

**NANYANG  
TECHNOLOGICAL  
UNIVERSITY  
SINGAPORE**

Masters in science (Artificial Intelligence)  
(MSAI)

AI-610

Assignment 1

Name of Student:  
Teo Lim Fong

Matriculation No: G2101964G

## Contents

Q learning.....	4
Plotting of Episode rewards vs. Episodes with learning rate of 0.5 .....	4
Plotting of Episode rewards vs. Episodes with learning rate of 0.9 .....	5
Value Table.....	6
Policy for Q learning.....	8
SARSA.....	9
Plotting of Episode rewards vs. Episodes with learning rate of 0.5 .....	9
Plotting of Episode rewards vs. Episodes with learning rate of 0.32 .....	10
Value Table .....	11
Policy for SARSA.....	12
Experiment with Epsilon Decay .....	13
Experiment with Learning Rate Decay.....	14
Experiment with Epsilon Decay & Learning rate Decay.....	14
Discussion .....	15
Evaluation .....	16
Extra (Q Learning without Box coordinate) .....	16
Plotting of Episode rewards vs. Episodes with learning rate of 0.5 .....	16
Plotting of Episode rewards vs. Episodes with learning rate of 0.9 .....	17
V Table.....	18
Policy and Value Table for Q learning without Box coordinate .....	19
ANNEX A (Q learning policy step by step) .....	20
Agent (5,0) move → to (5,1) then move ↑ to Box (4,1) .....	20
Agent (4,1) move ↑ to Box (3,1).....	20
Agent (3,1) move ↑ to Box (2,1).....	20
Agent (2,1) move ← to (2,0) then ↑ to (1,0) then → to Box (1,1) .....	21
Agent (1,1) move → to Box (1,2).....	21
Agent (1,2) move → to Box (1,3).....	21
Agent (1,3) move → to Box (1,4).....	22
Agent (1,4) move ↑ to (0,4) then → to (0,5) then move ↓ to Box (1,5).....	22
Agent (1,5) move ↓ to Box (2,5).....	22
Agent (2,5) move ↓ to Box (3,5).....	23
Agent (3,5) move ← to (3,4) the move ↓ to (4,4) then move → to Box (4,5).....	23
Agent (4,5) move → to Box (4,6).....	23
Agent (4,6) move → to Box (4,7).....	24
Agent (4,7) move ↓ to (5,7) then move → to (5,8) then move ↑ to Box (4,8).....	24
Agent (4,8) move ← to (4,7) then ↑ to (3,7) then → to Box (3, 8) .....	24

Agent (3,8) move → to Box (3,9).....	25
Agent (3,9) move → to Box (3,10).....	25
Agent (3,10) move ↓ to (4,10) then move → to (4,11) then move ↑ to Box (3,11).....	25
Agent (3,11) move ↑ to Box (2,11).....	26
Agent (2,11) move ← to (2,10) then ↑ to (1,10) then → to Box (1,11) .....	26
Agent (1,11) move → to Box (1,12).....	26
Agent (1,12) move ↑ to (0,12) then → to (0,13) then ↓ to Box (1,13) .....	27
Agent (1,13) move ↓ to Box (2,13).....	27
Agent (2,13) move ↓ to Box (3,13).....	27
<b>ANNEX B SARSA (Step-by-step) .....</b>	<b>28</b>
Agent (5,0) move → to (5,1) then move ↑ to Box (4,1) .....	28
Agent (4,1) move ↑ to Box (3,1).....	28
Agent (3,1) move ↑ to Box (2,1).....	29
Agent (2,1) move ← to (2,0) then ↑ to (1,0) then move → to Box (1,1).....	29
Agent (1,1) move → to Box (1,2).....	29
Agent (1,2) move → to Box (1,3).....	30
Agent (1,3) move → to Box (1,4).....	30
Agent (1,4) move ↑ to (0,4) then move → to (0,5) then move ↓ to Box (1,5).....	30
Agent (1,5) move ↓ to Box (2,5).....	31
Agent (2,5) move ↓ to Box (3,5).....	31
Agent (3,5) move ← to (3,4) then move ↓ to (4,4) then move → to Box (4,5).....	31
Agent (4,5) move → to Box (4,6).....	32
Agent (4,6) move → to Box (4,7).....	32
Agent (4,7) move ↓ to (5,7) then move → to (5,8) the move ↑ to Box (4,8).....	32
Agent (4,8) move ← to (4,7) then move ↑ to (3,7) then move → to Box (3,8).....	33
.....	33
Agent (3,8) move → to (3,9) .....	33
Agent (3,9) move → to (3,10) .....	33
Agent (3,10) move ↓ to (4,7) then move → to (4,11) then move ↑ to Box (3,11).....	34
Agent (3,11) move ↑ to Box (2,11).....	34
Agent (2,11) move ← to (2,10) then move ↑ to (1,10) then move → to Box (1,11).....	34
Agent (1,11) move ↑ to Box (1,12).....	35
Agent (1,12) move ↑ to (0,12) then move → to (0,13) then move ↓ to Box (1,13).....	35
Agent (1,13) move ↓ to Box (2,13).....	35
Agent (2,13) move ↓ to Box (3,13).....	36

# Q learning

Q learning is an off policy that updates its Q values by using the Q values of the next state and select the maximum action in that state. In off policy, the target policy is not equal to the behaviour policy. It evaluates and improve a different policy from that policy that is used for the action.

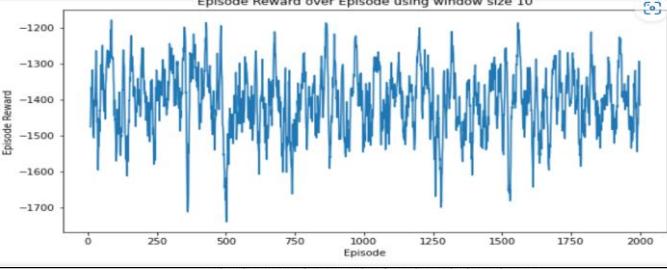
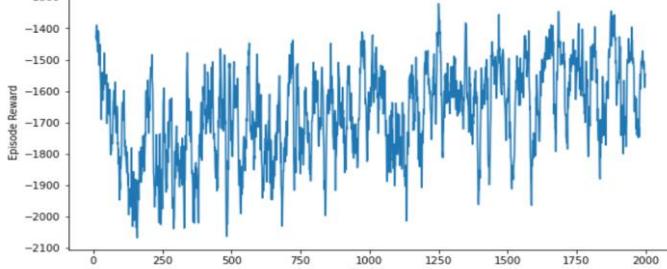
Nevertheless, we need to finetune the parameters such as the discounted value, the learning rate, and the epsilon. Different environments have different values of these parameters. For discounted value, gamma, it is normally the range of 0.9 to 0.99. For this assignment, I will assign the discounted value as 0.98 and experiment the other two parameters – learning rate and the epsilon.

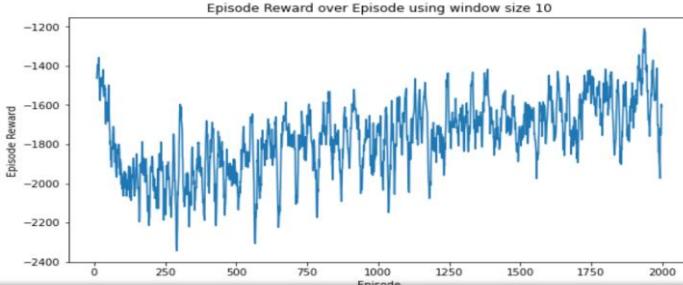
Learning rate can be defined as how fast you want the function to learn the slope. In deep learning, this learning rate is normally set to 0.01 – 0.001 because if the learning rate is too high, this mean it will cause overfitting and if the learning rate is too low it will cause underfitting. However, Q learning does not involve any deep learning. Therefore, we can try to set a range of value such as 0.9 and 0.5.

