begin{aligned} \frac{\delta e_k}{\delta v_k} &= -\frac{\delta e_k}{\delta v_k} \\ &= \frac{\delta s_k}{\delta v_k} \\ &= \varphi'(v_k) \\ \underline{v_k} &= \underline{\varphi(v_k)} \end{aligned}$$

Output layer (j) $\rightarrow s_j = e_j \varphi'(v_j)$

Hidden layer (j) $\rightarrow s_j = \sum_{k \in C} s_k w_{kj} \cdot \varphi'(v_j)$

Example



Connections

$$j_1 \rightarrow (p_1, p_2)$$

$$j_2 \rightarrow (p_2, p_4)$$

$$j_3 \rightarrow (p_2, p_3)$$

$$p_1 \rightarrow k_1$$

$$p_2 \rightarrow k_2$$

$$p_3 \rightarrow (k_1, k_2)$$

$$p_4 \rightarrow k_3$$

$$\delta k_1 = g'(k_1) e_{k_1}$$

$$\delta k_2 = g'(k_2) e_{k_2}$$

$$\delta k_3 = g'(k_3) e_{k_3}$$

$$\delta p_1 = g'(V_{p_1}) \delta k_1 w_{k_1 p_1}$$

$$\delta p_2 = g'(V_{p_2}) \delta k_2 w_{k_2 p_2}$$

$$\delta p_3 = g'(V_{p_3}) (\delta k_1 w_{k_1 p_3} + \delta k_3 w_{k_3 p_3})$$

$$\delta p_4 = g'(V_{p_4}) \delta k_3 w_{k_3 p_4}$$

$$\delta j_1 = g'(V_{j_1}) (\delta p_1 w_{p_1 j_1} + \delta p_2 w_{p_2 j_1})$$

$$= g'(V_{j_1}) \{ (g'(V_{p_1}) \delta k_1 w_{k_1 p_1}) w_{p_1 j_1} + (g'(V_{p_2}) \delta k_2 w_{k_2 p_2}) w_{p_2 j_1} \}$$

$$\delta j_2 = g(V_{j_2}) \left(\sum_{p_1, p_2 \in P_1} \delta p_i w_{p_i j_2} \right)$$

← short form
in the context of
computing

$$\delta j_3 = - - -$$

local gradient of induced error inside the cell

Computational Graph

$$f(x, y, z) = y \sin(x) + \cos(x \cdot z)$$

$$\frac{\delta f}{\delta x}, \frac{\delta f}{\delta y}, \frac{\delta f}{\delta z}$$

$$t_1 = x \cdot z$$

$$t_2 = \cos(t_1)$$

$$t_3 = \sin(x)$$

$$t_4 = y \cdot t_3$$

$$t_5 = t_4 + t_2$$

Operand \rightarrow  (Variables)

Operator \rightarrow  Operator
↓ Operand

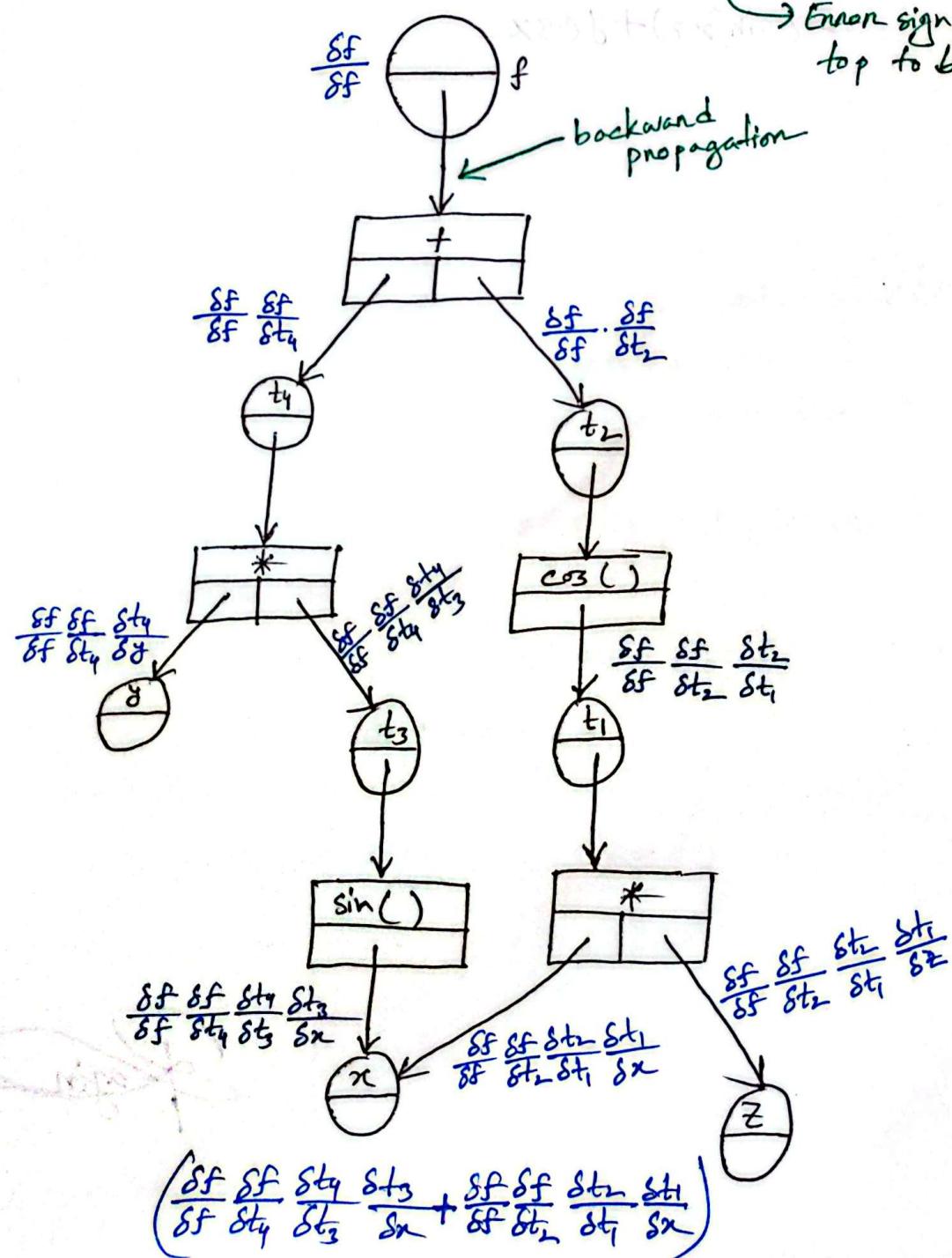
Operator \rightarrow No differentiation

Feed forwarding process

Computational result

Backward process

Error signal from top to bottom



$$\frac{\delta f}{\delta f} \frac{\delta f}{\delta t_4} \frac{\delta t_4}{\delta t_3} \frac{\delta t_3}{\delta x} + \frac{\delta f}{\delta f} \frac{\delta f}{\delta t_2} \frac{\delta t_2}{\delta t_1} \frac{\delta t_1}{\delta x}$$

