

## 6) Reinforcement Learning.

- \* Reinforcement learning is a feedback-based machine learning technique in which an agent learns to behave in an environment by performing actions & seeing the results of actions.
- For each good action, the agent gets +ve feedback, & for each bad action, the agent gets -ve feedback or penalty.
- ex) game-playing, robotics.

### \* Elements of RL:-

#### 1) Policy :-

- Policy defines the learning agent behavior for given time period.
- It is mapping from perceived states of env. to actions to be taken when in those states.

#### 2) Reward function:-

- Reward function is used to define a goal in a reinforcement learning problem.
- A reward fn is a function that provides a numerical score based on the state of env.

#### 3) Value function:-

- specify what is good in long run.
- the value of state is total amount of reward an agent can expect to accumulate over the future, starting from that state.

#### 4) Model:-

- with the help of model, one can make inference about how the environment will behave.

-such as, if a state and action are given, then a model can predict the next state

shaded at 20 & 21, so it is suppose to be

initially probability of transitioning to s1

\* Discount Factor:- ( $\gamma$ ) :- Ranges from 0 to 1

cost of travelling from one state to

another

↳  $\gamma = 0.5$

Ex)

$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$
100	0	0	0	0	40
1	2	3	4	5	6

consider:  $\gamma = 0.5$

From 1:-

i) To left:-

Return = 100

ii) To right:-

Return =  $0.5 \cdot 40 = 20$

From 2:-

i) To left:-

Return =  $0.5 \cdot 100 = 50$

ii) To right:-

Return =  $0.5 \cdot 50 = 25$

2) To right:-

$$\text{Return} = \gamma^4 \cdot 40$$

$$= 0.5^4 \cdot 40 = 1.25$$

To traverse total =  $2 + 5 + 3 = 10$

Successor of  $s_4$  are  $s_5$  and  $s_6$

\* State-action value function  $\alpha(s,a)$  is

$$\alpha(s,a) = \text{Return}$$

$s$  = state  $a$  = action

optimal behaviour!:-

$\rightarrow$  Ans normal +

100	5	5	5	5	40	$\gamma = 0.5$
-----	---	---	---	---	----	----------------

(0.2) two half of 876H = 438

$$\begin{aligned}
 1) Q(3, \rightarrow) &= 0 + \gamma \cdot 0 + \gamma^2 \cdot 0 + \gamma^3 \cdot 0 + \gamma^4 \cdot 100 \\
 &= \gamma^4 \cdot 100 \\
 &= 0.5^4 \cdot 100 \\
 &= 6.25
 \end{aligned}$$

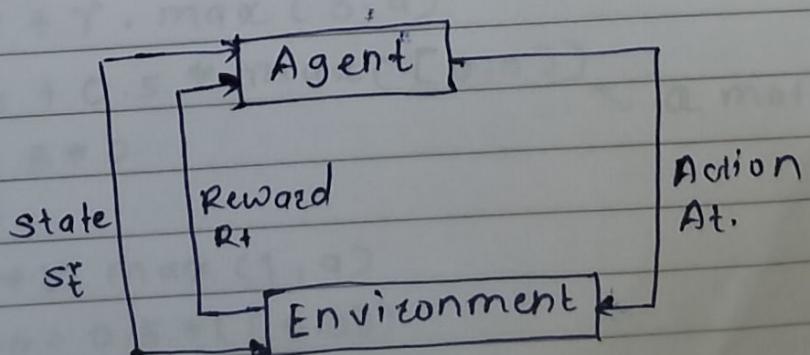
$$\begin{aligned}
 2) Q(5, \leftarrow) &= \gamma^4 \cdot 100 \text{ or } \text{Bridges} \\
 &= 0.5^4 \cdot 100
 \end{aligned}$$

$$\begin{aligned}
 3) Q(3, \leftarrow) &= \gamma^2 \cdot 100 \\
 &= 0.5^2 \cdot 100
 \end{aligned}$$

Now,  
 $\max Q(s, a)$  for all  $a$  it is called  $\underline{Q}^*$

\* Markov decision process:-

- a mathematical framework used for modelling decision-making processes in situations where outcomes are partly random & partly under the control of a decision-maker.