Epsilon is important in reinforcement learning. It is a trade-off between exploration and exploitation. Ideally, we need to explore the environment at the beginning which normally the epsilon will be high to choose random action then after a few hundred episodes the epsilon will decay and exploited by choosing the maximum action in that state. However, I want to experiment a static epsilon instead of decay epsilon to see how different epsilon performed. I will be using 0.5, 0.1 and 0.01.

## Plotting of Episode rewards vs. Episodes with learning rate of 0.5

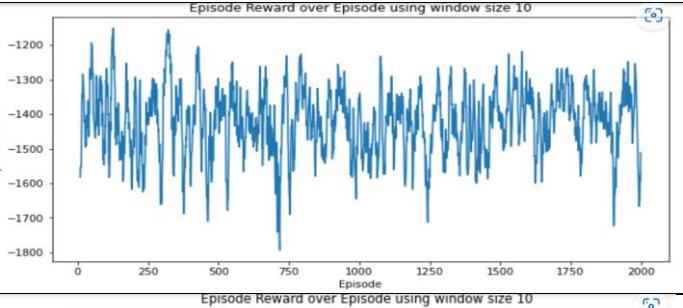
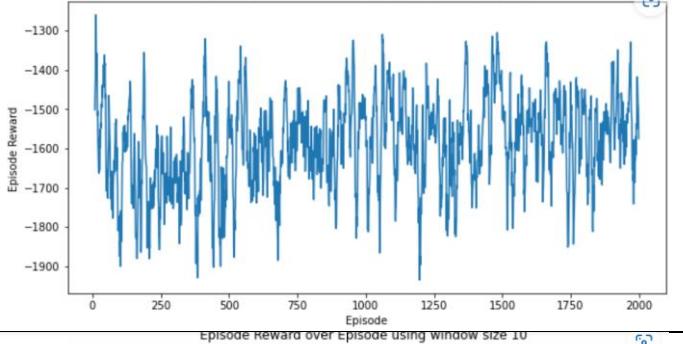
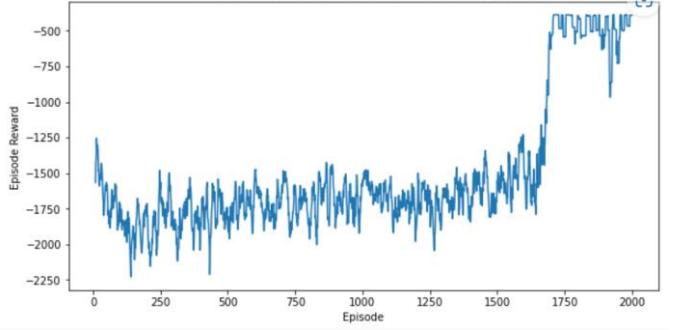
Table 1: Plotting Episode Reward over Episodes with learning rate of 0.5 and different epsilon values

	Epsilon	Episode Reward over Episodes	Push Box to Goal	Total Reward
Learning rate = 0.5	0.5		No	NIL
	0.1		No	NIL

	0.01		Yes	~ -1200
--	------	--	-----	---------

### Plotting of Episode rewards vs. Episodes with learning rate of 0.9

Table 2: Plotting Episode Reward over Episodes with learning rate of 0.5 and different epsilon values

	Epsilon	Episode Reward over Episodes	Push Box to Goal	Total Reward
Learning rate = 0.9	0.5		No	NIL
	0.1		No	NIL
	0.01		Yes	-388

From both Table1 and 2, we can clearly see that epsilon 0.01 are able to make the agent to push the box to the desire state (Goal). However, I choose the Table 2 epsilon 0.01 (highlighted in red) because from the plot we can analysis that in Table 1 the highest total reward of epsilon 0.01 is around -1200 compared to the epsilon 0.01 in Table 2 is -388 to be exact. If the total reward is low but still managed to reach the destination, this means two things. 1) The solution is still not coverage and can be improved or 2) It has many repeated routes that cause the total reward to be low.

## Value Table

The ideal path generated using Q-learning are shown below:

```
[4], [1], [1], [1], [3], [1], [4], [4], [4], [1], [4], [2], [2], [2],
[3], [2], [4], [4], [4], [2], [4], [1], [3], [1], [4], [4], [4], [2], [4],
[1], [1], [3], [1], [4], [1], [4], [2], [2], [2]]
```

Figure 1: Policy for Q learning

The optimal action is highlighted in red for each state.

Table 3: Value Table for Q learning

State	Actions				Best Action
	UP	Down	Left	Right	
((5,0), (4,1))	-294.956	-300.025	-298.063	<b>-293.473</b>	Right
((5,1), (4,1))	<b>-285.176</b>	-291.020	-298.886	-293.343	Up
((4,1), (3,1))	<b>-275.691</b>	-286.992	-278.913	-283.235	Up
((3,1), (2,1))	<b>-264.9905</b>	-270.643	-268.347	-274.807	Up
((2,1), (1,1))	-259.448	-264.598	<b>-253.051</b>	-253.628	Left
((2,0), (1,1))	<b>-239.84</b>	-257.511	-248.422	-258.853	Up
((1,0), (1,1))	-245.253	-251.098	-239.715	<b>-227.396</b>	Right
((1,1), (1,2))	-217.267	-233.183	-230.825	<b>-215.710</b>	Right
((1,2), (1,3))	-212.516	-211.737	-208.240	<b>-204.808</b>	Right
((1,3), (1,4))	-195.845	-1013.018	-208.518	<b>-194.700</b>	Right
((1,4), (1,5))	<b>-185.408</b>	-197.956	-191.669	-885.5	Up
((0,4), (1,5))	-185.340	-192.723	-182.698	<b>-174.907</b>	Right
((0,5), (1,5))	-174.634	<b>-165.211</b>	-177.613	-950.625	Down
((1,5), (2,5))	-165.358	<b>-156.338</b>	-165.441	-981.343	Down
((2,5), (3,5))	-157.951	<b>-148.304</b>	-151.924	-1011.012	Down
((3,5), (4,5))	-146.525	-145.437	<b>-141.126</b>	-1003.102	Left
((3,4), (4,5))	-138.946	<b>-132.782</b>	-948.75	-144.794	Down
((4,4), (4,5))	-138.246	-128.892	-995.203	<b>-125.288</b>	Right
((4,5), (4,6))	-122.729	-126.403	-122.038	<b>-118.661</b>	Right
((4,6), (4,7))	-1005.058	-113.016	-117.016	<b>-112.919</b>	Right
((4,7), (4,8))	-114.787	<b>-108.0813</b>	-112.605	-880.25	Down
((5,7), (4,8))	-108.147	-106.849	-109.003	<b>-102.128</b>	Right
((5,8), (4,8))	<b>-97.065</b>	-102.072	-107.167	-945	Up
((4,8), (3,8))	-91.781	-97.590	<b>-90.882</b>	-882.875	Left
((4,7), (3,8))	<b>-83.554</b>	-93.735	-94.031	-94.470	Up
((3,7), (3,8))	-78.507	-79.417	-882.875	<b>-77.095</b>	Right
((3,8), (3,9))	-75.762	-74.8	-80.300	<b>-71.526</b>	Right
((3,9), (3,10))	-975.531	-755.25	-67.758	<b>-66.864</b>	Right
((3,10), (3,11))	-63.971	<b>-63.121</b>	-64.325	-988.313	Down
((4,10), (3,11))	-63.688	-64.170	-755.25	<b>-58.291</b>	Right
((4,11), (3,11))	<b>-54.379</b>	-61.689	-58.825	-880.25	Up
((3,11), (2,11))	<b>-49.366</b>	-55.762	-52.251	-944.063	Up
((2,11), (1,11))	-47.652	-44.769	<b>-43.231</b>	-882	Left
((2,10), (1,11))	<b>-35.950</b>	-41.396	-756.75	-39.827	Up
((1,10), (1,11))	-30.908	-31.814	-1000.125	<b>-29.541</b>	Right
((1,11), (1,12))	-24.806	-27.657	-27.751	<b>-24.022</b>	Right
((1,12), (1,13))	<b>-19.401</b>	-880.25	-20.654	-23.421	Up
((0,12), (1,13))	-17.056	-15.702	-23.721	<b>-13.683</b>	Right
((0,13), (1,13))	-9.851	<b>-8.8608</b>	-18.777	-11.767	Down
((1,13), (2,13))	-5.913	<b>-4.96</b>	-6.158	-5.940	Down
((2,13), (3,13))	-3	<b>-2</b>	-753	-2.985	Down
((3,13), (4,13))	<b>GOAL</b>				

