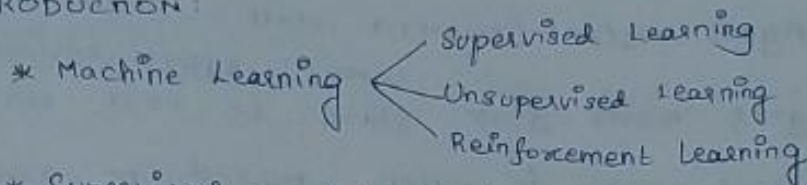


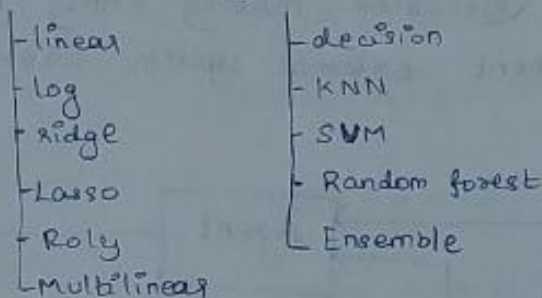
22/07/25

U - 1: INTRODUCTION

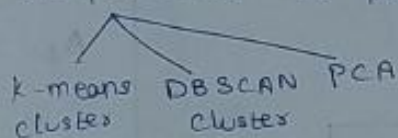
INTRODUCTION:



* Supervised - dataset Input data mapped to output
labelled data, regression, classification

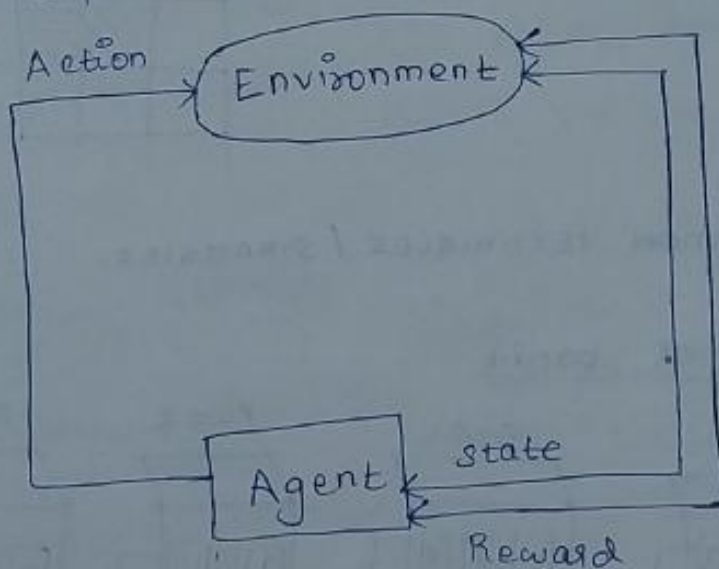


* Unsupervised - no proper label / map, input & output



* Why RL? → no data pts needed (model learn by itself)
→ too much data not knowing reliability
→ not practical for real env
→ Learn by trial & error