\* Bellman eqn: -

- Helps to find out Q(s,a)

$$Q(s,a) = R(s) + \gamma \max_{a'} Q(s', a')$$

$$0.1 \cdot R + 0.8 \cdot 0.5 + 0.8 \cdot R + 0.1 \cdot R + 0.8 \cdot 0 = (0.1 \cdot R + 0.8 \cdot 0.5) + 0.8 \cdot (0.1 \cdot R + 0.8 \cdot 0.5)$$

Ex)

	1	2	3	4	5	6	$\gamma = 0.5$
100	50	12.5	25	6.25	12.5	10	6.25
0	0	0	0	0	0	40	

check: checking values using Bellman eqn.

$$\textcircled{1} \quad Q(3, \rightarrow) = R(3) + \gamma \max([12.5, 10])$$

$$= 0 + 0.5 * 12.5$$

$$= 6.25$$

$$\textcircled{2} \quad Q(4, \rightarrow) = R(4) + 0.5 * \max([6.25, 20])$$

$$= 0 + 0.5 * 20$$

$$\textcircled{3} \quad Q(5, \leftarrow) = R(5) + 0.5 * \max([12.5, 10])$$

$$= 0 + 0.5 * 12.5$$

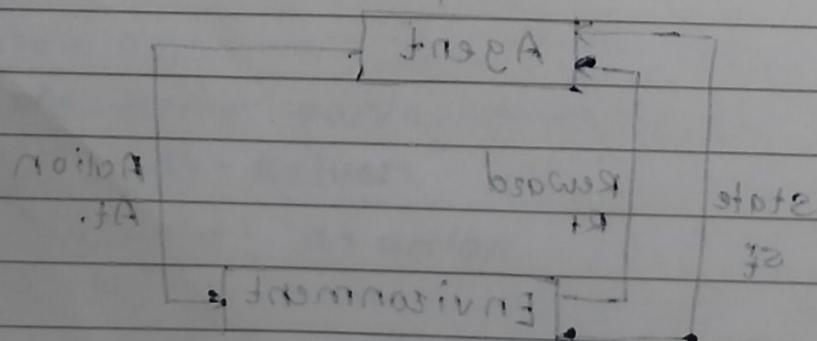
$$= 6.25$$

so below is 10 & 10 so 0.5 \* 10 = 5

so 6.25 & 6.25 so 0.5 \* 6.25 = 3.125

so 10 & 10 so 0.5 \* 10 = 5

so 3.125 & 5 so 0.5 \* 3.125 = 1.5625



## \* Q-Learning :-

- enables model to iteratively learn & improve over time by taking the correct action.

Ex)

Reward matrix:

	1	2	3	4	5	6
1	-1	0	-1	-1	-1	-1
2	100	-1	0	-1	-1	-1
3	-1	0	-1	-1	-1	-1
4	-1	0	-1	20	-1	-1
5	-1	-1	-1	0	40	-1
6	-1	-1	-1	0	0	-1

i) Initialize with all zeros ( $Q$ -matrix)

	1	2	3	4	5	6
1	0	0	0	0	0	0
2	<u>100</u>	0	0	0	0	0
3	0	<u>50</u>	0	0	0	0
4	0	0	<u>25</u>	0	0	0
5	0	0	0	<u>12.5</u>	0	<u>40</u>
6	0	0	0	0	0	0

Now, by Bellman eqn.

