

# Reinforcement Learning: Grid world

Charuka Amarappulige | Chinonso Imochukwu

## Part 1

In this section we take 5 x 5 grid world with specific characteristics to compare different reinforcement learning methodologies. Property of grid world environment as below. Agent in environment can take a step up, down, left or right actions. However, If the agent attempts to step off the grid, the location of the agent remains unchanged with reward of  $-0.5$ .

The blue, green, red and yellow squares represent special states. At the blue square, any action yields a reward of 5 and causes the agent to jump to the red square. At the green square, any action yields a reward of 2.5 and causes the agent to jump to either the yellow square or the red square with probability 0.5. Taking any action in white, red and yellow squares will not yield any reward unless it is not trying to step off the grid.

For location identifying purpose each state is numbered as below figure in most of the cases unless different requirement arises.

### 5x5 Grid World

(0,0)	(0,1)	(0,2)	(0,3)	(0,4)
(1,0)	(1,1)	(1,2)	(1,3)	(1,4)
(2,0)	(2,1)	(2,2)	(2,3)	(2,4)
(3,0)	(3,1)	(3,2)	(3,3)	(3,4)
(4,0)	(4,1)	(4,2)	(4,3)	(4,4)

## Question 1

### 1. Solving the system of Bellman equations explicitly

In order to solve bellman equation fixed policy  $\pi$  and solve it explicitly, we have broken down original Bellman equation.

$$V(s) = \sum_{s' \in S} P(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma V(s')]$$

Writing the equation in matrix from

$$V = R + \gamma PV$$

where R and P are expected reward matrix and transition probability matrix according to given policy.

$$R(s) = \sum_{s' \in S} P(s'|s, \pi(s)) R(s, \pi(s), s')$$

$$P(s, s') = P(s, \pi(s), s')$$

By calculating the R and P we get.

Transition probability from one state another  $P = \begin{bmatrix} 0.5 & 0.25 & \cdots & 0 \\ 0.25 & \ddots & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0.5 \end{bmatrix}_{25 \times 25}$

Expected reward for each state  $R = \begin{bmatrix} -0.25 \\ 5 \\ \vdots \\ -0.25 \end{bmatrix}_{25 \times 1}$

Rearranging the equation

$$V = R + \gamma PV \Rightarrow (I - \gamma P)V = R \Rightarrow V = (I - \gamma P)^{-1}R$$

### Functions used

*get\_transition\_rewards*

Objective of this function is to calculate transition probabilities, rewards, expected value and individual values for each action once it is once it was given the state, policy and given value function.

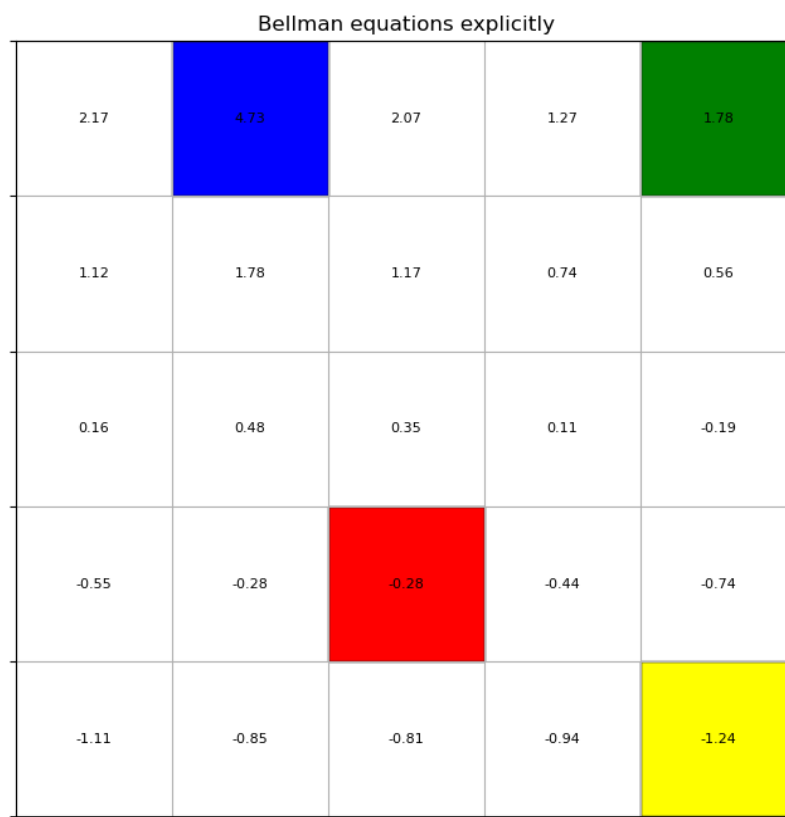
*plot\_gridworld*

This function is used to plot the gridworld and its values accordingly once finalized value given.