$$\Rightarrow (1 \cdot 1 \cdot y \cdot \cos x) + (1 \cdot 1 \cdot (-\sin(xz))z)$$

$$\Rightarrow -z\sin(xz) + y\cos x$$

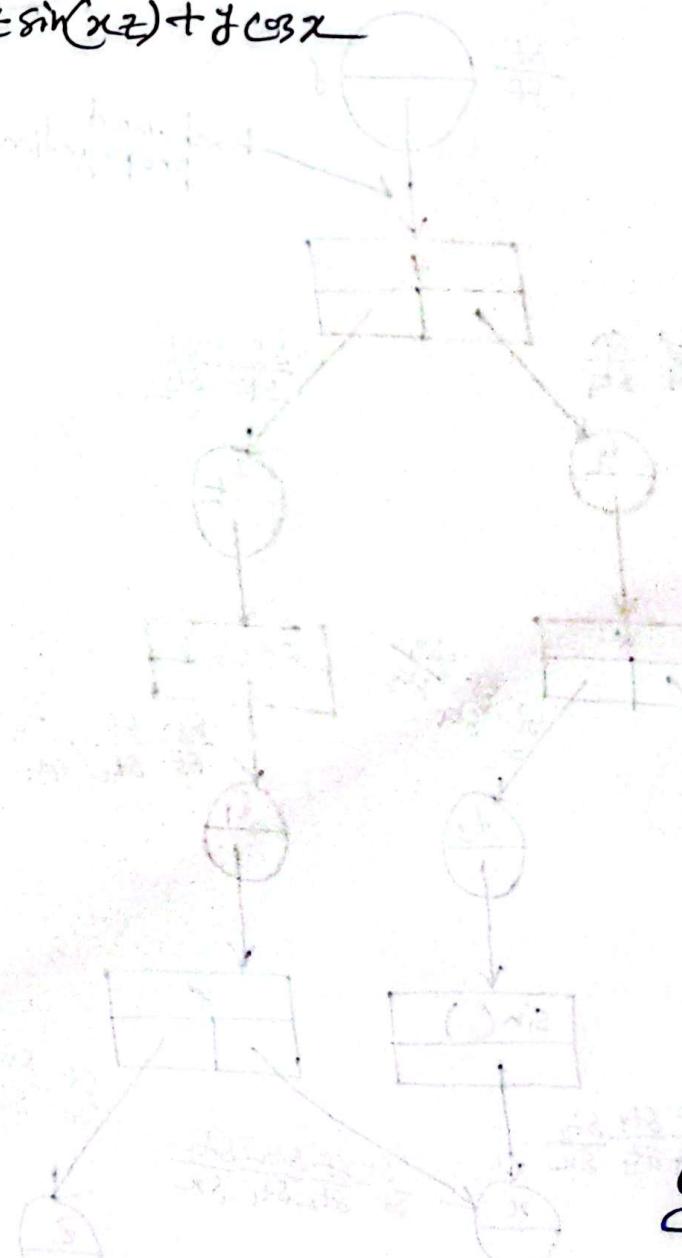
Check

$$f = y\sin(x) + \cos(xz)$$

$$\frac{df}{dx} = y \cdot \frac{\sin(x)}{\sin} + (-\sin(xz)) \cdot \frac{\partial(xz)}{\partial x}$$

$$= y\cos x + (-\sin(xz))z$$

$$= -z\sin(xz) + y\cos x$$



Rajin

Recurrent Neural Network (RNN)

① Sequence or pattern, we can apply the RNN

$\text{if } (\text{weather} == \text{sunny})$ $\quad \quad \quad \text{eat chicken}$ $\text{else if } (\text{weather} == \text{Rainy})$ $\quad \quad \quad \text{eat Burger}$	Control flow program No sequence is getting generated We can not apply RNN
--	--

input: Weather (sunny, rainy)

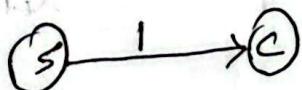
output: Food (Chicken, Burger, pizza)

Weather \rightarrow Food

MLP \rightarrow One hot vector representation — Only 1's value could be 1.

$$W = \begin{bmatrix} S \\ R \end{bmatrix} \Rightarrow S = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, R = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$F = \begin{bmatrix} C \\ B \\ P \end{bmatrix} \Rightarrow C = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, B = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, P = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$



Adj. Matrix



inputs

(P)
Food

MLP

$$\begin{array}{ccc} & C & B & P \\ S & \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} & & \\ R & \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} & & \\ & (2 \times 3) & & \end{array}$$

Transpose

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \quad (3 \times 2)$$

This problem shows how a boolean functionality can be represented with the help of MLP

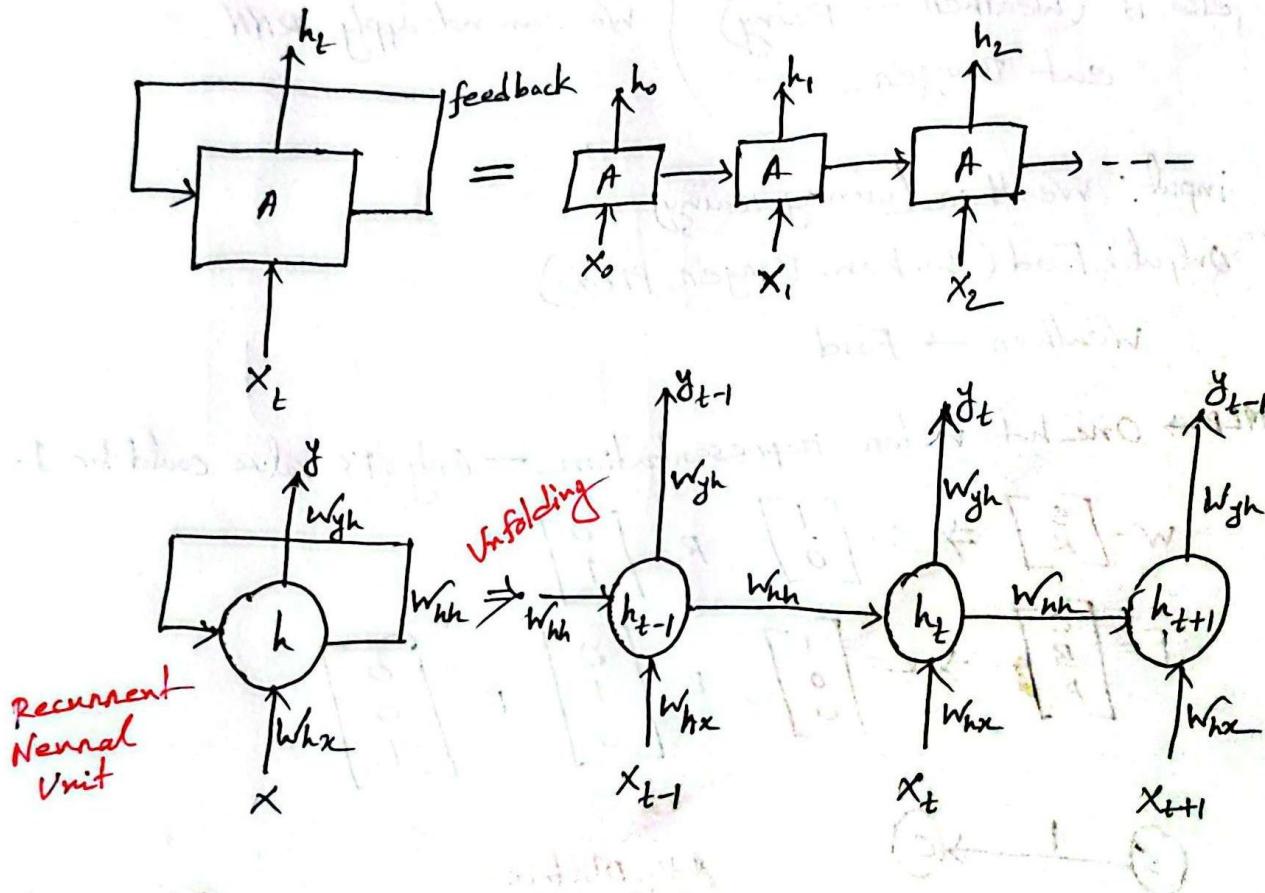
This weight matrix basically computed by the learning procedure of MLP

