

We find value iteratively till it converges.

→ discount factor
classmate

$$v_{11} = 0.9$$

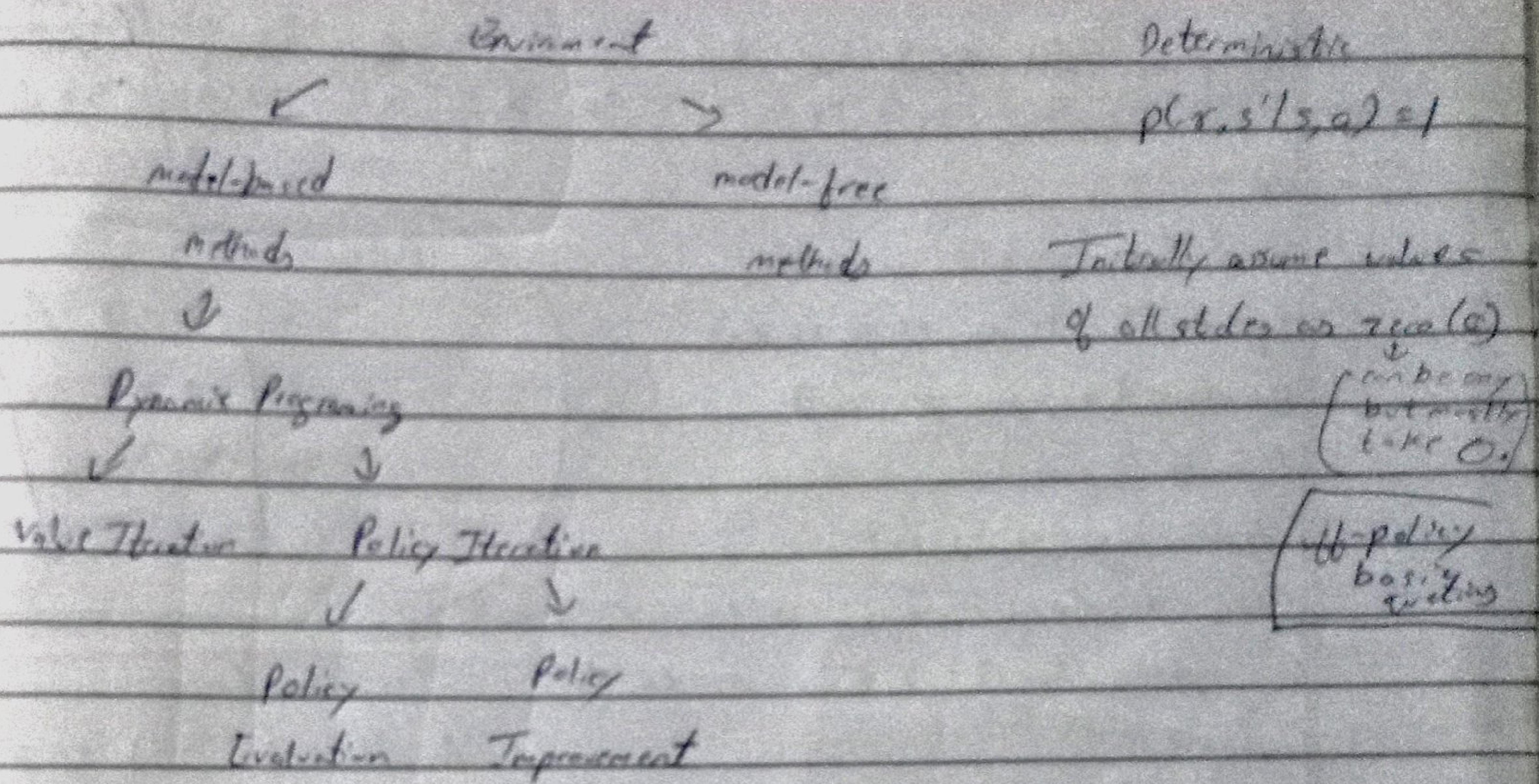
$$v_{12} = 0.9$$

$$v_{13} = 0.9$$

} thus no change, so stop.

Date _____

Page _____



Q. Why Bellman eq. are recursive in nature?

↳ grid problem ($T \in \mathbb{R}^{S \times S \times A}$) (ESE - Lecture)

↳ reward + maximum step to reach the goal.