*SciPy.linalg*

Module is used to calculate matrix inverse and multiplication.

Once value function is solved finalized values ( $V_\pi$ ) rounded up to two decimal places is looks like in below figure. Highest valued state came as blue state while rest of the values distributed toward rest of the area.



## 2. Iterative policy evaluation

In order to calculate ( $V_\pi$ ) using iterative policy evaluation we have used the below algorithm extracted from Chapter 4.1 of Reinforcement Learning book by Richard S. Sutton and Andrew G. Barto. The same 'get\_transition\_rewards' function being used to calculate expected values in each state

### Iterative Policy Evaluation, for estimating $V \approx v_\pi$

Input  $\pi$ , the policy to be evaluated

Algorithm parameter: a small threshold  $\theta > 0$  determining accuracy of estimation

Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop:

$\Delta \leftarrow 0$

Loop for each  $s \in \mathcal{S}$ :

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until  $\Delta < \theta$

The output ( $V_\pi$ ) came as below figure which is identical to the output in pervious case.

Iterative policy evaluation				
2.17	4.73	2.07	1.27	1.78
1.12	1.78	1.17	0.74	0.56
0.16	0.48	0.35	0.11	-0.19
-0.55	-0.28	-0.28	-0.44	-0.74
-1.11	-0.85	-0.81	-0.94	-1.24

### 3. Value iteration

In this section we have tried to calculate value function using value iteration methodology similarly we have used algorithm extracted by the Reinforcement Learning book chapter 4.4. However it was obvious that value iteration method is not suitable to evaluate a policy. It is much suitable to calculate optimum value function.

### Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold  $\theta > 0$  determining accuracy of estimation  
Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop:

```
|  $\Delta \leftarrow 0$   
| Loop for each  $s \in \mathcal{S}$ :  
|    $v \leftarrow V(s)$   
|    $V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$   
|    $\Delta \leftarrow \max(\Delta, |v - V(s)|)$   
until  $\Delta < \theta$ 
```

Output a deterministic policy,  $\pi \approx \pi_*$ , such that  
 $\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$

The results for value iteration function can be illustrate as below figure

Value iteration				
21.0	22.1	21.0	19.95	18.38
19.95	21.0	19.95	18.95	18.0
18.95	19.95	18.95	18.0	17.1
18.0	18.95	18.0	17.1	16.25
17.1	18.0	17.1	16.25	15.43

From all the three algorithm it was very clear that highest valued state is state (0,1) which is blue demarcated state. It was not a surprise that since any action taken from blue state incur reward of 5 to the agent. Due to this reason, it is quite obvious the highest valued state is the blue state.

You can recreate the above results from referring to the file A2P1Q1.ipynp and it is available in the below git hub repository

Link: [https://github.com/CharukaP/RL\\_Gridworld](https://github.com/CharukaP/RL_Gridworld)

## Question 2

### 1. Optimum policy by explicitly solving the Bellman optimality equation

This segment explain how optimum policy is obtained by explicitly solving bellman equation. Two methods being used to check the possibility of solving this problem. Two new function being used in order to support this evaluation. In method 1 we have carried out the solution in four steps.

#### Functions used

*get\_next\_state\_and\_reward*

This function work on calculating the reward and the next state according to the environment characteristics once it was given the current state and the action.

*action\_search*

This function will return concatenated actions string if more than one action is optimum.

#### Method 1

**Step 1:** Adjustment to Bellman optimality equation

The Bellman equation for a state  $s$  can be written as:

$$V(s) = \max_a (R(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s'))$$

Adjusted Bellman equation.

$$V(s) - \max_a (R(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s')) = 0$$

**Step 2:** Genarate 25 equation for each 25 states.

Using the adjusted Bellman equation generate 25 equations.

**Step 3:** Solve the set of equations

Solving the system of equations using the 'fsolve' function from SciPy will give you the  $V^*$

#### Step 4: Solve the set of equations

Use calculated  $V^*$  and perform one step look ahead action to identify optimum action  $\pi^*$

Once calculated the optimum actions can be illustrate as below figure. Optimum policy is leading the agent towards the blue state obviously.