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}_{(3 \times 2)} \quad S = \begin{bmatrix} 1 \\ 0 \end{bmatrix}_{(3 \times 1)} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}_{(3 \times 1)}$$

This is not applicable for RNN, because it demands sequence and also needs to produce sequence or pattern.

* How it is different from MLP

① Architectural difference



② MLP can not be applied upon sequence

all the information passes in single pass so. that's why we can not apply long contextual sequences in MLP for prediction.

Example

I am travelling to Europe. Now I'm at Greece, after travelling from Germany to Rome, France and Switzerland.

In France, I need to communicate with people in --- language-
contextual information

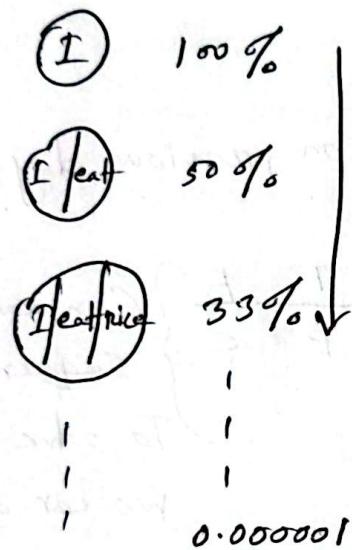
Upper hand RNN is used to capture the sequence which hold the context.

There is a limitation

It can not process or provide correct prediction if the sequence become too long.

Human memories are work like shuttle points.

Example : I eat nice everyday.



Longer sequences will wipeout
details from the context and
that's why RNN suffers from
Vanishing Gradient Problem

Exploded Gradient Problem

A diagram showing a fraction $\frac{\delta E}{\delta w_{ji}}$ enclosed in a circle. An arrow labeled "big" points to the numerator δE , and another arrow labeled "small" points to the denominator δw_{ji} . A third arrow labeled "very small" points to the overall result of the division.

Vanishing Gradient Problem

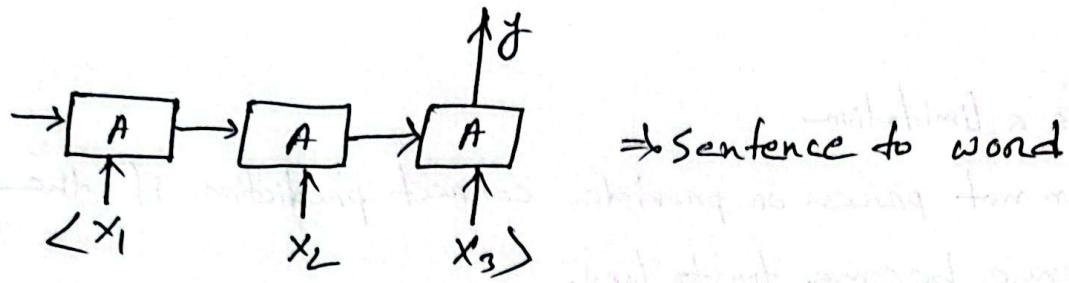
For removing vanishing gradient problem

RNN takes help of LSTM (Long short Term Memory)

Short term memory will be recalled for Long period

* 9 different formats of RNN

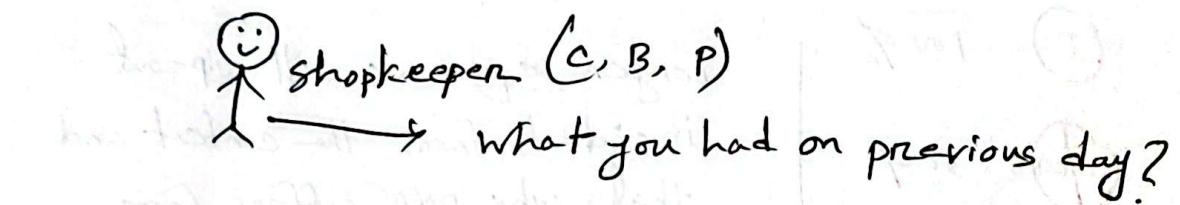
- | | | |
|--------------|---|----------------------|
| One to one | { | Word to word |
| One to many | | Word to sentence |
| Many to one | | Sentence to word |
| Many to Many | | Sentence to sentence |



\Rightarrow Sentence to word

Weather \rightarrow Food using MLP

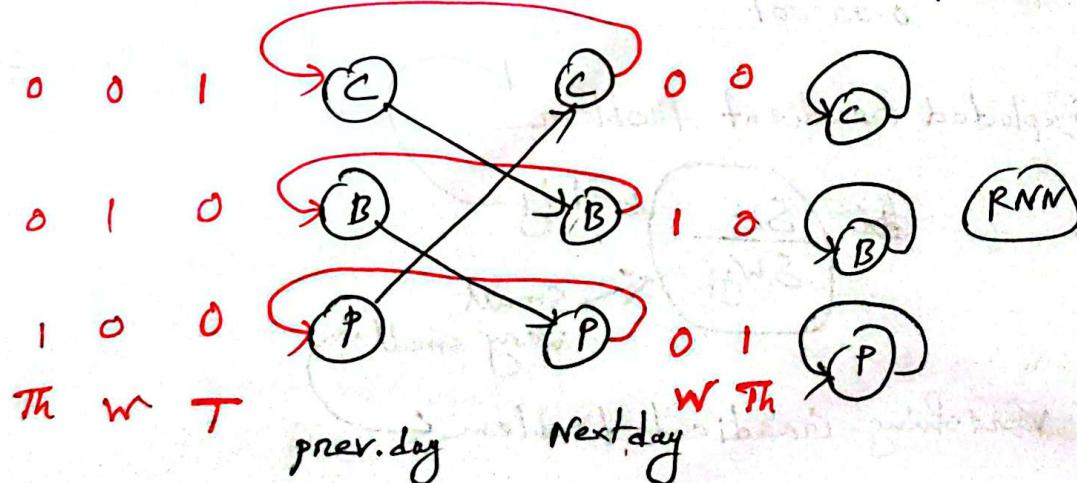
Let's think about a scenario.