$$Q(2,3) = R + \gamma \cdot \max(3, 9)$$

$$= 0 + 0.5 * \max([0, 0]) \quad \leftarrow Q\text{-matrix}$$

$$= 0.5 * 0$$

$$= 0.$$

$$Q(2,1) = R + \gamma \cdot \max(1, 9)$$

$$= 100 + 0.5 * ([0, 0])$$

$$= 100 + 0 = 100 \text{ update}$$

$$Q(3,2) = R + \gamma \cdot \max(2, a)$$

$$= 0 + 0.5 * \max([100, 0])$$

$$= 0.5 * 100 = 50 \text{ - update.}$$

$$Q(4,3) = R + \gamma \cdot \max(3, a)$$

$$= 0 + 0.5 * \max([50, 0])$$

$$= 0.5 * 50$$

$$= 25 \text{ - update.}$$

$$Q(5,4) = R + \gamma \cdot \max(4, a)$$

$$= 0 + 0.5 * \max([25, 0])$$

$$= 0.5 * 25$$

$$= 12.5 \text{ - update.}$$

$$Q(5,6) = R + \gamma \cdot \max(5, a)$$

$$= 40 + 0.5 * \max([0, 0])$$

$$= 40 + 0$$

$$= 40 \text{ - update.}$$

0	0	0	0	0	0	0	1
5	0	0	0	0	0	0	1
0	0	0	0	0	0	0	1
0	0	0	0	0	0	0	1

new possible path

$$(0, 8) \xrightarrow{\gamma = 0.5} R + \gamma \cdot \max(8, a) = (8, 8)$$

$$\xrightarrow{(0, 0)} (0, 0) \xrightarrow{\gamma = 0.5} R + \gamma \cdot \max(0, 0) = 0 + 0 = 0$$

$$0 + 0 = 0$$

$$0 + 0 = 0$$

$$(0, 1) \xrightarrow{\gamma = 0.5} R + \gamma \cdot \max(1, a) = (1, 1)$$

$$\xrightarrow{(0, 0)} (0, 0) \xrightarrow{\gamma = 0.5} R + \gamma \cdot \max(0, 0) = 0 + 0 = 0$$

$$0 + 0 = 0 + 0 = 0$$

## MOD 4: Supervised learning Advanced.

CLASSMATE

Date \_\_\_\_\_

Page \_\_\_\_\_

class imbalance problem: →

- this problem occurs when one class in a classification problem has significantly fewer samples than the others.

- this can lead to the model being biased towards the majority class, resulting in poor performance, especially for the minority class.

Example: How many times did

suppose we have a dataset with 1000 transactions, out of which only 20 are fraudulent & 980 are non-fraudulent. This represents a class imbalance prob.

\* ~~other~~ more on evaluation matrix:-

1) Sensitivity:-

- known as Recall, True positive rate

- measures the proportion of actual true cases that were correctly identified by the model.

$$\text{sensitivity} = \frac{TP}{TP + FN}$$

2) Specificity:- (True negative rate)

- measures the proportion of actual -ve cases that were correctly identified by the model.

$$\text{specificity} = \frac{TN}{TN + FP}$$

### 3) Area Under the ROC curve (AUC-ROC) :-

- the ROC curve is graphical representation of the true positive rate vs. false positive rate for different classification thresholds.
- AUC-ROC quantifies the model's ability to distinguish between classes.
- An AUC-ROC score of 1 indicates perfect model, while score 0.5 indicates random model.

#### \* Logistic regression with single neuron :-

edit part - supervised learning algorithm used for classification. It divides into two categories using sigmoid fn.

#### \* Sigmoid fn

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

edit output to no output -

#### Logistic regression model :-

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

where,

$x$  = input feature vector.

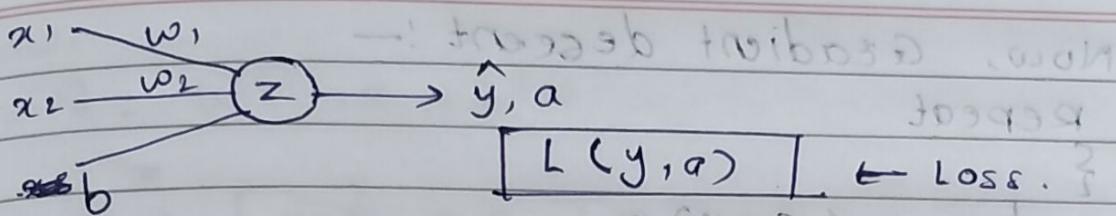
$w$  = weight vector

$b$  = bias

$z$  = combination of input features & weights.