The figures below demonstrate a step-by-step agent location to the box location. The full step-by-step policy will be placed under Annex A. The state contains two different coordinates, the agent coordinates and box coordinates. The agent needs to reach the box location to update the next location for the box. For example, in Figure 2, the agent location is (5,0) and the Box location is (4,1). Figure 2 show the correct direction to move from agent (5,0) to Box (4,1) by first moving → to (5,1) then move ↑ to reach Box (4,1). I also added different colours for each action for better visualization.

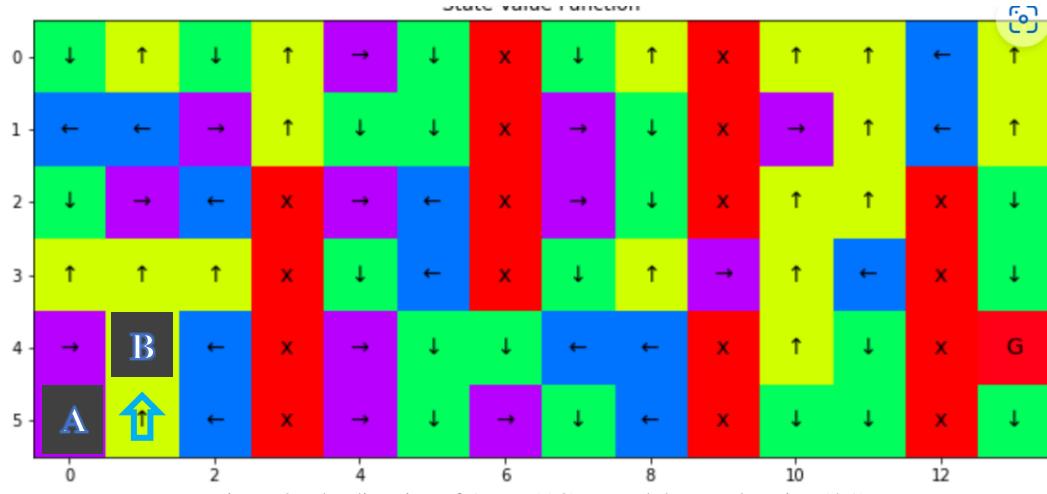


Figure 2: The direction of Agent (5,0) toward the Box location (4,1)

Next the agent (4,1) will move ↑ to Box (3,1).

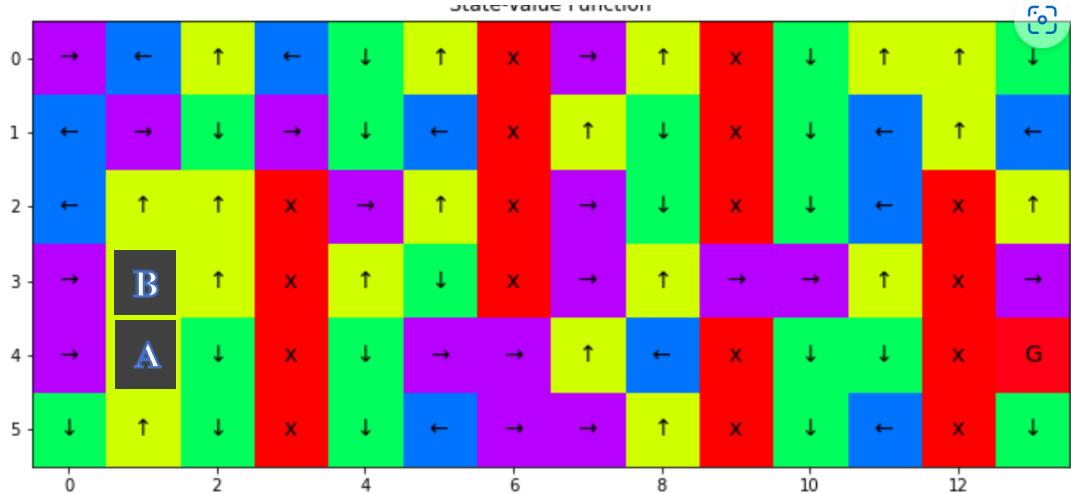


Figure 3: The direction of Agent (4,1) toward the Box location (3,1)

## Policy for Q learning

This is the ideal policy for Q learning. I divide the trajectory into two sections because it bypasses the point (4,7) twice with different action. The red arrow in coordinate (4,8) in Figure 4 mean that it will continue the remaining trajectory at coordinate (4,8) in Figure 5.

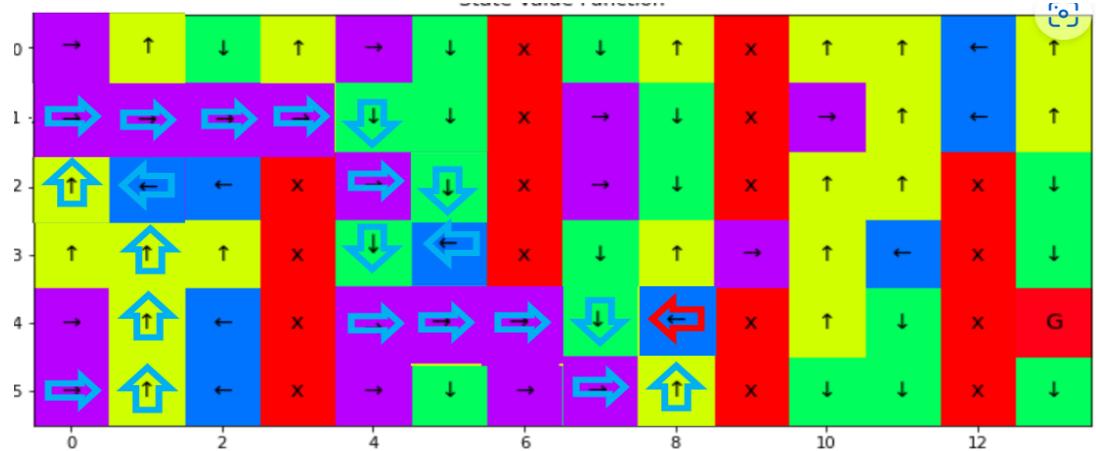


Figure 4: Policy for Q learning with arrows from (5,0) to (4,8)