Q. Give a 2x2 matrix- ($\gamma = 0.9$)

States Rewards

s_1 -1

s_2 -1

s_3 -1

s_4 10

3	6
3	6

Sol- Policy is Deterministic

Let us assume the following actions (a₁ and a₂)

Current state	Action	Next state
---------------	--------	------------

s_1

a₁

s_2

s_1

a₂

s_3

s_2

a₃

s_4

s_3

Terminal

\therefore formula :-

$$v_n(s) = \sum_{s'} (\gamma + r v_n(s'))$$

- Iteration - 1

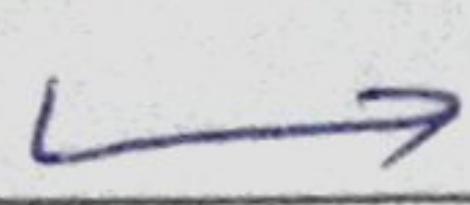
$$v_n(s_0) = \gamma + r v_n(s')$$

$$\therefore v_n(s_1) = -1 + 0.9 \times 0 = -1$$

$$v_n(s_2) = -1 + 0.9 \times 0 = -1$$

$$v_n(s_3) = -1 + 0.9 \times 0 = -1$$

$$v_n(s_4) = 10 \cancel{-1} = 10$$



Let Assume the value of all states as 0.

$$\therefore v_n(s_1) = 0$$

$$v_n(s_2) = 0$$

$$v_n(s_3) = 0$$

$$v_n(s_4) = 0$$

T- 2

$$v_n(s_1) = -1 + 0.9 \times (-1) = -1.9$$

$$v_n(s_2) = -1 + 0.9 \times (-1) = -1.9$$

$$v_n(s_3) = -1 + 0.9 \times \cancel{10} = \cancel{-8}$$

$$v_n(s_4) = 10$$

T- 3

$$v_n(s_1) = -1 + 0.9 \times (-1.9) = -2.71$$

$$v_n(s_2) = -1 + 0.9 \times (8) = 6.2$$

$$v_n(s_3) = -1 + 0.9 \times 10 = 8$$

$$v_n(s_4) = 10$$

T- 4

$$v_n(s_1) = -1 + 0.9 \times \cancel{6.2} = \cancel{5.58}$$

$$v_n(s_2) = -1 + 0.9 \times 8 = \cancel{6.2}$$

$$v_n(s_3) = -1 + 0.9 \times 10 = 8$$

$$v_n(s_4) = 10$$

T- 5

$$v_n(s_1) = -1 + 0.9 \times 6.2 = 5.58$$

$$v_n(s_2) = -1 + 0.9 \times 8 = 6.2$$

$$v_n(s_3) = -1 + 0.9 \times 10 = 8$$

$$v_n(s_4) = 10$$

10
6
Date _____
Page _____

SC

Delta Rule:- (obtained)

$$\Delta v = \gamma (q - v) r_j$$

\uparrow \uparrow \uparrow

Learning rate γ (input value)

discrete

data
page

1) objective

2) formula

3) initial policy π

4) initial value of all states zeroed

$$v(s) = 0$$

→ RL

→ Generalized Policy Iteration (GPI):-

Q.

	0	1	2	3
0	0			+10
1		1/4		-10
2				
3	0			

Terminal States:-

$s(3,3)$

$s(2,3)$

Rewards:-

1) Target state $(3,3) = +10$

2) Terminal State $(2,3) = -10$

3) Rest of the state = -1

Policy is stochastic

Greedy action probability = 0.8

Non-greedy action probability = 0.1

Discount factor = 0.6

State Formula:-

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} P(s', r|s, a) [R + \gamma V_{\pi}(s')]$$

\uparrow
value
function
 \uparrow
probability of
taking an
action given
a state

\uparrow \uparrow \uparrow \uparrow \uparrow
environment
and
state
from next state

∴ If take this option the system it will get high probability result

4) Objective:-

To reach to the goal in minimum makes of steps

\therefore Greedy Action
Following will be the initial policy

	0	1	2	3
0	\rightarrow	\rightarrow	\rightarrow	+10
1	\uparrow	high	\uparrow	-10
2	\uparrow	\rightarrow	\uparrow	\leftarrow
3	*	\uparrow	\uparrow	\uparrow