$\sigma$  = sigmoid fn

$a, \hat{y}$  = predicted output.



So,

$$z = w_1 x_1 + w_2 x_2 + b \quad (z - a) = db$$

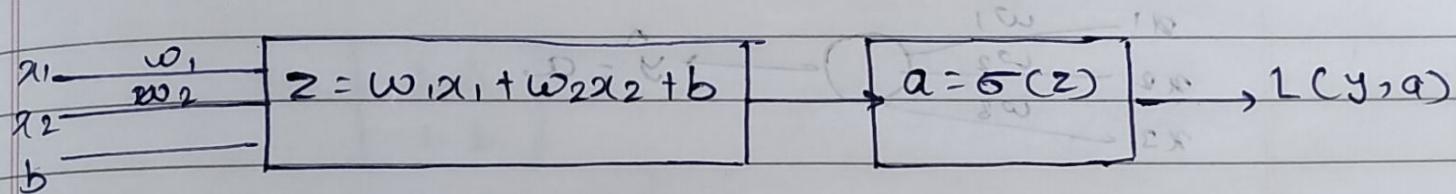
$$a = \sigma(z) \quad \sigma(b) = aw$$

$$\omega b x_1 + \omega x_2 = \omega w$$

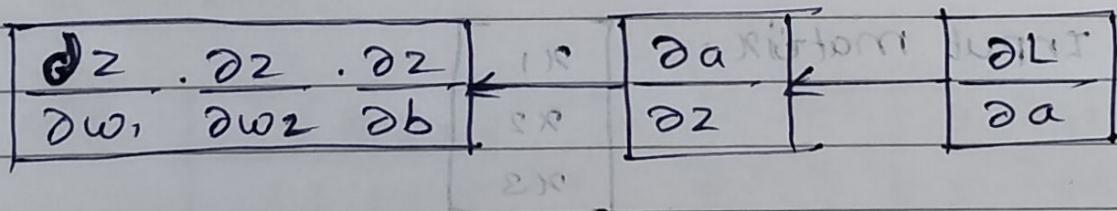
$$w_1 b x_1 + w_2 x_2 = \omega w$$

$$L(\hat{y}, a) = -y \log_2 a - (1-y) \log_2 (1-a)$$

It can be shown like this:-



For back propagation, = gradient descent



Updating weights :-

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_1} \quad (x_1)$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_2} \quad (x_2)$$

$$\frac{\partial L}{\partial w_1} = dw_1, \quad \frac{\partial L}{\partial w_2} = dw_2, \quad \frac{\partial L}{\partial b} = db.$$

$$\begin{bmatrix} 5.0 & 2.0 \end{bmatrix} = [s] \omega \quad \begin{bmatrix} 2.0 & 2.0 \end{bmatrix} = [s] \omega \quad \begin{bmatrix} 1.0 & 1.0 \end{bmatrix} = [s] \omega$$

$$\begin{bmatrix} 8.0 \end{bmatrix} = [s] \omega \quad \begin{bmatrix} 1.0 \end{bmatrix} = [s] \omega \quad \begin{bmatrix} 1.0 \end{bmatrix} = x$$

Now, Gradient descent! —

Repeat

$$\{ \rightarrow (w, b) \}$$

$$dw_1 = (a - y)x_1$$

$$dw_2 = (a - y)x_2$$

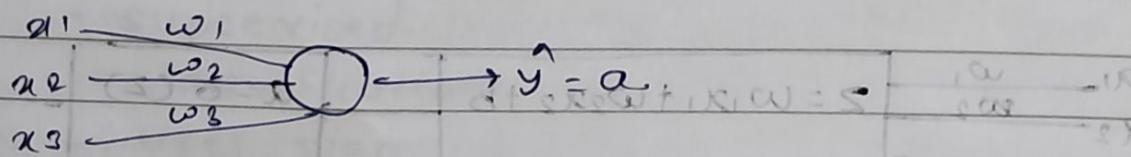
$$db = (a - y)$$

$$w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$b = b - \alpha db$$

?.