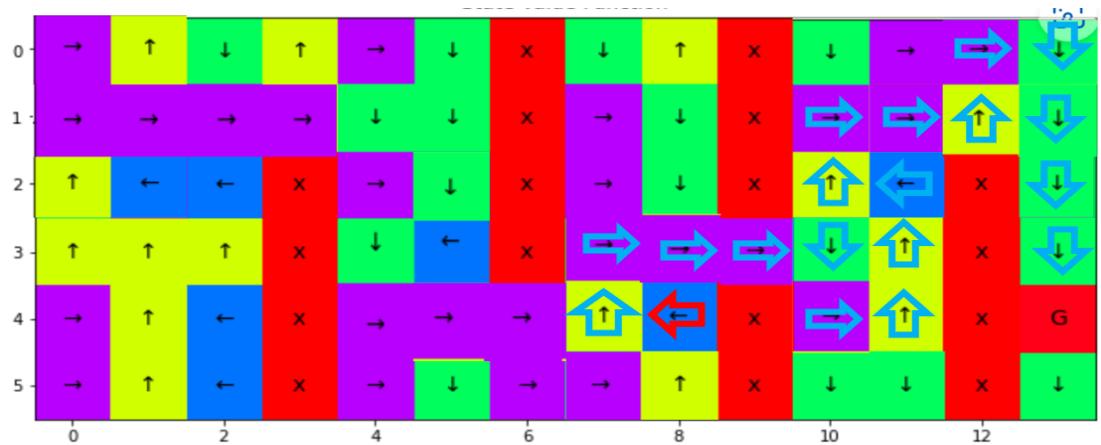


Figure 5: Policy for Q learning with arrows from (4,8) to Goal (4,13)

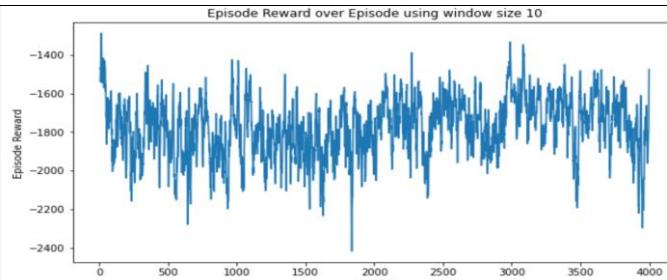
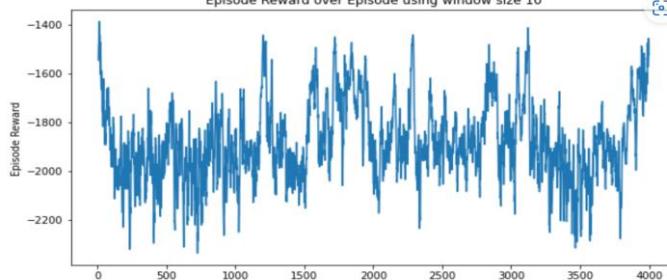
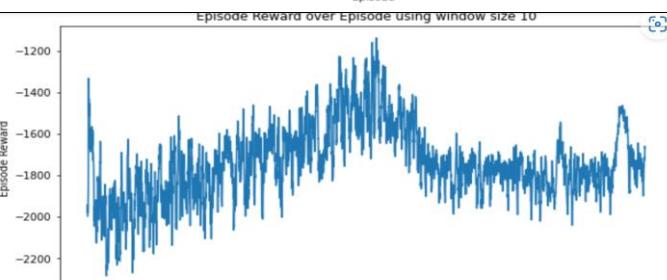
# SARSA

SARSA is an on policy that updates its Q values by using the Q values of the next state and the current policy action. For on policy, the target policy is equal to the behaviour policy. It evaluates and improve the same policy from that policy that is used for the action.

I tested the same parameters in Table 1 and Table 2 for SARSA. However, those parameters could not make SARSA to push the box to the goal. A lot of times, the agent ignores the box to reach the goal by itself. After many finetuning, I figure out that epsilon 0.001 is effective for SARSA. So, I readjusted the epsilon and learning rate for Table 4 and 5. In Table 5, I used learning rate 0.32. It is a weird figure, but it works. I also realised SARSA need longer duration to coverage even though in theory SARSA should be faster than Q learning and SARSA is easily stuck in local minimum.

## Plotting of Episode rewards vs. Episodes with learning rate of 0.5

Table 4: Plotting Episode Reward over Episodes with learning rate of 0.5 and different epsilon values

	Epsilon	Episode Reward over Episodes	Push Box to Goal	Total Reward
Learning rate = 0.5	0.1		No	NIL
	0.01		No	NIL
	0.001		Yes	~1200

## Plotting of Episode rewards vs. Episodes with learning rate of 0.32

Table 5: Plotting Episode Reward over Episodes with learning rate of 0.32 and different epsilon value

	Epsilon	Episode Reward over Episodes	Push Box to Goal	Total Reward
Learning rate = 0.32	0.1	<p>Episode Reward over Episodes</p> <p>Y-axis: Episode Reward (from -2400 to -1400)</p> <p>X-axis: Episode (from 0 to 4000)</p>	No	NIL
	0.01	<p>Episode Reward over Episodes using window size 10</p> <p>Y-axis: Episode Reward (from -2400 to -1400)</p> <p>X-axis: Episode (from 0 to 4000)</p>	No	NIL
	0.001	<p>Episode Reward over Episodes using window size 10</p> <p>Y-axis: Episode Reward (from -2250 to -500)</p> <p>X-axis: Episode (from 0 to 4000)</p>	Yes	~500

From Table 4 and 5, the ideal parameters for SARSA are learning 0.32 and epsilon 0.001 highlighted in red. The total reward is around -500 while in Table 4 is around -1200.

## Value Table

The ideal path generated using SARSA are shown below:

```
print the historical actions: [[4], [1], [1], [1], [3], [1], [4], [4], [4], [4], [1], [4], [2], [2], [2], [3], [2], [4], [4], [4], [2], [4], [1], [3], [1], [4], [4], [4], [2], [4], [1], [1], [3], [1], [4], [4], [1], [4], [2], [2], [2]]
```

Figure 6: Policy for SARSA

The optimal action is highlighted in red for each state.