S S M T W T F
C B P C B P C
B

} Generation of Sequence.

To solve this generation we can apply RNN



$$\begin{array}{c}
 \text{C} \quad \text{B} \quad \text{P} \\
 \begin{matrix}
 C & \left[\begin{matrix} 0 & 1 & 0 \end{matrix} \right] \\
 B & \left[\begin{matrix} 0 & 0 & 1 \end{matrix} \right] \\
 P & \left[\begin{matrix} 1 & 0 & 0 \end{matrix} \right]
 \end{matrix}
 \xrightarrow{\text{Transpose}}
 \begin{matrix}
 \left[\begin{matrix} 0 & 0 & 1 \end{matrix} \right] \\
 \left[\begin{matrix} 1 & 0 & 0 \end{matrix} \right] \\
 \left[\begin{matrix} 0 & 1 & 0 \end{matrix} \right]
 \end{matrix}
 \end{array}$$

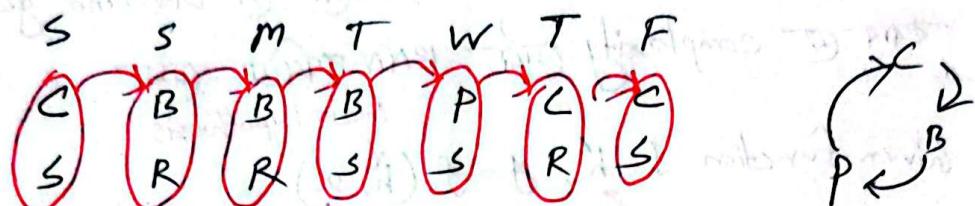
(3×3)

$$\begin{matrix}
 \left[\begin{matrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{matrix} \right] \\
 \left[\begin{matrix} 1 \\ 0 \\ 0 \end{matrix} \right]
 \end{matrix} = \begin{matrix}
 \left[\begin{matrix} B \end{matrix} \right] \\
 \left[\begin{matrix} 0 \\ 1 \end{matrix} \right]
 \end{matrix}$$

(3×3) (3×1)

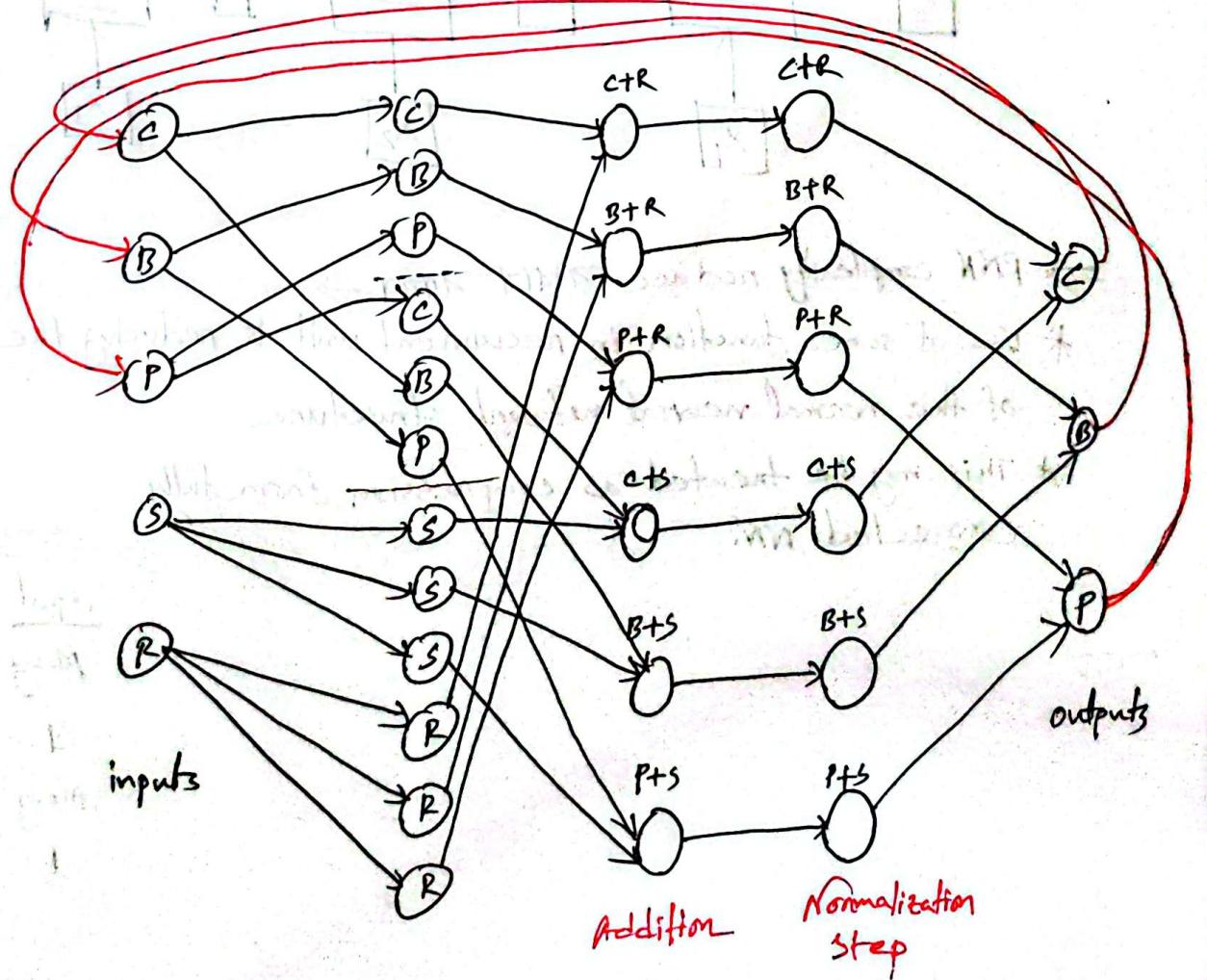
This weight matrix are computed similar fashion like MLP

Let's introduce impact of weather in the problem



- (1) If weather is sunny I can go to shopkeeper and he will provide me the next dish on the cycle
- (2) If weather is Rainy I can not go, so I had to eat previous dish.