Let initially value of all states be zero

$$\therefore V(s) = 0$$

~~Environment~~ \because The environment is deterministic

$$\therefore P(s', r | s, a) = 1$$

Iteration :-

lets consider -

$$V(0,1) = 0.8 \times (1 \times (-1 + 0.6 \times 0)) \text{ (Right)} \quad \begin{pmatrix} \text{(greedy)} \\ \text{left} \\ \text{so} \\ \text{at 1} \end{pmatrix} \rightarrow \begin{pmatrix} R_{\text{right}}(\text{greedy}) \\ 10, 1 \\ 0.0 \cdot 8 \end{pmatrix}$$

$$+ 0.1 \times (1 \times (-1 + 0.6 \times 0)) \text{ (left)} \quad \text{No up. down because all}$$

$$= -0.9$$

Discount factor γ = 1

Episode :- (s_1, a_1, s_2) (s_2, a_2, s_3) (s_3, a_3, s_4) (s_4, a_4, s_5) (s_5, a_5, s_6)

Formula:-

$$q_n(s, a) = \delta + \gamma v_n(s')$$

$$\therefore q_n(s, A_1) = 5 + 5 + 10 + 13 + 6 = 79$$

$$q_n(s_2, A_{10}) = 10 + 3 + 6 = 19$$

$$q_n(s_1, a_1) = 5 + (1 \times v_n(s_2))$$

$$= 5 + (1 \times 5)$$

$$v_n(s_1) = 5$$

$$= 5 + 5 = 10$$

$$q_n(s_2, a_2) = 5 + (1 \times v_n(s_3))$$

$$v_n(s_2) = 10$$

$$q_n(s_3, a_3)$$

when episodes are repeated:-
Epis.- (s_1, a_1, s_2) (s_2, a_2, s_3) (s_3, a_3, s_4) (s_4, a_4, s_5) (s_5, a_5, s_6)

Q.F.:-

Epis.- (s_1, a_1, s_2) (s_2, a_2, s_3) (s_3, a_3, s_4) (s_4, a_4, s_5) (s_5, a_5, s_6)

$$\therefore q_n(s_1, a_1)$$

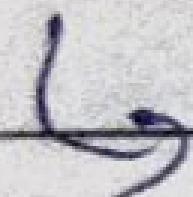
\hookrightarrow It is a hyperparameter which keeps balance between exploration & exploitation.

high \downarrow
explore low \downarrow
exploit

Date _____
Page _____
 \Rightarrow avg \Rightarrow total no. of episodes.

E-Greedy:-

Episode-3 :- $(S_1, A_3, 6) (S_1, a_3, 5) (S_2, a_3, 3) (S_1, A_1, 2) (S_2, A_2, 4)$



First visit-

$$q_{\pi}(S_2, a_1) = 5 + 3 + 2 + 4 \\ = 14$$

mostly used
best in every visit

Every Visit

$$q_{\pi}(S_1, a_1) = (5 + 3 + 2 + 4) + (2 + 4) \\ = 20$$

calculation of every ~~state-action~~ pair that is repeated
for odd biasness to the selection.

Total reward

For (S_2, a_1) :-

$$\begin{aligned} \text{1st visit: } & 5 + 5 + 10 + 3 + 6 = 29 \\ \text{2nd: } & 5 + 6 + 4 + 3 = 18 \\ \text{3rd: } & 5 + 3 + 2 + 4 = 14 \end{aligned}$$

every visit

$$(29 + 18 + 14) / 3 = 20$$

$$\therefore q_{\pi}(S_2, a_1) = 20.33$$

$$\therefore q_{\pi}(S_1, a_1) = 20.33$$

G-Monte-Carlo (ab Aiver Question ($\gamma=1$))

rev.d = -1

$$q_{\pi}(S=(2,3), R_{j+1}) =$$

$\gamma + \gamma C$

$$q(S=(0,0), a=\text{Right})$$

Formula:-

$$q_{\pi}(s, a) = \gamma + \gamma v_{\pi}(s')$$

$$v_{\pi}(s) = \sum \pi(a|s) \cdot q_{\pi}(s, a)$$

Initials:-

biased
deterministic

$$q_{\pi}(S=(0,0), \text{Right}) = \gamma + \gamma(v(S=(0,1)))$$

$$v(S=(0,1)) = 1 \times q_{\pi}(S=(0,1), \text{Right})$$

$$q_{\pi}(S=(0,1), \text{Right}) = \gamma + \gamma(v(S=(0,2)))$$

$$v(S=(0,2)) = 1 \times q_{\pi}(S=(0,2), \text{Down})$$

$$v(S=(0,2)) = -\gamma$$

$$q_{\pi}(S=(1,2), \text{Down}) = \gamma + \gamma(v(S=(2,2)))$$

$$v_{\pi}(S=(1,2)) = 1 \times q_{\pi}(S=(1,2), \text{Right})$$