Table 6: Value Table for Q learning

State	Actions				Best Action
	UP	Down	Left	Right	
((5,0), (4,1))	-339.208	-340.030	-339.518	<b>-292.133</b>	Right
((5,1), (4,1))	<b>-284.152</b>	-334.614	-335.510	-333.086	Up
((4,1), (3,1))	<b>-274.888</b>	-324.817	-324.467	-325.731	Up
((3,1), (2,1))	<b>-264.296</b>	-315.349	-314.839	-316.547	Up
((2,1), (1,1))	-309.004	-305.280	<b>-252.356</b>	-306.159	Left
((2,0), (1,1))	<b>-239.073</b>	-297.128	-299.167	-296.129	Up
((1,0), (1,1))	-287.468	-290.896	-287.348	<b>-226.508</b>	Right
((1,1), (1,2))	-279.619	-282.431	-280.623	<b>-214.720</b>	Right
((1,2), (1,3))	-271.941	-270.599	-270.867	<b>-203.758</b>	Right
((1,3), (1,4))	-241.938	-545.664	-243.040	<b>-194.962</b>	Right
((1,4), (1,5))	<b>-190.097</b>	-236.456	-236.564	-544.051	Up
((0,4), (1,5))	-234.418	-231.728	-232.139	<b>-180.165</b>	Right
((0,5), (1,5))	-227.389	<b>-164.419</b>	-219.966	-324.48	Down
((1,5), (2,5))	-222.630	<b>-155.668</b>	-227.915	-324.16	Down
((2,5), (3,5))	-202.446	<b>-147.755</b>	-202.812	-544.051	Down
((3,5), (4,5))	-197.090	-197.373	<b>-140.690</b>	-543.513	Left
((3,4), (4,5))	-194.359	<b>-132.447</b>	-323.84	-194.430	Down
((4,4), (4,5))	-188.161	-185.788	-543.516	<b>-125.041</b>	Right
((4,5), (4,6))	-182.631	-181.893	-182.511	<b>-118.487</b>	Right
((4,6), (4,7))	-542.438	-168.837	-168.837	<b>-112.802</b>	Right
((4,7), (4,8))	-165.757	<b>-108.006</b>	-165.618	<b>-321.92</b>	Down
((5,7), (4,8))	-161.357	-163.366	-162.804	<b>-102.078</b>	Right
((5,8), (4,8))	<b>-97.038</b>	-158.901	-159.352	-322.56	Up
((4,8), (3,8))	-133.214	-134.525	<b>-80.540</b>	-322.88	Down
((4,7), (3,8))	<b>-73.399</b>	-129.096	-129.164	-130.078	Up
((3,7), (3,8))	<b>-126.486</b>	-126.319	-328.88	<b>-67.143</b>	Right
((3,8), (3,9))	-144.082	-143.325	-142.936	<b>-71.524</b>	Right
((3,9), (3,10))	-322.24	-541.363	-67.104	<b>-66.863</b>	Right
((3,10), (3,11))	-63.647	<b>-63.125</b>	-63.805	-321.28	Down
((4,10), (3,11))	-58.685	-59.235	-322.24	<b>-58.291</b>	Right
((4,11), (3,11))	<b>-54.379</b>	-54.448	-55.443	-321.92	Up
((3,11), (2,11))	<b>-49.367</b>	-49.907	-49.980	-322.24	Up
((2,11), (1,11))	-46.491	-43.760	<b>-43.231</b>	-322.56	Left
((2,10), (1,11))	<b>-35.950</b>	-36.278	-322.88	-40.116	Up
((1,10), (1,11))	-30.067	-32.663	-322.56	<b>-29.541</b>	Right
((1,11), (1,12))	-24.645	-24.916	-24.523	<b>-24.022</b>	Right
((1,12), (1,13))	<b>-19.409</b>	-321.92	-20.907	-20.317	Up
((0,12), (1,13))	<b>-15.983</b>	-14.742	-15.001	<b>-13.683</b>	Right
((0,13), (1,13))	-10.650	<b>-8.860</b>	-11.727	-11.258	Down
((1,13), (2,13))	-5.651	<b>-4.96</b>	-5.694	-6.757	Down
((2,13), (3,13))	-2.150	<b>-2</b>	321.28	-3.423	Down
((3,13), (4,13))			Goal		

Same for SARSA, the step-by-step policy is placed under ANNEX B.

### Policy for SARSA

Surprisingly, the ideal policy for SARSA is similar as Q learning. I divide the trajectory into two sections because it bypasses the point (4,7) twice with different action. The red arrow in coordinate (4,8) in Figure 7 mean that it will continue the remaining trajectory at coordinate (4,8) in Figure 8.

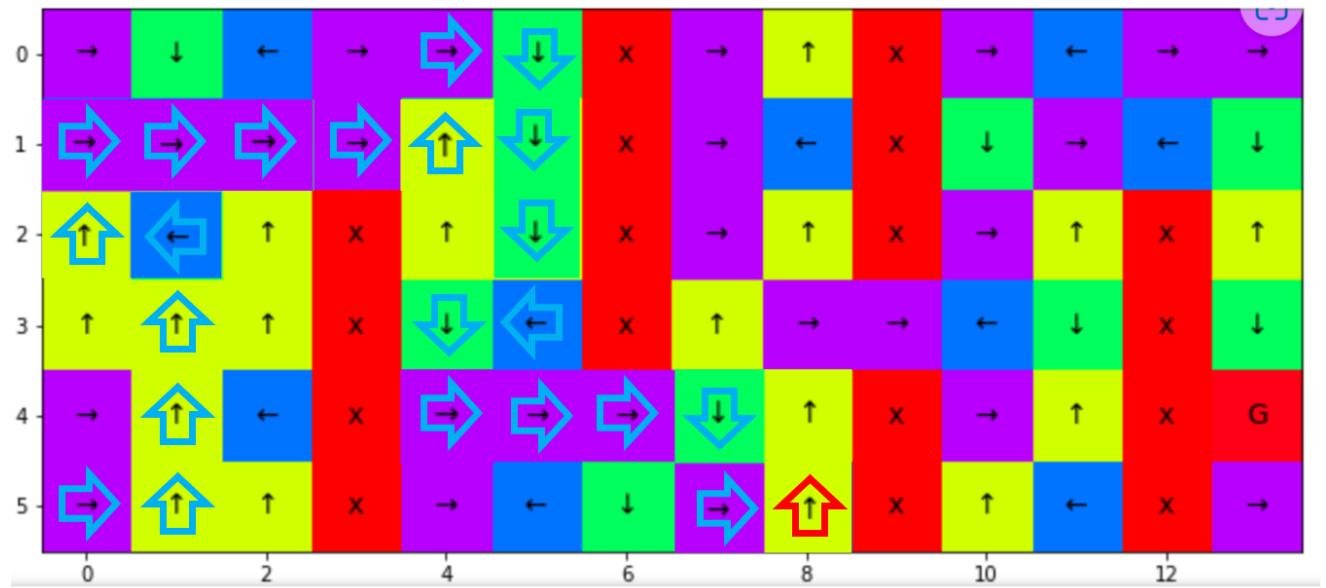


Figure 7: Policy for Q SARSA with arrows from (5,0) to (4,8)