Optimum policy solving Bellman optimality equations explicitly

right	up	left	right	up
up	up	up	up	up
up	up	up	up	up
up	up	left	right	up
up	up	left	up	up

#### Method 2

In method 2 we used the same matrix solution we used for part1 question 1 section with slight change to reward matrix here we used reward as maximum reward instead of expected reward. As below

$$R(s) = \max_a R(s, a)$$

Solving  $V = (I - \gamma P)^{-1} R$  equation with SciPy.linalg module to calculate matrix inverse and multiplication. We were able to achieve the below results. Code for the second methodology is uploaded in git hub link under special file name A2P1Q2\_part1\_optional.ipynb

Optimum policy solving Bellman optimality equations explicitly

right	up	left	left	down
up	up	up	up	up
up	up	up	up	up
up	up	up	up	up
up	up	up	up	up

## 2. Policy iteration with iterative policy evaluation

Calculation of optimum policy using iterative policy evaluation. Was carried out according to algorithm extracted by Chapter 4.3 in Reinforcement Learning book where the algorithm can be formulated as below figure.



### Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization  
 $V(s) \in \mathbb{R}$  and  $\pi(s) \in \mathcal{A}(s)$  arbitrarily for all  $s \in \mathcal{S}$
2. Policy Evaluation  
 Loop:  
 $\Delta \leftarrow 0$   
 Loop for each  $s \in \mathcal{S}$ :  
 $v \leftarrow V(s)$   
 $V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$   
 $\Delta \leftarrow \max(\Delta, |v - V(s)|)$   
 until  $\Delta < \theta$  (a small positive number determining the accuracy of estimation)
3. Policy Improvement  
 $policy-stable \leftarrow true$   
 For each  $s \in \mathcal{S}$ :  
 $old-action \leftarrow \pi(s)$   
 $\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$   
 If  $old-action \neq \pi(s)$ , then  $policy-stable \leftarrow false$   
 If  $policy-stable$ , then stop and return  $V \approx v_*$  and  $\pi \approx \pi_*$ ; else go to 2

Once  $V^*$  being calculated we used multiple action searching function to identify action with same values. The results are illustrated in below figure.

Policy iteration with iterative policy evaluation

Right	up Down Left Right	Left	Left	up Down Left Right
up Right	up	up Left	up Left	up Left
up Right	up	up Left	up Left	up Left
up Right	up	up Left	up Left	up Left
up Right	up	up Left	up Left	up Left

### 3. Policy improvement with value iteration

Identify the optimum policy via value iteration just extension of last sub question of Question 1. Since we have already calculated the  $V^*$  values it's just a matter of carrying out one step look had and get the optimal policy. We have used the extension of algorithm extracted from same 4.4 chapter in RL book.

#### Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold  $\theta > 0$  determining accuracy of estimation  
Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop:

```
|  $\Delta \leftarrow 0$ 
| Loop for each  $s \in \mathcal{S}$ :
|    $v \leftarrow V(s)$ 
|    $V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$ 
|    $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
until  $\Delta < \theta$ 
```

Output a deterministic policy,  $\pi \approx \pi_*$ , such that

$$\pi(s) = \arg \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$$

The results show the similar to the second part where  $\pi^*$  leads all the action towards the blue state.

Policy improvement with value iteration

Right	up Down Left Right	Left	Left	up Down Left Right
up Right	up	up Left	up Left	up Left
up Right	up	up Left	up Left	up Left
up Right	up	up Left	up Left	up Left
up Right	up	up Left	up Left	up Left

You can recreate the above results from referring to the file A2P1Q2.ipynp and it is available in the below git hub repository

Link: [https://github.com/CharukaP/RL\\_Gridworld](https://github.com/CharukaP/RL_Gridworld)

## Part 2

There have been slight changes into environment is introduced in this section. Two terminal states were introduced in black where episode terminate if agent visits that cell. Previous general movement incurred zero reward has changed to -0.2 compared to part 1. Same state identifying notations were carried out in here as well which is illustrated in below figure.

### 5x5 Grid World

(0,0)	(0,1)	(0,2)	(0,3)	(0,4)
(1,0)	(1,1)	(1,2)	(1,3)	(1,4)
(2,0)	(2,1)	(2,2)	(2,3)	(2,4)
(3,0)	(3,1)	(3,2)	(3,3)	(3,4)
(4,0)	(4,1)	(4,2)	(4,3)	(4,4)

### Question 1

Apart from 'get\_next\_state\_and\_reward' function and 'plot\_gridworld' functions there are new supporting functions have been used to fulfill the task of finding optimum policy.