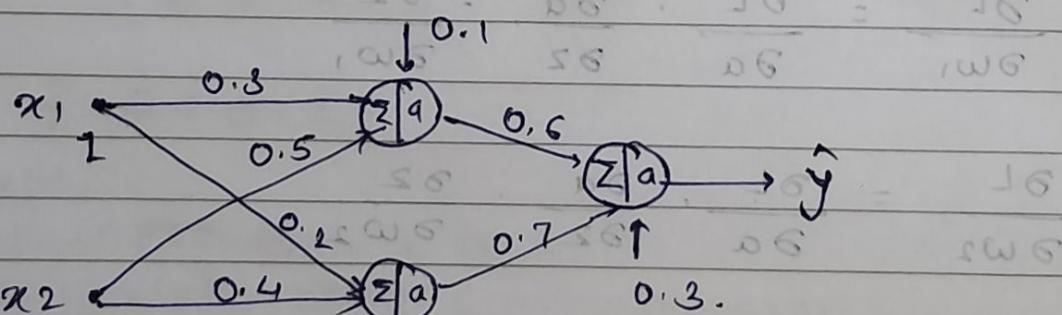


Weight matrix =  $[w_1, w_2, w_3]$

$$\text{Input matrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 56 \\ 56 \\ 56 \end{bmatrix}$$

$\therefore$  Weight vector

Ex)



$$db^2 = 0.6 \quad \omega^1 = 0.2 \quad \omega^2 = 0.6$$

$$\omega^{[1]} = \begin{bmatrix} 0.8 & 0.5 \\ 0.2 & 0.4 \end{bmatrix} \quad \omega^{[2]} = \begin{bmatrix} 0.1 & 0.6 \\ 0.3 & 0.7 \end{bmatrix}$$

$$x = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix} \quad b^{[2]} = [0.6]$$

Now,

$$z^{[1]} = w^{[1]} \cdot x + b^{[1]}$$

=

$$= \begin{bmatrix} 0.3 & 0.5 \\ 0.2 & 0.4 \end{bmatrix} * \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}$$

$$= \begin{bmatrix} 0.3+1 \\ 0.2+0.8 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix} = \begin{bmatrix} 1.4 \\ 1.2 \end{bmatrix}$$

$$= \begin{bmatrix} 1.4 \\ 1.2 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix} = \begin{bmatrix} 1.5 \\ 1.4 \end{bmatrix}$$

Now apply sigmoid,

$$a^{[1]} = \frac{1}{1 + e^{-1.4}} = \begin{bmatrix} 0.8021 \\ 0.19785 \end{bmatrix}$$

Now,

$$z^{[2]} = w^{[2]} \cdot a^{[1]} + b^{[2]}$$

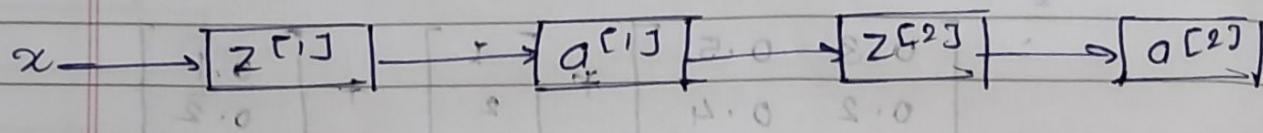
$$= \begin{bmatrix} 0.6 & 0.7 \end{bmatrix} * \begin{bmatrix} 0.8021 \\ 0.19785 \end{bmatrix} + \begin{bmatrix} 0.3 \\ 0.3 \end{bmatrix}$$