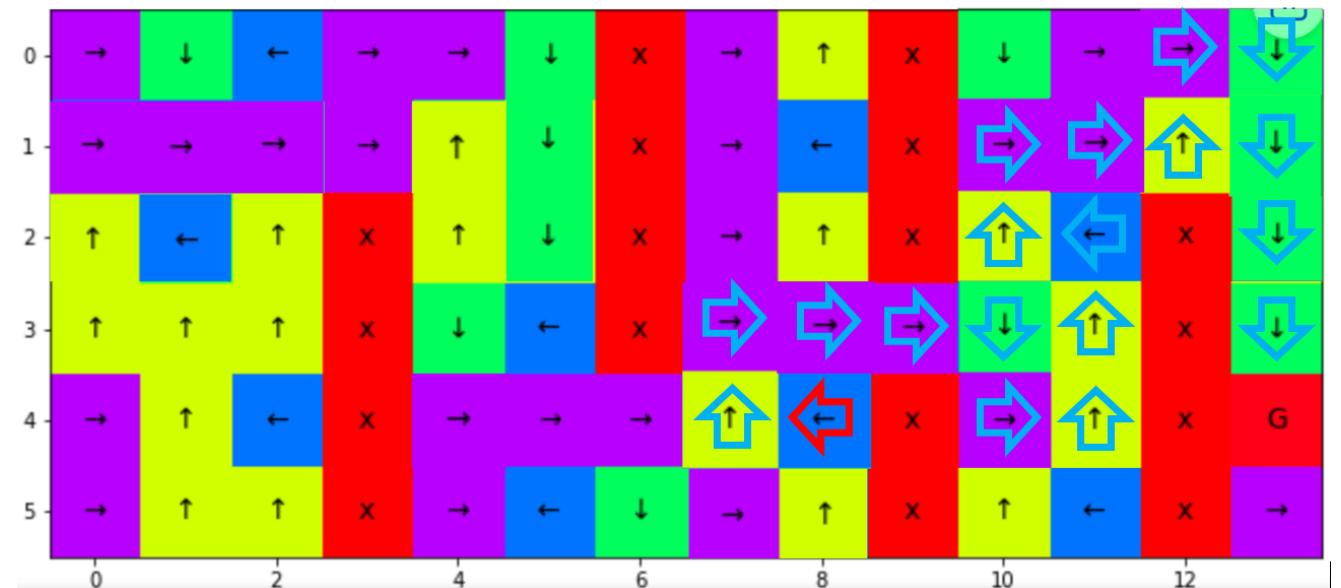


Figure 8: Policy for SARSA with arrows from (4,8) to Goal (4,13)

## Experiment with Epsilon Decay

In theory, we often read about setting the epsilon to be high to let the agent to explore and slowly decay the epsilon. In this experiment, I will be applying epsilon decay to Q learning and SARSA to experiment if epsilon decay indeed improves the performance. Table 7 show the target we need to break to get better result for either Q learning or SARSA. Will by applying epsilon require less step to reach goal?

Table 7: Result before applying epsilon decay

	No. of Episode to reach goal	Step to reach goal
Q learning	~1.6k	41
SARSA	~3.5k	41

The formula of the epsilon decay I will be using is:

$$\epsilon' = \epsilon_{end} + (\epsilon_{start} - \epsilon_{end}) \cdot e^{\left(\frac{-1 \times N_i}{\epsilon_{decay rate}}\right)}$$

where  $\epsilon_{start}$  is the initial epsilon value,  $\epsilon_{end}$  is the final epsilon,  $\epsilon_{decay rate}$  is the rate of the initial epsilon being decayed and  $N_i$  is the  $i^{th}$  number of episodes.

Table 8: Experiment with different initial epsilon, different final epsilon and different decay rate for Q learning and SARSA

With decay epsilon	No. of episode	$\epsilon_{start}$	$\epsilon_{end}$	Decay rate	No. of Episode to reach goal	Steps to reach goal
Q learning	2k	1	0.01	200	~1.9k	41
Q learning	2k	0.5	0.01	200	~1.7k	41
Q learning	2k	0.1	0.01	200	~1.71k	41
Q learning	10k	0.1	0.01	500	~1.7k	41
Q learning	2k	0.01	0.001	200	~1.5k	41
Q learning	10k	0.01	0.001	500	~1.5k	41
SARSA	4k	1	0.001	500	NIL	NIL
SARSA	4k	0.5	0.001	200	NIL	NIL
SARSA	4k	0.01	0.001	200	~3.5k	41
SARSA	10k	0.01	0.001	50	~3.6k	41
SARSA	4k	0.001	0.0005	100	~2.6k	41
SARSA	10k	0.001	0.0005	100	~2.6k	41

From the Table 8, I experiment many different decays rate with different initial epsilon, and I also discovered that Q learning for this environment the epsilon must be around 0.01 to reach the destination. And for SARSA the epsilon must be as low as 0.001. I tried exploration by setting initial epsilon as high as 1 and slowly decay. It only starts to optimal when the epsilon is around 0.001 for SARSA and 0.01 for Q learning. The reasons for this will be discuss in the Discussion section.

For Q learning, 0.01 initial epsilon, 0.001 end epsilon with 200 decay rate is the best parameters. It uses 1.5k episodes to reach the goal which is slightly better than without epsilon decay. However, no matter how I finetuned the epsilon the optimal steps to reach goal is still 41.

For SARSA, 0.001 initial epsilon, 0.0005 with decay rate of 100 is the best parameters. It uses around 2.6k episodes to reach goal which is relatively improve compared to one without epsilon decay. Similarly, the number of steps is still 41.

## Experiment with Learning Rate Decay

There are still other parameters we can experiment with, which is the learning rate. In deep learning, we often use learning rate scheduler to slow down the learning when converging. The purpose of applying this is to prevent the gradient being stuck at the local minima. Therefore, I will be using the concept to decay the learning rate decay using step decay.

The formula for the step decay is as follow:

$$LR_{Final} = LR_{Initial} \cdot (decay\ rate)^i$$

where  $LR_{Initial}$  is the initial learning rate,  $decay\ rate$  is the rate of learning rate to decay and  $i$  is the number of iterations.

Table 9: Experiment with different initial learning rate and different decay rate for Q learning and SARSA

	No. of episode	Initial lr	Final lr	Decay rate	No. of Episode to reach goal	Step to reach goal
<b>Q learning</b>	<b>2k</b>	<b>1.0</b>	<b>0.833</b>	<b>0.0001</b>	<b>~1.6k</b>	<b>41</b>
Q learning	2K	1.0	0.505	0.0005	~1.9k	41
Q learning	10k	1.0	0.5	0.0001	~1.5k	41
Q learning	10k	1.0	0.09	0.001	~2.3k	41
Q learning	10k	0.5	0.25	0.0001	~2.3k	41
SARSA	4k	1.0	0.294	0.0006	NIL	NIL
SARSA	4k	0.5	0.118	0.0009	~3.9k	52
<b>SARSA</b>	<b>4k</b>	<b>0.35</b>	<b>0.257</b>	<b>0.00009</b>	<b>~3.5k</b>	<b>41</b>

Learning rate is important parameters. I experiment different range of learning rate and Q learning require a high learning rate to learn faster. Slower learning rate need longer episodes to reach the goal. For SARA, the learning rate is around 0.3. If the learning rate is higher than that, it will not reach to the goal.

Therefore, from Table 9 the ideal learning rate decay for Q learning start from 1 then apply 0.0001 decay rate. It required 1.6k episodes to reach the goal which is around the same as without learning rate decay. For SARSA, the ideal learning rate decay start from 0.35 and end of 0.257.

## Experiment with Epsilon Decay & Learning rate Decay

Next, we combined both epsilon decay and learning rate decay to Q Learning and SARSA.

Table 10: Result after applying epsilon decay and learning rate decay for Q learning and SARSA

	No. of episode	No. of Episode to reach goal	Step to reach goal
Q learning	2k	~1.5k	41
SARSA	4K	~2.7k – 3.1k	41

From Table 8 and Table 9, I select the best epsilon decay and learning rate decay highlighted in red to Q learning and SARSA. Apparently, SARSA improve from 3.6k episode to 2.7k episode to reach the goal and Q learning improve from 1.6k to 1.5k.