#### Functions used

Random\_start

This function is used to get a random starting state except the terminal states

generate\_episode

This function is used to generate series of steps according to equiprobable actions till agent reaches to terminal point and finally return the episode to main function.

choose\_action

soft action selection was carried out from this function where it will receive state and policy, and it returns the action chosen

generate\_episode\_soft

This function used to generate series of steps till agent reach to terminal point while taking action from choose action function with soft policy.

## 1. Monte Carlo method with exploring starts

Usage of Monte Carlo method to identify optimal policy is explained in chapter 5.3 in Reinforcement Learning book and I have extracted the algorithm for it as in below figure.

### Monte Carlo ES (Exploring Starts), for estimating $\pi \approx \pi_*$

Initialize:

$\pi(s) \in \mathcal{A}(s)$  (arbitrarily), for all  $s \in \mathcal{S}$

$Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$  empty list, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose  $S_0 \in \mathcal{S}, A_0 \in \mathcal{A}(S_0)$  randomly such that all pairs have probability  $> 0$

Generate an episode from  $S_0, A_0$ , following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :

Append  $G$  to  $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$

Using  $\gamma=0.95$  and running the algorithm will give you the optimum policy as below grid.

Monte Carlo ES (Exploring Starts), for estimating optimum policy

right	down	left	right	left
up	up	up	up	up
up	up	up	right	up
down	up	up	up	up
up	left	left	up	up

## 2. Without exploring starts but the $\epsilon$ -soft approach

We have used algorithm from chapter 5.4 as illustrated in below figure.

### On-policy first-visit MC control (for $\epsilon$ -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small  $\epsilon > 0$

Initialize:

$\pi \leftarrow$  an arbitrary  $\epsilon$ -soft policy

$Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$  empty list, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :

Append  $G$  to  $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken arbitrarily)

For all  $a \in \mathcal{A}(S_t)$ :

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \epsilon + \epsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \epsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

In order to determine appropriate epsilon value multiple exercises were carried out with values 0.1, 0.5, 0.8 and 0.9 for 1000 episodes.

Epsilon = 0.1

Monte Carlo control epsilon-soft approach (epsilon = 0.1)

right	left	left	down	down
up	up	left	left	left
up	up	left	up	up
up	up	up	up	left
up	up	left	up	left

Epsilon = 0.5

Monte Carlo control epsilon-soft approach (epsilon = 0.5)

right	up	left	left	right
up	up	left	up	up
up	up	up	up	up
up	up	up	left	left
up	up	up	up	left

Epsilon = 0.8

Monte Carlo control epsilon-soft approach (epsilon = 0.8)

right	left	left	left	left
up	up	up	left	up
up	up	up	up	up
up	up	up	up	up
up	left	up	up	up

Epsilon = 0.9

Monte Carlo ES (Exploring Starts), for estimating optimum policy

right	down	left	right	up
up	up	up	up	up
up	up	up	right	down
down	left	up	up	up
up	left	left	up	up

From the multiple runs it was observed that Epsilon value closer to 1 (0.9) give the most accurate results compared to other evaluations.

You can recreate the above results from referring to the file A2P2Q1.ipynp and it is available in the below git hub repository

Link: [https://github.com/CharukaP/RL\\_Gridworld](https://github.com/CharukaP/RL_Gridworld)



## Question 2

This question was addressed using off policy Monte Carlo control methodology extracted from chapter 5.7 in Reinforcement learning book. Algorithm steps are mentioned in the below figure.

Due to constructed methodology in algorithm the learning rate is slow, and we need to use higher episode count to get near optimal policy. Running episodes for below 100,000 will not give an optimum policy. Furthermore, arbitrary starting value for Q table is taken as -3 after comparison with multiple experiments.

You can recreate the above results from referring to the file A2P2Q2.ipynp and it is available in the below git hub repository

Link: [https://github.com/CharukaP/RL\\_Gridworld](https://github.com/CharukaP/RL_Gridworld)

### Off-policy MC control, for estimating $\pi \approx \pi_*$

```
Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :  
   $Q(s, a) \in \mathbb{R}$  (arbitrarily)  
   $C(s, a) \leftarrow 0$   
   $\pi(s) \leftarrow \operatorname{argmax}_a Q(s, a)$  (with ties broken consistently)  
  
Loop forever (for each episode):  
   $b \leftarrow$  any soft policy  
  Generate an episode using  $b$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$   
   $G \leftarrow 0$   
   $W \leftarrow 1$   
  Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :  
     $G \leftarrow \gamma G + R_{t+1}$   
     $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$   
     $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$   
     $\pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$  (with ties broken consistently)  
    If  $A_t \neq \pi(S_t)$  then exit inner Loop (proceed to next episode)  
     $W \leftarrow W \frac{1}{b(A_t|S_t)}$ 
```

The results obtained by running algorithm for 10,000,000 episodes will give a policy similar to below grid. Due to high number of iteration required for get an accurate result this algorithm is not suitable for quick evaluations.

