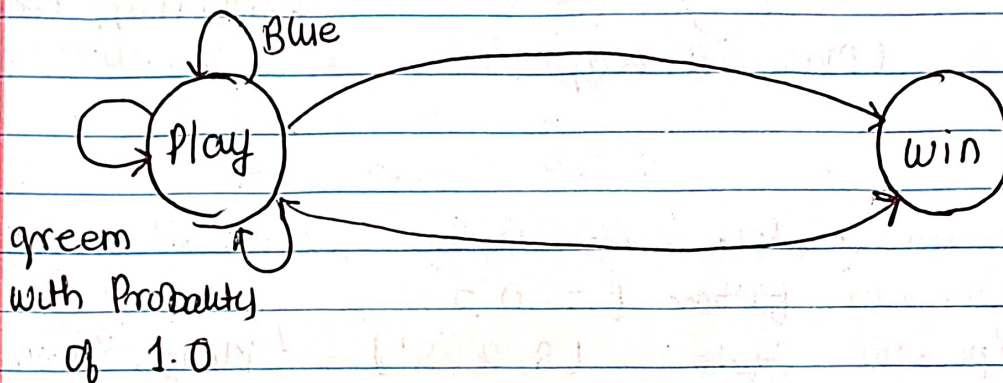NAME : PAURAVI UDAY NAGARKAR.

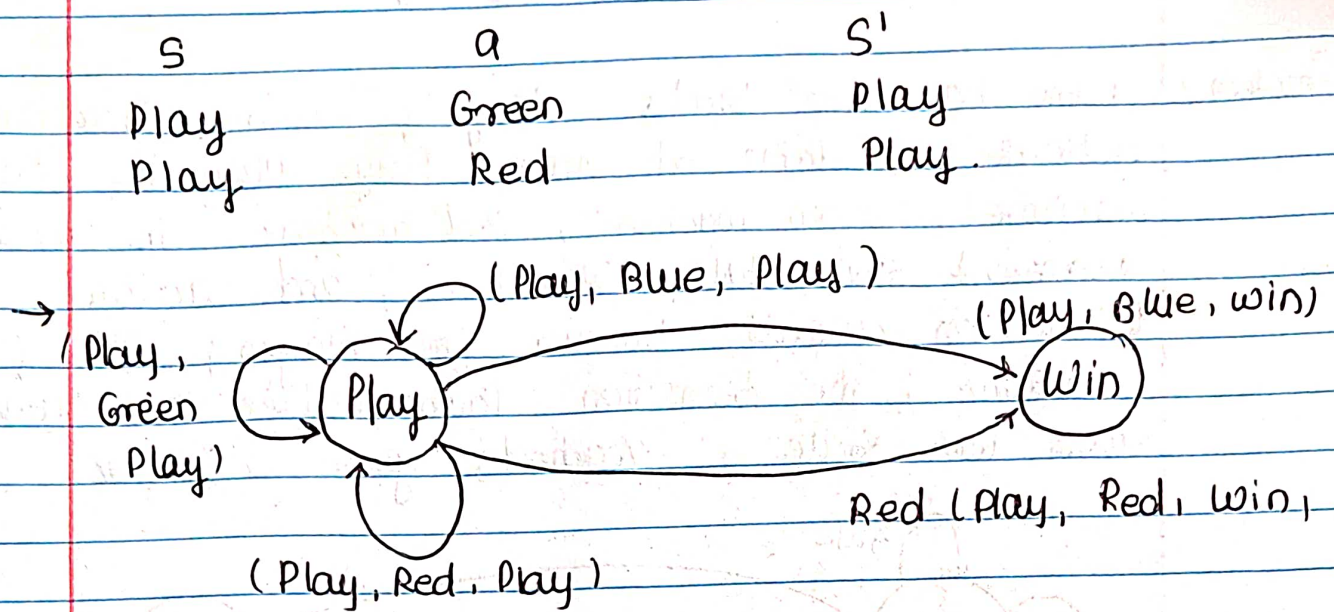JD : WIG50209

## Homework No : 5

**Problem 2** Game has two States Play & Win. There are three actions available at State Play : Play the Blue machine, Green machine, Red machine. The arrows represent the State transitions, and actions are labelled on arrows. however for Playing Blue & Red machine, the transaction Probabilities are unknown. Once wen State is reached, game is over.



greem
with Probality
of 1.0

The reward are given :

| S | a | S' | $R(s, a, s')$ |
|------|-------|------|------|
| Play | green | Play | 2 |
| Play | Blue | Play | 4 |
| Play | Blue | Win | 10 |
| Play | Red | Play | 0 |
| Play | Red | Win | 50 |