Ultimately, our goal is to improve the number of steps to reach goal. It seemed that 41 steps are the optimal minimum required steps for both Q learning and SARSA.

## Discussion

Now, we need to discuss why exploration does not work. In this environment, exploration does not work. The diagram below, Figure 9, has labelled five different sections. For example, we train 1k episodes and epsilon is set to 0.9 with decay. First the agent will explore around Section 1 in Figure 9 and let's say it managed to reach Section 2. The problem arises. The epsilon already decayed and there are still 4 more sections. The agent is not able to explore another four sections unless the epsilon is reset every time it reaches to the different section. Then it will start to explore. But if we reset the epsilon, isn't it take longer episodes to reach the goal? Thus, setting epsilon to be very low is the optimal solution for this environment as it chooses the best actions in every states.

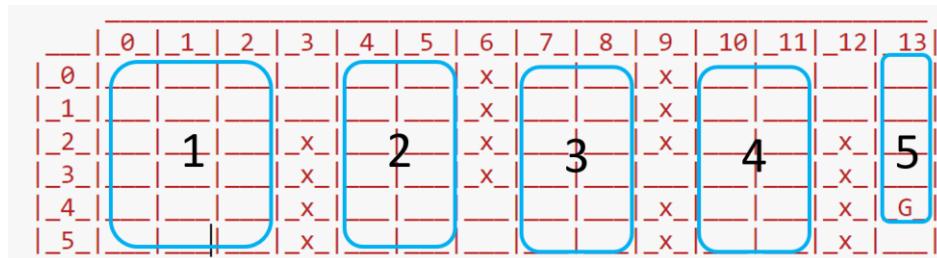


Figure 9: Environment divided into 5 sections

Unlike the Cliffwalking example, the exploration method works because the agent can explore all the states (from (3,0) to (2,11)) before exploitation.

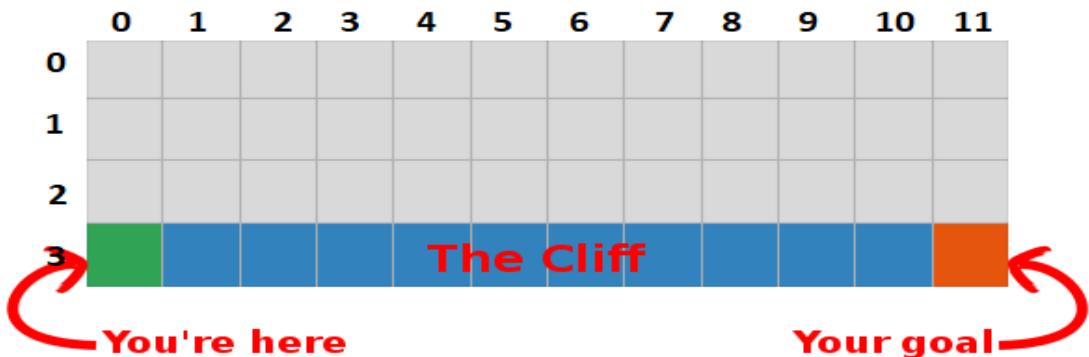


Figure 10: Example of an environment from gym - Cliffwalking

## Evaluation

Table 11: The table that consists of all the optimal results experimented using Q learning and SARSA

	Learning Rate	Gamma Value	Epsilon	No. of episodes to reach G (minimum steps)	No. of steps to reach G	Total Rewards
Q learning	0.9	0.98	0.01	~1600	41	-388
SARSA	0.32	0.98	0.001	~3600	41	-388
Q learning (with learning rate decay and epsilon decay)	1.0 – 0.83	0.98	0.01-0.001	~1500	41	-388
SARSA (with learning rate decay and epsilon decay)	0.35-0.257	0.98	0.001-0.0005	~2700	41	-388

In conclusion, the minimum steps to reach goal is 41 for both Q learning and SARSA. The only different is the number of episodes to reach the destination. For Q learning, by applying learning rate decay and epsilon decay increase the performance by reaches the goal earlier and for SARSA by applying learning rate decay and epsilon decay it reduces the number of episodes from 3600 to 2700. Since all the algorithms took 41 steps to reach, we must determine which algorithm reaches the goal first which is Q learning with learning rate decay and epsilon decay.

## Extra (Q Learning without Box coordinate)

Out of curiosity, I decided to convert state into grid number instead of coordinates. For example, the state initially is (5,0) will be converted to 70. I converted by using the x-axis of the coordinates of the state multiply by the row of the environment which is 14 plus y-axis.

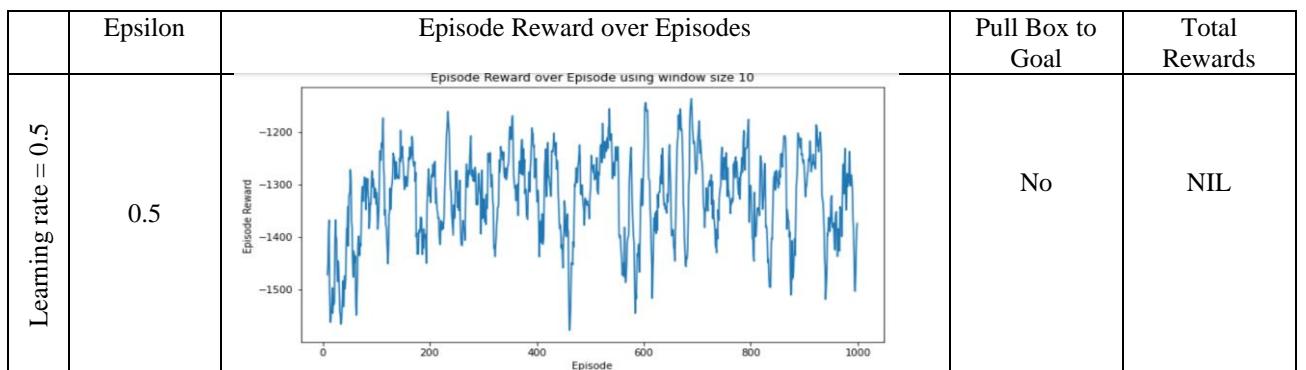
I removed the box coordinates and train using Q learning,