1. Exploration 2. Exploitation



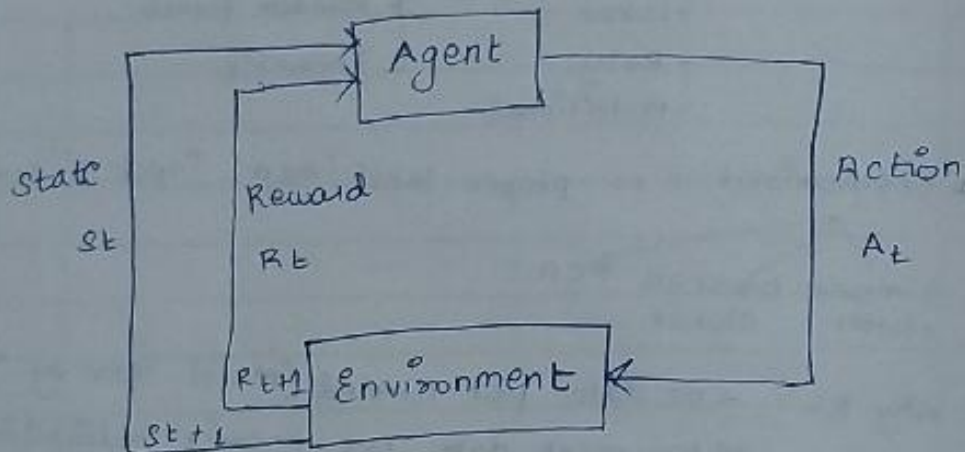
RL is
sequence
decision-making

Exploration:

Exploitation:

- * **State**: Agent resides / current state
- * **Action**: agent moves from start to next state ^{choices}, affects environment
- * **Reward**: agent gets in return, based on goals / actions, maybe +ve, -ve
- * **Agent**: Decision-making entity that acts - max rewards
- * **Environment**: External system where agent interacts and learns

24/07/25



- * **Policy**: strategy followed to maximize the rewards.



EXPLORATION TECHNIQUES / STRATEGIES:

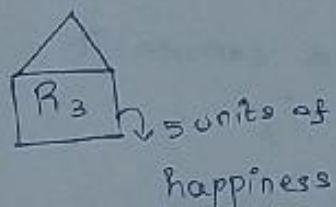
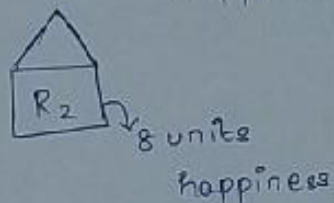
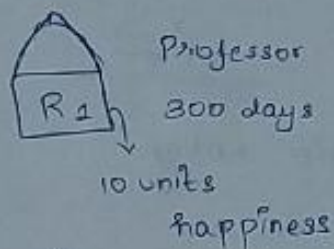
1. One armed bandit



- * Without any previous knowledge, trying everything, trial - exploration

2. Multi-armed bandit

① only exploration



Professor explores the 3 restaurants for 300 days

$$100 \times 10 + 100 \times 8 + 100 \times 5$$

$$= 1000 + 800 + 500$$

$$= 2300 \text{ units of happiness}$$

$$\text{max}_{\text{happ}} \rightarrow \text{happ} [3000 - 2300]$$

$$700 = 700 \text{ units}$$

② only Exploitation

* 297 days, he uses the same restaurants which he has explored

* first day R_1 , second day R_2 , third day R_3 . that day then he decides and uses those restaurant

③ ϵ - greedy strategy

$$\epsilon = 10\%$$

300 days \rightarrow 10 days - R_1
10 days - R_2
10 days - R_3

best 270 days selected restaurant

29/07/25

\rightarrow finding many possible solution - exploration

\rightarrow finding one solution and using it repeatedly without

finding other solutions - exploitation.

$\epsilon \rightarrow$ exploration

$1 - \epsilon \rightarrow$ exploitation

Q - LEARNING ALGORITHM:



for each S , initialize table entry

$\hat{Q}(S, a)$ to zero

observe current state S

Do forever

- Select an action 'a' and execute it
- Receive immediate reward 'r'
- observe the new state S'
- update the table new entry for $\hat{Q}(S, a)$

$$\hat{Q}(S, a) \leftarrow \gamma + r + \max_{\hat{a}} \hat{Q}(S', \hat{a})$$

$S \leftarrow S'$ discounting factor

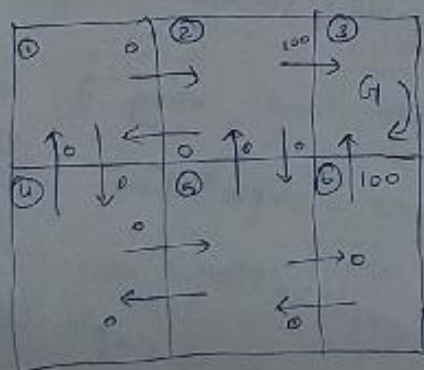
can be 0 to [0.9]

Eg:-

* higher the Q value - closer to the goal

* Grid, goal is door, action is up, down, right, left.