Weather + Food \rightarrow Food



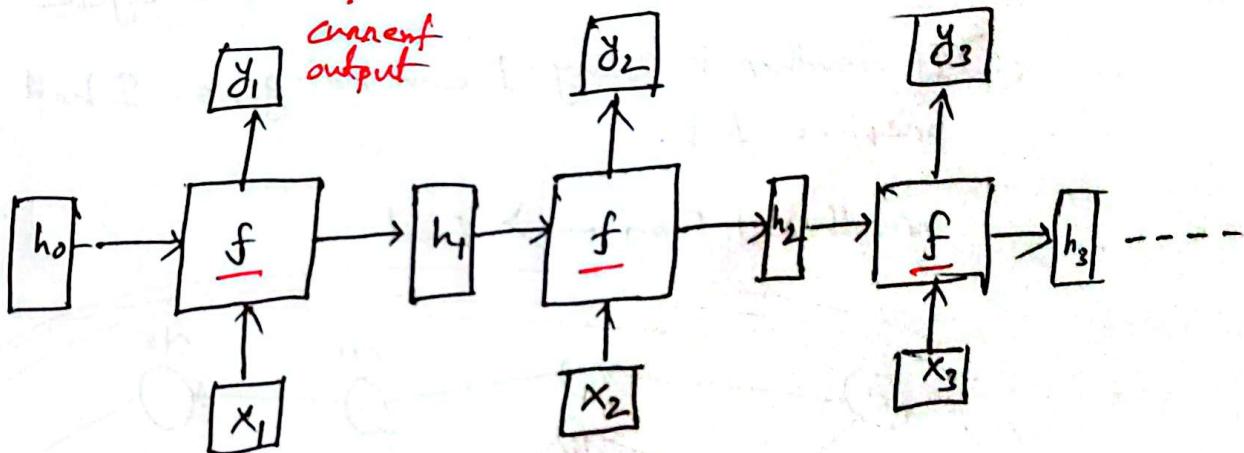
Problem
Friend +
Weather +
Food

* RNN ৰেখাৰ normal MLP দেখা হৈছে কিন্তু complex structures বলোৱা বাধা
 বাধা। Normal MLP দেখা - দেখা RNN দেখা - calculation complexity বাধা
 আৰু কিন্তু দেখা - computation দেখা ফলৰ ফলৰ outcome generate
 - বাধা - কৈ complexify কৈ RNN কৈকৈ দেখা

prev. hypothesis

Given function $f: h' \rightarrow y = f(h, x)$

current hypothesis *RNN function* *current input*



⇒ RNN complexity reduce বাধাৰ বাধা →

* Use of same function in recurrent unit is reducing the complexity
 of the normal neural network structure

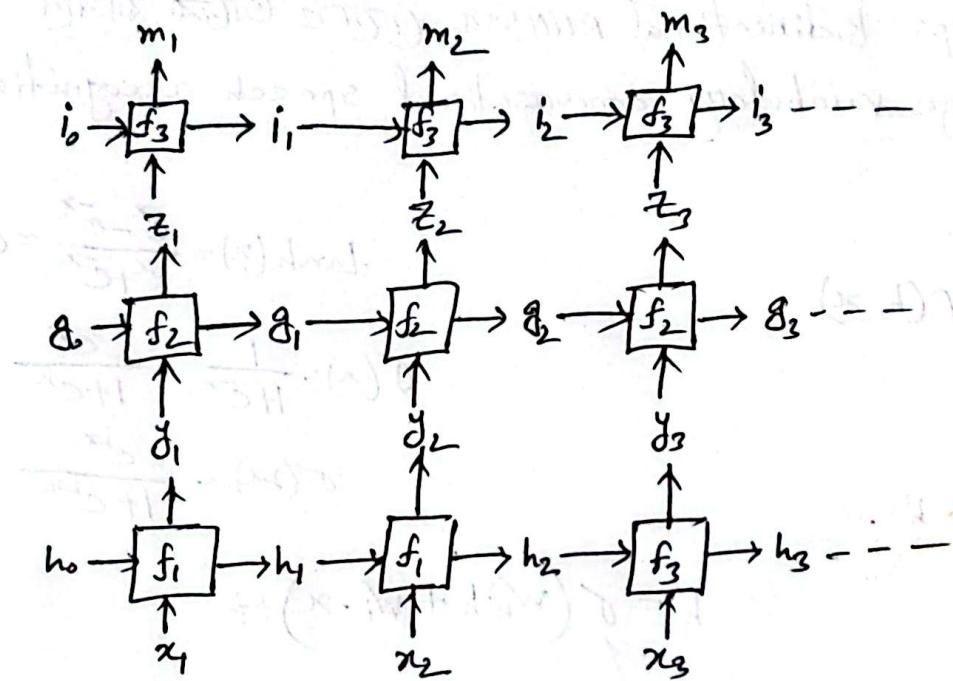
* This may be treated as compression from fully
 connected NN.

Sequences

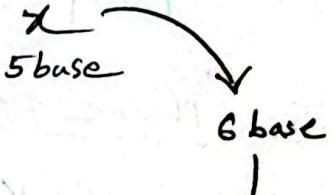
<u>input</u>	<u>Output</u>
Many to	Many
1 to	Many
Many to	1
1 to	1

Deep RNN

$$h, y = f_1(h, x) \quad g, z = f_2(g, y) \quad i', m = f_3(i, z)$$



Use case



10 base

Outcome

Language translator

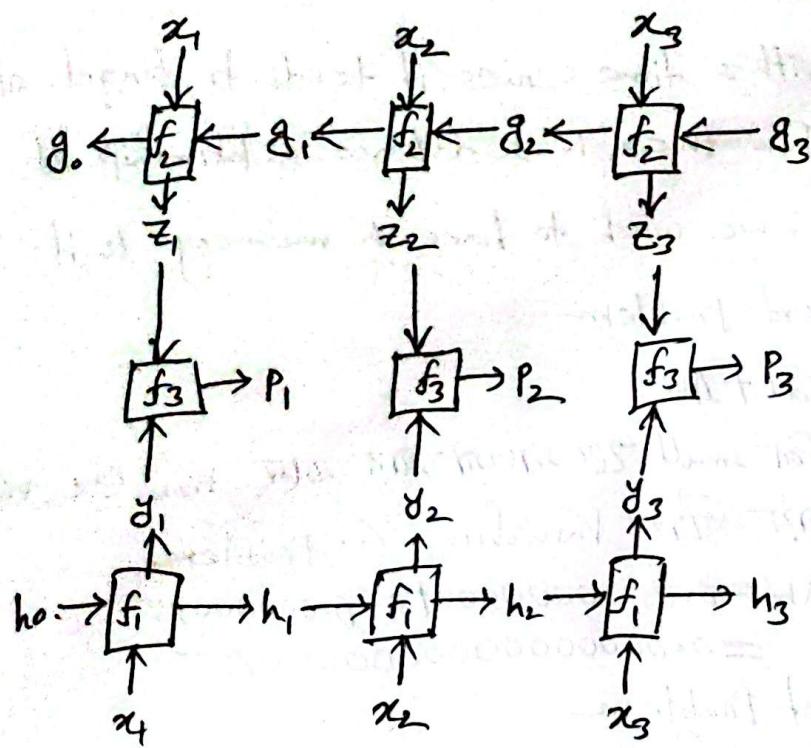
Bidirectional RNN

Use case → Palindrome checking

Normal RNN \Rightarrow Complexity হিচাবে Bidirectional RNN \Rightarrow Complexity হারাবাবা,