Behaviour policy with equiprobable moves

right	up	left	left	up
up	up	up	left	up
right	up	up	up	up
down	up	up	up	up
up	left	up	up	up

### Question 3

In this segment new characteristic performance was introduced to gridworlds where green and blue states interchange with each other with probability of 0.1 every time special function is called. It can be either steps or episode. We have introduced new function in order to facilitate this requirement. We have use two methods to evaluate policy

1. Policy iteration with iterative policy evaluation
2. Monte Carlo method with exploring starts

#### Functions used

`permute_location`

The primary objective of this function is to get current blue and green stats as variables and swap it. Once swap completed rewrite the characteristics of those special states.

#### Method 1

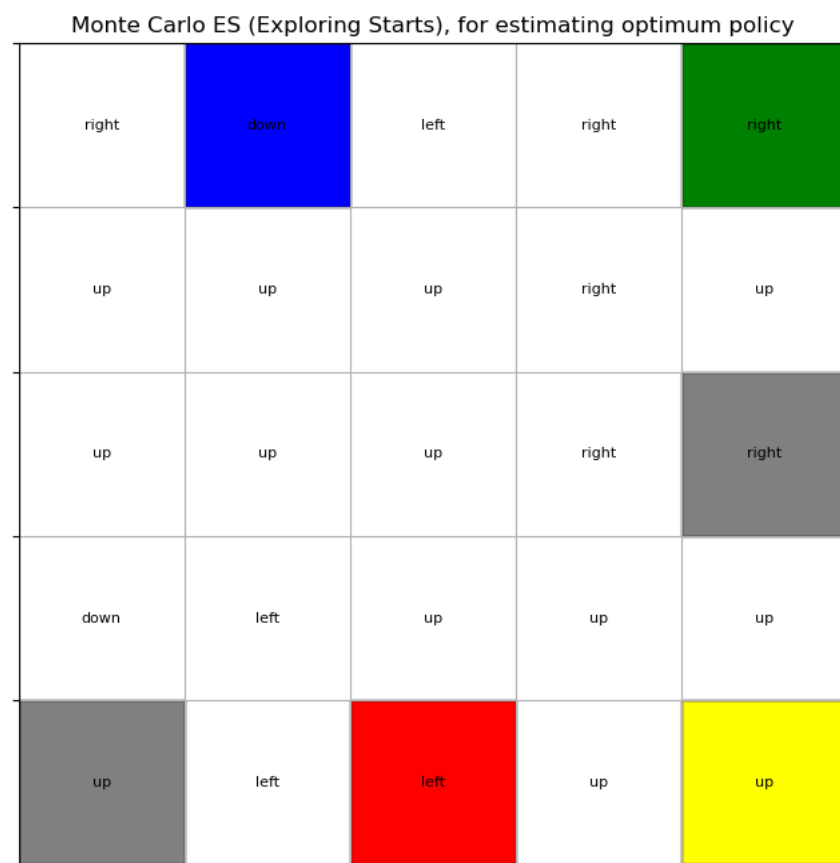
The requirement of the sub section is to determine suitable policy for this dynamic environment via policy iteration as we explained earlier policy iteration methodology, we discussed in part 1 question 2 sub section 2. The specially adjusted algorithm for the evaluation is



## Method 2

Here we have used Monte Carlo method with exploring start algorithm used in Part 2 Question 1 sub section 1 with adjustment to episode generation function where in each step of the episode permutation function is called. By running the algorithm for 10,000 episodes gives you some accurate representation of the optimum policy. According to results it was observed that agent tried to move towards either blue and green states whenever they are close by. And also when the agent is closer to terminal points instead of going to blue and green states agent move to terminal state where reward is 0. Since taking further steps incur -0.2 reward it is ideal to get 0 rather than -0.2.

From the both methodologies suggested it can be observed that methodology 2 is more suitable for this kind of dynamic environments. Policy obtained by method 2 is illustrated in below figure.



You can recreate the above results from referring to the file A2P2Q3.ipynp and it is available in the below git hub repository

Link: [https://github.com/CharukaP/RL\\_Gridworld](https://github.com/CharukaP/RL_Gridworld)