* Each action gives a Q value, higher the Q value, more close to reaching the goal.



near the goal
Higher the reward
so value near goal is 100
and rest is 0

$Q(S, a)$ reward values

Reward Matrix 6x6

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

* no connection b/w

(3,1)/(1,3) so [-1]

* reward matrix is 100

* connected room (1,2)/

(3,2) ... but no

reward so zero.

* initialize Q ^{matrix} value to zero.

	1	2	3	4	5	6
1	0	0	0	0	0	0
2	0	0	0	0	0	0
3	0	0	0	0	0	0
4	0	0	0	0	0	0
5	0	0	0	0	0	0
6	0	0	0	0	0	0

→ Q table

$\alpha(=0.9)$

$$Q(s, a) = R(s, a) + \gamma \max_b [Q(s, b)]$$

$$\text{value} = 100 + (0.9)(0) = 100$$

then update Q table

updated Q table

0	0	0	0	0	0
0	0	100	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0

value from
reward
matrix
then value from
Q table

$$\begin{aligned}
 Q(6,3) &= R(6,3) + \gamma \max [Q(3,3)] \\
 &= 100 + (0.9)(0) \\
 &= 100
 \end{aligned}$$

update (6,3) value in Q table.

2

updated

	1	2	3	4	5	6
1	0	90	0	0	0	0
2	81	0	100	0	0	0
3	0	0	0	0	0	0
4	0	0	0	0	0	0
5	0	0	0	0	0	0
6	0	0	100	0	0	0

$$\begin{aligned}
 Q(1,2) &= R(1,2) + \gamma \max [Q(2,1), Q(2,3), Q(2,5)] \\
 &= 0 + 0.9 \max [0, 100, 0] \\
 &= 0 + 0.9 \times 100 \\
 &= 90
 \end{aligned}$$

update $Q(1,2)$ as 90 in Q table.

see updated table.

$$\begin{aligned}
 Q(2,1) &= R(2,1) + \gamma \max [Q(1,2), Q(1,4)] \\
 &= 0 + 0.9 \max [90, 0] \\
 &= 0 + 0.9 \times 90 \\
 &= 81
 \end{aligned}$$

$$\begin{aligned}
 Q(1,4) &= R(1,4) + \gamma \max [Q(4,1), Q(4,5)] \\
 &= 0 + 0.9 \max [0, 0] \\
 &= 0
 \end{aligned}$$