A* search

huge depth \Rightarrow
problem কঠো
Complexity O
 \Rightarrow solve কৰা যাবে



Pyramid RNN

Use case → Sequence to word generation

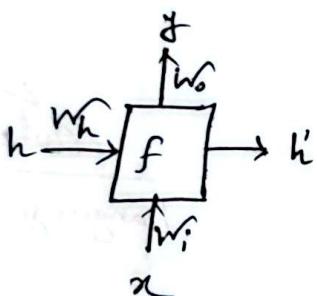
④ learning/training speed significantly faster

Recurrent steps Bidirectional RNN ଏହି ଟେକ୍ସ୍ ଓ ଗୋଟିକ ରାମିଳ୍ଲ ଦ୍ୱାସି

Example: Large vocabulary conversational speech recognition

Naive RNN

$$f: h, y = f(h, x)$$



y is computed from h'

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \sigma(2x)$$

$$\sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{1+e^x}$$

$$\sigma(2x) = \frac{e^{2x}}{1+e^{2x}}$$

$$h' = \sigma \left(\hat{w_h} \cdot h + \hat{w_i} \cdot x \right) + b$$

Activation function

$$y = \sigma(w \cdot h) + b$$

Softmax process

Problem

- When dealing with a time series, it tends to forget old information. When there is a distance relationship of unknown length, we wish to have a memory to it.
 - Vanishing Gradient Problem

$$\hat{w}_{new} = \hat{w}_{old} + \lambda w$$

Now Gurusakha's small \overline{GCF} \overline{GCD} \overline{GMR} - \overline{GCM} \overline{GMC} \overline{GCR} \overline{GCS} \overline{GCV}
 same \overline{GCF} \overline{GCD} , \overline{GCR} vanishing G. Problem.

Example: $\Delta W = 0.0000000001 \times 0.00000075$
 $= 0.000000000000000075$

- ## • Exploding Gradient Problem

LSTM

Neural Network Layer

Pointwise Operation

Matrix A
5xA

$A+B$

$$\begin{bmatrix} Ax \\ Ay \\ Az \end{bmatrix} + \begin{bmatrix} Bx \\ By \\ Bz \end{bmatrix} = \begin{bmatrix} Ax + Bx \\ Ay + By \\ Az + Bz \end{bmatrix}$$

Point wise Multiplication

$$\begin{bmatrix} Ax \\ Ay \\ Az \end{bmatrix} \odot \begin{bmatrix} Bx \\ By \\ Bz \end{bmatrix}$$

Matrix multiplication

possible $\text{if } (3 \times 1) \cdot (3 \times 1)$

wrong

& Pointwise multip.
possible

$$= \begin{bmatrix} Ax \cdot Bx \\ Ay \cdot By \\ Az \cdot Bz \end{bmatrix}$$

\odot vector

Point wise operation
Dot product

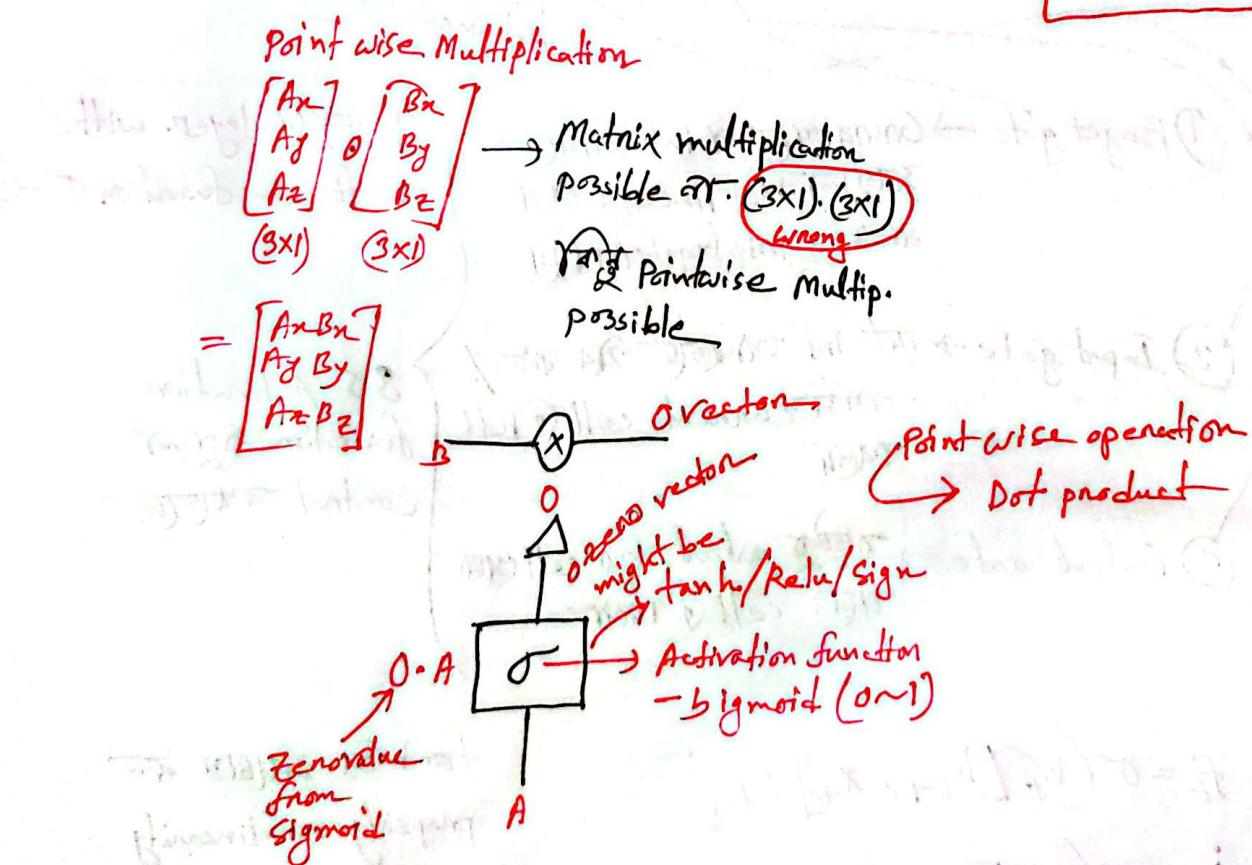
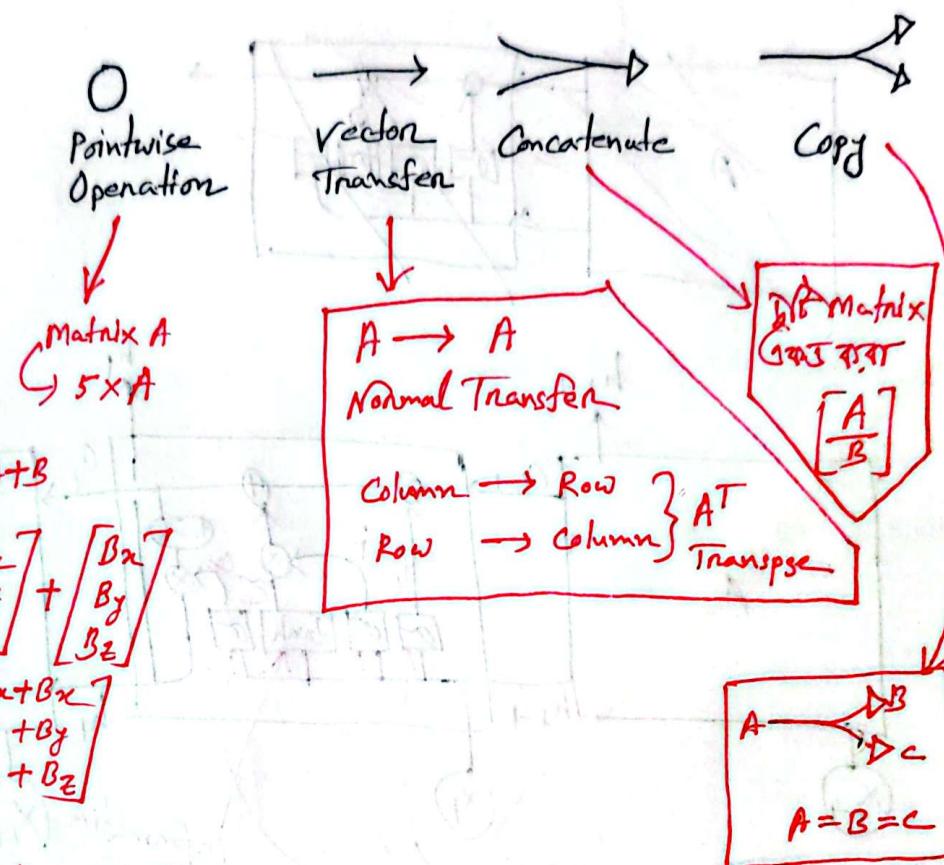
\odot vector
might be
tanh/Relu/Sign