Use temporal difference learning to learn the value of Play State, by Successively applying the two Episode below. The initial values of all States from 0. With learning rate $\alpha = 0.5$ and discount factor $\gamma = 0.5$. what value of State Play do we learn?

|  S | a | S' |
|----|-------|-------|
| Play | Green | Play |
| Play | Red | Play. |



(Play, Blue, Play)

(Play, Blue, win)

(Play, Green, Play)

Win

(Play, Red, Play)

Red (Play, Red, win)

→ Learning Rate $\alpha = 0.5$
Discount factor $\gamma = 0.5$
Episodes Given :- $(S, a, S') = ($ Play, green, Play $)$
                                    $($ Play, Red, Play $)$

using temporal difference learning to learn the value of the 'Play state'

updated $V(S) \leftarrow V(S) + \alpha [$ Sample $- V(s)]$ ↙
Sample $= R(S, a, S') + \gamma V(S')$

Initializing the state's value as 0
$V_0(S) = 0$,     for $S =$ Play win.

First Episode :
Sample $= R($ Play, green, Play $) + \gamma V($ Play $)$
       $= 2 + 0.5 * 0 = 2$

updated Value
$V($ Play $) = V($ Play $) + \alpha ($ Sample $- V($ Play $)]$
          $= 0 + 0.5 [2 - 0]$
$V($ Play $) = 1$     after Episode 1.

Second Episode :

$$Sample = R(Play, Red, Play) + \gamma V(Play)$$
$$= 0 + 0.5 * 1 = 0.5$$

updated

$$V(Play) = V(Play) + \alpha[Sample - V(Play)]$$
$$= 1 + 0.5 * (0.5 - 1)$$
$$V(Play) = 0.75 \quad (after\ 2^{nd}\ episode)$$

Value of State Play for Episode 1 : 1
Value of State Play for Episode 2 : 0.75

Problem 2  for same game shown in figure 1, we observed three episodes:

| S | a | S' |
|---|---|----|
| Play | Green | Play |
| Play | Red | Play |
| Play | Blue | win |

use Q-learning to update the values of Q-State by applying the above three episodes one after another. use a learning of 0.5 and a discount of 0.5. Intialize all Q states values as 0

| S | a | Q(S,a) |
|---|---|--------|
| Play | Green | |
| Play | Red | |
| Play | Blue | |

→ learning Rate $\alpha = 0.5$
Discount factor $\gamma = 0.5$
Episodes given - $(S, a, S') \rightarrow$ (Play, Green, Play)
(Play, Red, Play)
(Play, Blue, Win)

using Q learning to update the values of Q-states
Intializing Q states values as '0'
$Q_0(S) = 0$ for S = Play, win.
updated $Q(S, a) \leftarrow Q(S, a) + \alpha \cdot (\text{Sample} - Q(S, a))$
where
Sample = $R(S, a, S') + \gamma \max_a Q'(S', a')$

episode 1:
Sample = $R(\text{Play, green, Play}) + 0.5 * \text{Max}$
$(Q(\text{Play } a'))$

$= 2 + 0 * 0.5$
$= 2$.

updated Play.
$Q(S, a) = 0 + 0.5 (2 - 0)$
$= 1$.
$Q(\text{Play, green}) = 1$

Episode 2
Sample = $R(\text{Play, Red, Play}) + 0.5 * \max_{a'} (Q(\text{Play, } a'))$

Sample = $0 + 0.5 * 1$
$= 0.5$

updated
$Q(S, a) = Q(\text{Play, Red}) + 0.5 * (\text{Sample} - Q(\text{Play, Red}))$
$= 0 + 0.5 * (0.5 - 0)$
$Q(S, a) = 0.25$

$Q$ (Play , Red) = 0.25

Episode 3
  Sample = $R$( Play, Blue , Win) + $\gamma$ Max $Q$ ( win, a')
                                              a'

      = 10+ 0.5 * 0

      = 10

updated
   $Q(S, a)$ = $Q$ (Play, Blue)
         = $Q$( Play, Blue) + $\alpha$ (Sample - $Q$ (Play, Blue)
         = 0 + 0.5 * (10 - 0)
   $Q$( Play , Blue) = 5

| S | a | $Q(S,a)$ |
|------|-------|------|
| Play | green | 1 |
| Play | Red | 0.25 |
| Play | Blue | 5 |