$$Q(5,2) = R(5,2) + \gamma \max [Q(2,1), Q(2,3), Q(2,5)]$$

$$= 0 + 0.9 \max [81, 100, 0]$$

$$= 0 + 0.9 \times 100$$

$$= 90$$

$$Q(5,4) = R(5,4) + \gamma \max [Q(4,1), Q(4,5)]$$

$$= 0 + 0.9 \max [0]$$

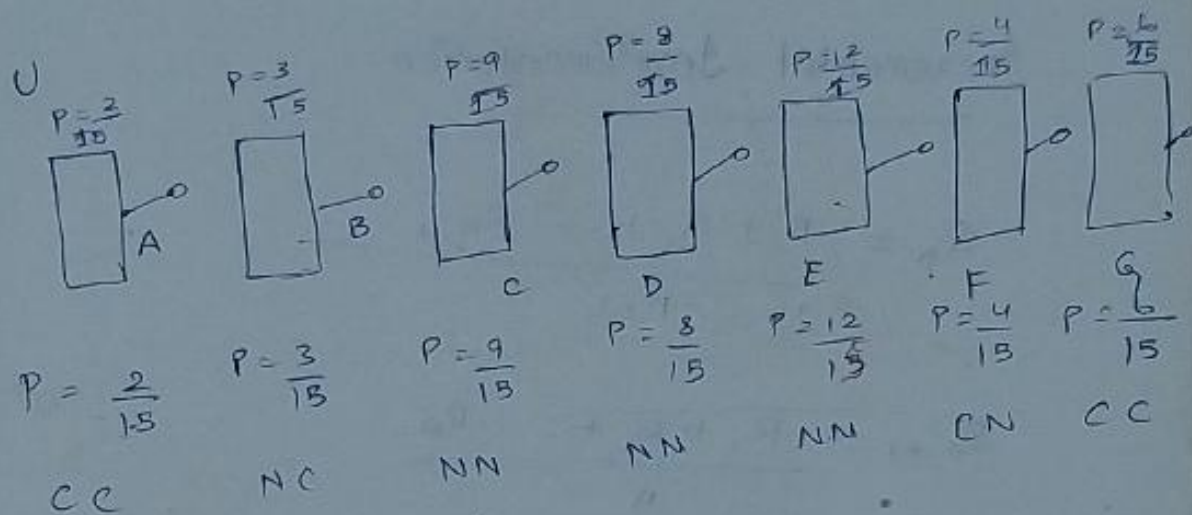
$$= 0$$

Final updated Q value after 30

	1	2	3	4	5	6
1	0	90	0	0.72	0	0
2	81	0	100	0	81	0
3	0	0	0	0	0	0
4	81	0	0	0	81	0
5	0	90	0	72	0	90
6	0	0	100	0	81	0

Model based RL | Model free RL

22/7/25

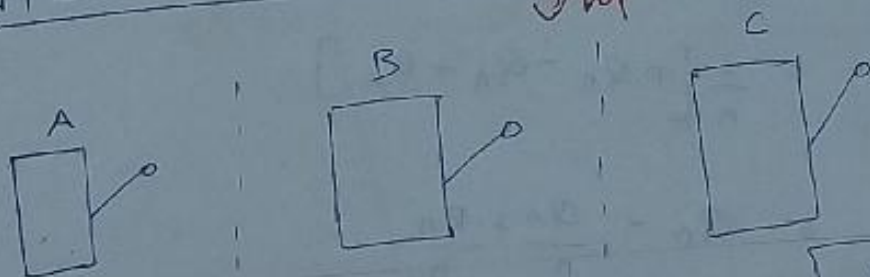


coin, No coin. \downarrow [Multi-armed bandit problem]

$\epsilon \rightarrow$ exploration

10-1. $1 - \epsilon [0.9]$
 \downarrow
 based on the result \rightarrow exploitation

UPPER CONFIDENCE BOUND \rightarrow [note]



bound value \leftarrow $N = n_A + n_B + n_C$ (no. of times playing for machine A, B, C)

* Bound + Estimation Q value should be maximum

* Deterministic approach [fixed action at a state]

($\ln \rightarrow$ log arithmetic)

Incremental Implementation:

$$Q_n = \frac{R_1 + R_2 + \dots + R_{n-1}}{n-1}$$

$$Q_{n+1} = \frac{R_1 + R_2 + \dots + R_n}{n}$$

$$Q_{n+1} = \frac{1}{n} \sum_{i=1}^n R_i = \frac{1}{n} \sum_{i=1}^{n-1} R_i + \frac{R_n}{n}$$

$$= \frac{1}{n} \times \frac{n-1}{n-1} \sum_{i=1}^{n-1} R_i + R_n \Rightarrow \frac{n-1}{n-1} \sum_{i=1}^{n-1} R_i$$

\swarrow
 Q_n

$$= \frac{1}{n} [(n-1) \times Q_n + R_n]$$

$$= \frac{1}{n} [nQ_n - Q_n + R_n]$$

$$= Q_n - \frac{Q_n}{n} + \frac{R_n}{n}$$

$$Q_{n+1} = Q_n + \frac{1}{n} [R_n - Q_n]$$

ϵ
Epsilon greedy

G
0.1

30°
10°
1-ε
1-0.1
0.9