Zeros
from
Sigmoid

$$0 \cdot A$$

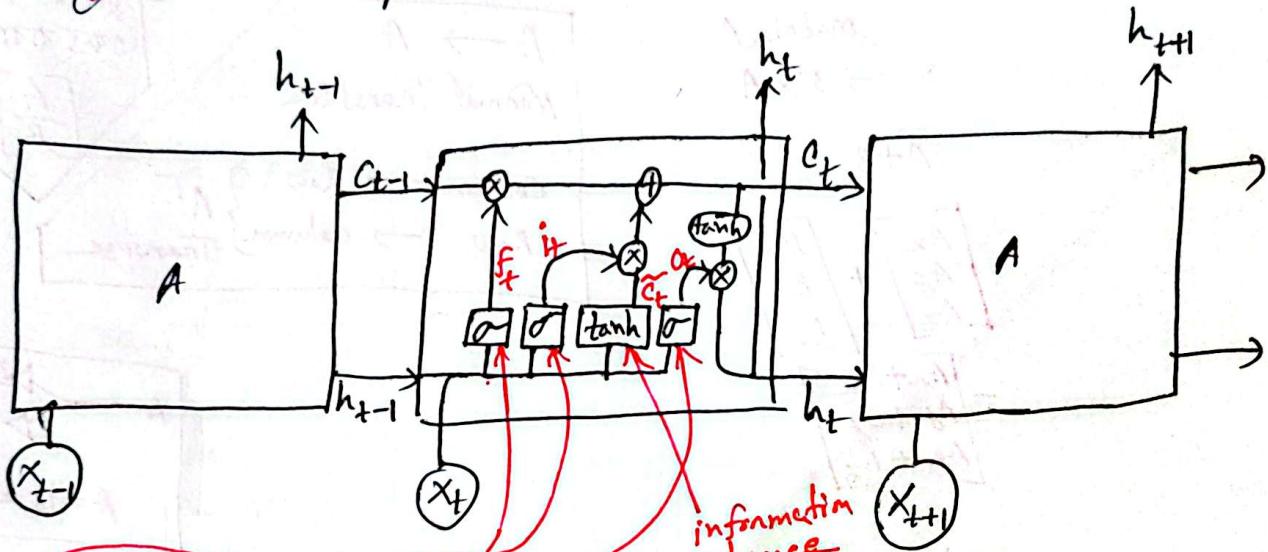
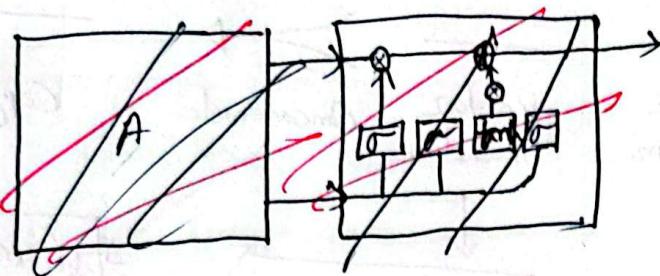
Activation function
- sigmoid (0~1)

A



The sigmoid layer outputs numbers between 0-1 determine how much each component should be let through.

LSTM



① Forget gate → গুরুত্ব দিতে নির্দিষ্ট সবচেয়ে প্রারম্ভিক পদক্ষেপ হল ফরগেট গেট। এটি পুরাতন জ্ঞান নির্দিষ্ট করে।

gate/NN layer with Activation function σ - ②

② Input gate → যদি জ্ঞান আপডেট করা হচ্ছে তবে এটি করে। এটি করে বর্তমান ক্ষেত্রে কাউন্ট করে।

3rd Activation function \tanh করে। Control - করে।

③ Output gate → সামুদ্রিক output করে। ক্ষেত্রে কাউন্ট করে।

$$f_t = \sigma(W_f [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i [h_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(W_{\text{information}} [h_{t-1}, x_t] + b_{\text{information}})$$

