**Reinforcement Learning**
**Assignment 2**
Summary

Q1) Exercise 3.4.

Please find the pdf file submitted for Question - 1.

Q2) The final value function of the equiprobable random policy. Comments added in python notebook.

```
values.round(1).reshape(5,5)
```
```
array([[ 3.3,  8.8,  4.4,  5.3,  1.5],
       [ 1.5,  3. ,  2.3,  1.9,  0.5],
       [ 0.1,  0.7,  0.7,  0.4, -0.4],
       [-1. , -0.4, -0.4, -0.6, -1.2],
       [-1.9, -1.3, -1.2, -1.4, -2. ]])
```

Q4) Optimal solutions to the grid world. Comments added in python notebook.

```
print(np.around(value_star, decimals=1))
```
```
[[22.  24.4 22.  19.4 17.5]
 [19.8 22.  19.8 17.8 16. ]
 [17.8 19.8 17.8 16.  14.4]
 [16.  17.8 16.  14.4 13. ]
 [14.4 16.  14.4 13.  11.7]]
```

```
for x in action_star:
    for y in x:
        print(y, end= ',')
    print()
```
```
['>'],['^' '<' 'd' '>'],['<'],['^' '<' 'd' '>'],['<'],
['^' '>'],['^'],['^' '<'],['<'],['<'],
['^' '>'],['^'],['^' '<'],['^' '<'],['^' '<'],
['^' '>'],['^'],['^' '<'],['^' '<'],['^' '<'],
['^' '>'],['^'],['^' '<'],['^' '<'],['^' '<'],
```
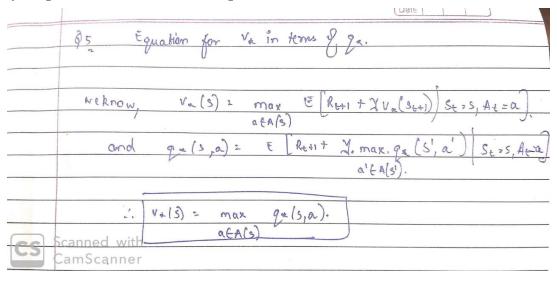
# Q3)

- Exercise 3.15 and 3.16 (Images also added in the 'Images' folder)

## Exercise – 3.15

\# Adding a constant $c$ to all rewards.

$$V_\pi(s) = E_\pi \left[ G_t \mid S_t = S \right]$$

$$= E_\pi \left[ \sum_{k=0}^{\infty} r_{t+1+k} \, \gamma^k \mid S_t = S \right]$$

\# $r' = r + c$

then,

new value of states

$$V_\pi'(s) = E_\pi \left[ \sum_{k=0}^{\infty} (r_{t+1+k} + C) \gamma^k \mid S_t = S \right]$$

$$= E_\pi \left[ \sum_{k=0}^{\infty} r_{t+1+k} \gamma^k \mid S_t = S \right] + E_\pi \left[ \sum_{k=0}^{\infty} C \gamma^k \mid S_t = S \right]$$

$$= V_\pi(s) + \frac{C}{1-\gamma}$$

$$\therefore V_\pi'(s) = V_\pi(s) + \frac{C}{1-\gamma}$$

$$V_c = \frac{C}{1-\gamma} \longrightarrow \text{constant added to all states.}$$

It does not affect relative value of states under any policies.

---

## Ex – 3.16

If the task were an episodic one, then,

$$V_\pi(s) = E_\pi \left[ G_t \mid S_t = S \right]$$

$$= E_\pi \left[ \sum_{k=0}^{T(s)} r_{t+1+k} \, \gamma^k \mid S_t = S \right]$$

$T(s)$ is the number of steps, starting from current state until the terminal state is reached.

$$V_\pi'(s) = E_\pi \left[ \sum_{k=0}^{T(s)} (r_{t+k+1} + C) \, \gamma^k \mid S_t = S \right]$$

$$= E_\pi \left[ \sum_{k=0}^{T(s)} r_{t+1+k} \gamma^k \mid S_t = S \right] + E_\pi \left[ \sum_{k=0}^{T(s)} C \gamma^k \mid S_t = S \right]$$

$$V_\pi'(s) = V_\pi(s) + C \left( \frac{1-(\gamma)^{T(s)}}{1-\gamma} \right)$$

∵ the additional term depends on $T(s)$, the relative value of states might change.

Example → This is because the states that are 'closer' to terminal states will have low $T(s)$ which will increase $(\gamma)^{T(s)}$ which will decrease the additional term.

$V_\pi'(s)$ will have small value of additional term for states that are closer to the terminal states.

∴ it might change relative valuation of states.

Q5) Equation for v* in terms of q*



Q5. Equation for $V_*$ in terms of $q_*$.

We know,

$$V_*(s) = \max_{a \in A(s)} \mathbb{E}\left[R_{t+1} + \gamma V_*(s_{t+1}) \mid S_t = s, A_t = a\right].$$

and

$$q_*(s,a) = \mathbb{E}\left[R_{t+1} + \gamma \max_{a' \in A(s')} q_*(s', a') \mid S_t = s, A_t = a\right]$$

$$\therefore \boxed{V_*(s) = \max_{a \in A(s)} q_*(s,a).}$$

Scanned with CamScanner

Q6)
- policy iteration and value iteration (VI) to solve the Gridworld in Example 4.1
- the fix to the bug mentioned in Exercise 4.4. (in jupyter notebook)

Q7) Exercise 4.7

Original Example                                Exercise - 4.7