$$= [0.6 * 0.8021 + 0.7 * 0.19785] + [0.3]$$

$$= 1.31$$

$$a^{[2]} = \frac{1}{1 + e^{-1.31}} = 0.7885 \therefore \hat{y} = 0.7885$$

Now,

For backpropagation,  $y = x \cdot w + b$ 

$$z^{[2]} = w^{[2]T} \cdot a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

$$L = -y \log_2 a^{[2]} - (1-y) \log_2 (1-a^{[2]})$$

$$dz^{[2]} = a^{[2]} - y$$

$$d\omega^{[2]} = dz^{[2]} \cdot a^{[1]}$$

$$db^{[2]} = dz^{[2]}$$

Now,

$$z^{[1]} = w^{[1]T} \cdot x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$L = -y \log_2 a^{[1]} - (1-y) \log_2 (1-a^{[1]})$$

$$dz^{[1]} = a^{[1]} - y w^{[2]}. d\omega^{[2]} * g^{[1]}'(z^{[1]})$$

$$d\omega^{[1]} = dz^{[1]} \cdot x^T$$

$$db^{[1]} = dz^{[1]}$$

$$w_{new}^{[i]} = w_{old}^{[i]} - \alpha d\omega^{[i]}$$

$$b_{new}^{[i]} = b_{old}^{[i]} - \alpha db^{[i]}$$

\* For above example,

If  $y_{actual} = 1, \alpha = 0.1$ . Find  $w_{new}^{[1]}, w_{new}^{[2]}$ ,

$$b_{new}^{[1]} = b_{old}^{[1]} + 1808.0 * 0.1$$

→ Finding  $w_{new}^{[2]}$ :

$$dz^{[2]} = a^{[2]} - y$$

$$= 0.78 - 1.0$$

$$= -0.22$$

$$\begin{aligned}
 dw_1^{[2]} &= dz^{[2]} \cdot a_1^{[1]} & dw_2^{[2]} &= dz^{[2]} \cdot a_2^{[1]} \\
 &= -0.22 * 0.8021 & & = -0.22 * 0.7685 \\
 &= -0.1764 & & = -0.1690
 \end{aligned}$$

$$\begin{aligned}
 db^{[2]} &= dz^{[2]} \\
 &= -0.1764 - 0.22
 \end{aligned}$$

Now,

$$\begin{aligned}
 w_1^{[2] \text{ new}} &= w_1^{[2] \text{ old}} - \alpha dw_1^{[2]} \\
 &= 0.6 - 0.1(-0.1764) \\
 &= 0.6176
 \end{aligned}$$

$$\begin{aligned}
 w_2^{[2] \text{ new}} &= w_2^{[2] \text{ old}} - \alpha dw_2^{[2]} \\
 &= 0.7 - 0.1(-0.1690) \\
 &= 0.7169
 \end{aligned}$$

$$\begin{aligned}
 b^{[2] \text{ new}} &= b^{[2] \text{ old}} - \alpha db^{[2]} \\
 &= 0.3 - 0.1(-0.22) \\
 &= 0.322
 \end{aligned}$$

$$\therefore w^{[2] \text{ new}} = [0.6176 \quad 0.7169]$$

$$b^{[2] \text{ new}} = [0.322]$$

Finding  $w^{[1] \text{ new}}$  &  $b^{[1] \text{ new}}$ ,

$$\begin{aligned}
 dz^{[1]} &= w^{[2] \cdot dz^{[2]} + g^{[1]'}(z^{[1]})} \\
 &= [0.6 \quad 0.7] * (-0.22) * \begin{bmatrix} 0.8021 \\ 0.7685 \end{bmatrix} \\
 &\approx [-0.132 \quad -0.154] * \begin{bmatrix} 0.8021 \\ 0.7685 \end{bmatrix} \\
 &\approx -0.2242
 \end{aligned}$$