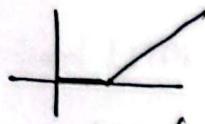
$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

$$h_t = o_t * \tanh(c_t)$$

\tanh করে ক্ষেত্রে কাউন্ট করে। non-linearity

Information enhanced and Shrink করে। করতে পারে।

ReLU



$$\max(z, 0)$$

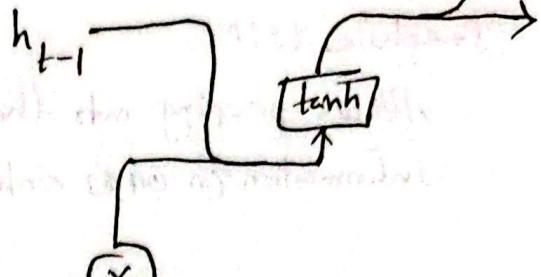
Linear
in manner

Information কে সুব গোলাত

enhanced / shrink করতে পার

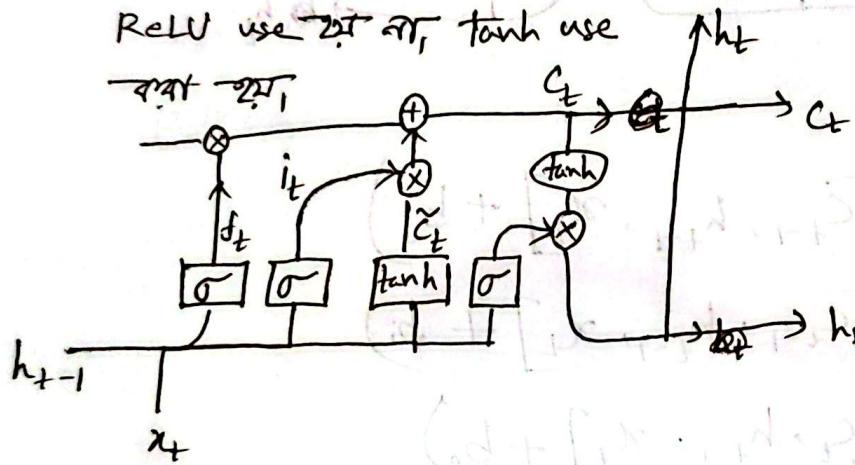
না, তবে LSTM কি-structure

ReLU use কর, না, tanh use



vanilla / Naive
RNN

- ① No cell state in RNN
- ② There is no sigmoid locking key.



LSTM

- ① Computationally LSTM হ্যাঁ একটি expensive but LSTM can hold a longer memory
- ② It is able to solve vanishing gradient problem

$$f_t = \sigma(w_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(w_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(w_c \cdot [h_{t-1}, x_t] + b_c)$$

$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

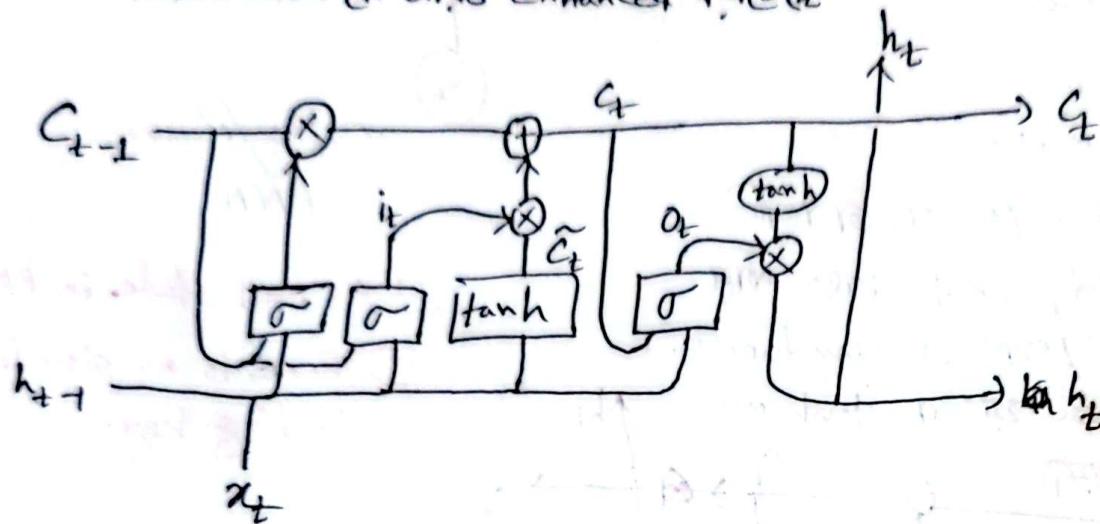
Information
Update

Updating
the cell state

$$\begin{cases} o_t = \sigma(w_o \cdot [h_{t-1}, x_t] + b_o) \\ h_t = o_t * \tanh(c_t) \end{cases}$$

Peephole LSTM

Allows peeping into the memory.
information can easily enhanced through



$$f_t = \sigma(W_f [C_{t-1}, h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i [C_{t-1}, h_{t-1}, x_t] + b_i)$$

$$O_t = \sigma(W_o [C_t, h_{t-1}, x_t] + b_o)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$h_t = O_t * \tanh(C_t)$$

Naive RNN VS LSTM

C changes slowly $\rightarrow C^t$ is C^{t-1} added by something

h changes faster $\rightarrow h^t$ and h^{t-1} can be very different

GRU

* ~~LSTM~~ GRU is computationally faster.

* GRU and RNN have cell memory.

GRU and RNN hypothesis is that they are same.

* To omit Vanishing Gradient Problem, GRU is being introduced.

2 gates in GRU

① Update gate (z_t)

② Reset gate (r_t)

Variables

$x_t \rightarrow$ Input vector

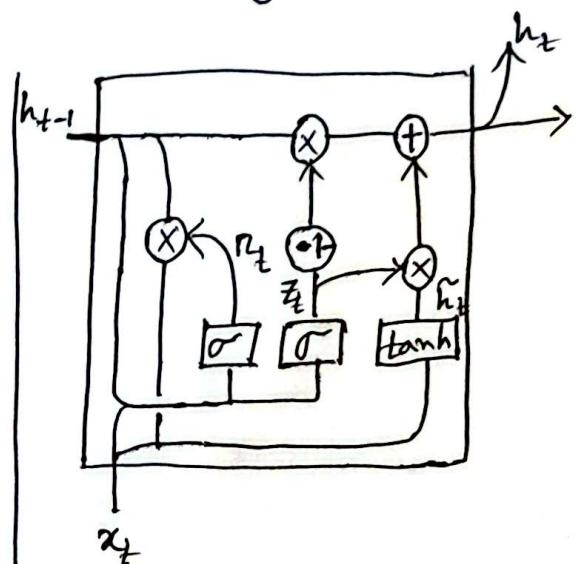
$h_t \rightarrow$ Output vector

$\tilde{h}_t \rightarrow$ Candidate activation vector

$z_t \rightarrow$ Update gate vector

$r_t \rightarrow$ Reset gate vector

$W, U \rightarrow$ Weight matrices; $b \rightarrow$ bias



Update gate

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z)$$

Reset gate

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

$$x_t \rightarrow (5 \times 1)$$

$$\begin{matrix} W_z x_t \\ V_z h_{t-1} \\ b_z \end{matrix} \left. \right\} \begin{matrix} \text{dimension} \\ \text{same} \\ \text{example: } (5 \times 1) \end{matrix}$$

$$\begin{matrix} z_t \\ r_t \end{matrix} \left. \right\} (5 \times 1)$$

$$\tilde{h}_t = \tanh(W_h x_t + U_h (h_{t-1} \cdot r_t) + b_h)$$

$$h_t = (1 - z_t) \cdot h_{t-1} + (z_t \cdot \tilde{h}_t)$$

h_{t-1} \rightarrow $h_t \times$ dimension
same as h_t

same as h_t

Rajin