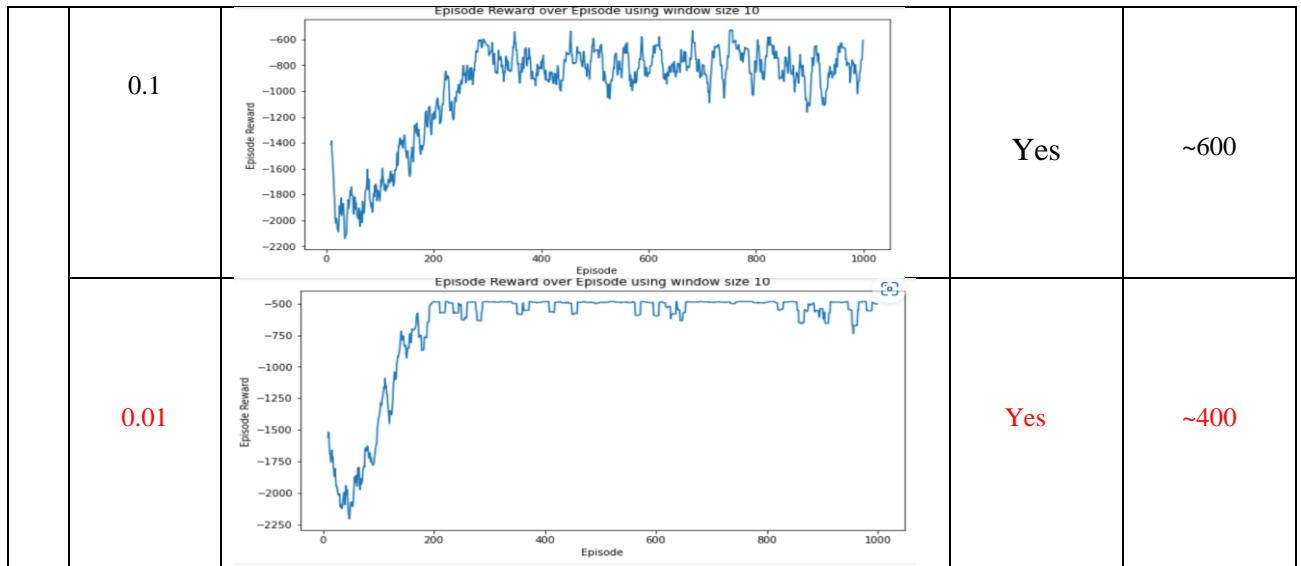
The purpose of this experiment is to analysis which trajectory does the agent move to the destination without the box location. Will it take 41 steps too?

Similarly, I will be using Q learning with learning rate 0.5 and 0.9.

### Plotting of Episode rewards vs. Episodes with learning rate of 0.5

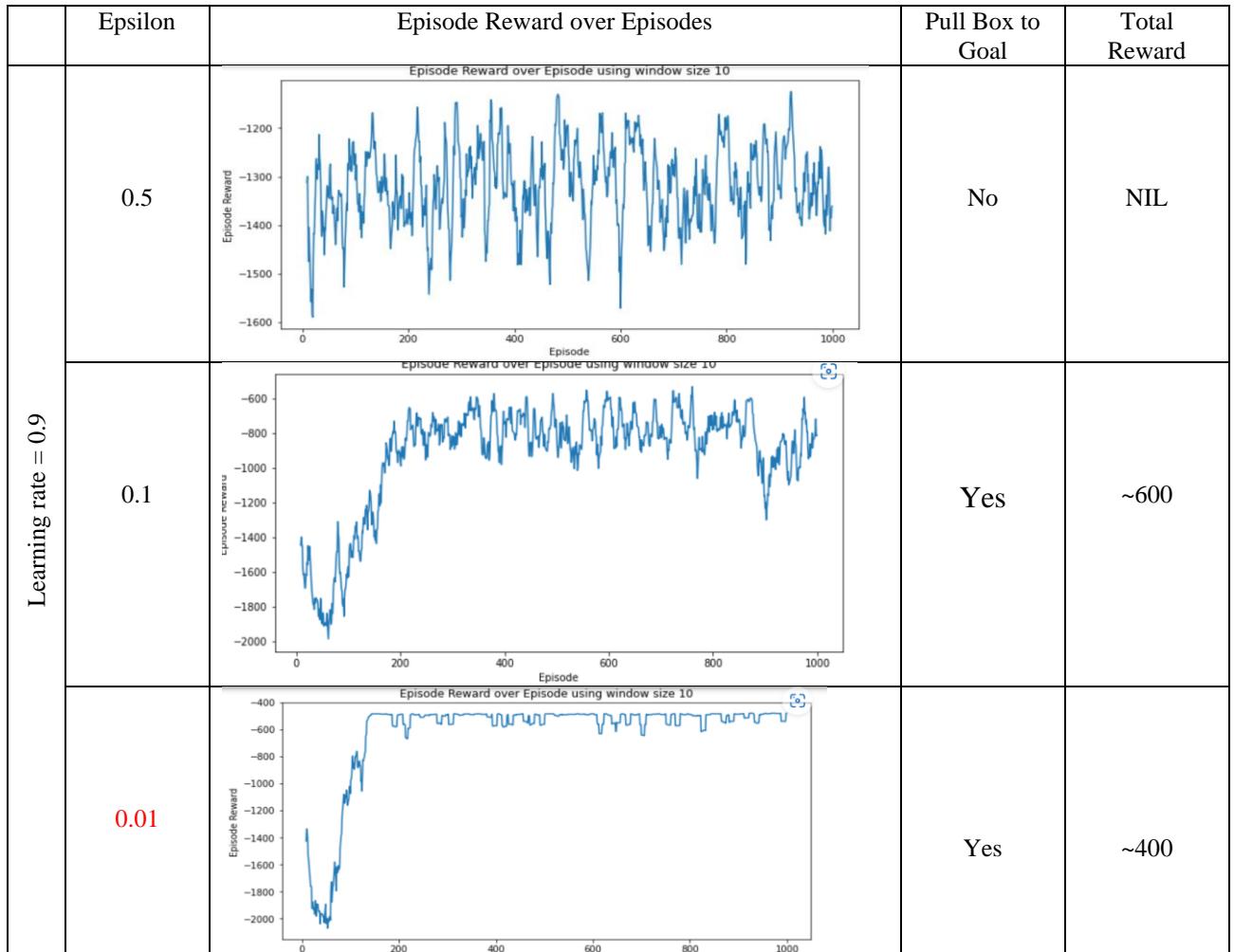
Table 12: Plotting Episode Reward over Episodes with learning rate of 0.5 and different epsilon values





### Plotting of Episode rewards vs. Episodes with learning rate of 0.9

Table 13: Plotting Episode Reward over Episodes with learning rate of 0.9 and different epsilon values



From Table 13, 0.01 epsilon and 0.9 learning rate took 200 episodes and only took 26 steps to reach the goal.

## V Table

```
print the historical actions: [[1], [4], [1], [1], [4], [1], [4], [4], [2], [2], [2], [4], [4], [1], [4], [4], [1], [4], [4], [4], [2], [2], [2]]]
```

Table 14: V Table for Q learning without Box

	Actions				Best Action
	UP	Down	Left	Right	
70 (5,0)	-372.707	-401.168	-398.185	-401.606	Up
56 (4,0)	-395.518	-405.249	-396.591	-363.824	Right
57 (4,1)	-355.964	-400.171	-393.818	-910.8	Up
43 (3,1)	-347.576	-386.819	-391.968	-990.415	Up
29 (2,1)	-376.251	-385.431	-386.431	-338.643	Right
15 (1,1)	-372.437	-372.544	-381.557	-366.634	Right
16 (1,2)	-366.244	-370.359	-361.638	-322.305	Right
17 (1,3)	-355.142	-917.1	-352.014	-312.525	Right
18 (1,4)	-356.293	-301.549	-343.019	-1211.667	Down
32 (2,4)	-342.408	-291.375	-1005.03	-334.212	Down
46 (3,4)	-329.761	-282.015	-918	-320.748	Down
60 (4,4)	-318.522	-311.521	-1010.07	-273.484	Right
61 (4,5)	-310.547	-310.965	-307.778	-263.759	Right
62 (4,6)	-917.1	-289.225	-290.518	-252.816	Right
63 (4,7)	-240.628	-283.720	-285.410	-267.769	Up
49 (3,7)	-254.425	-261.537	-913.5	-227.172	Right
50 (3,8)	-259.835	-236.533	-239.662	-212.421	Right
51 (3,9)	-918.9	-918.9	-223.388	-196.347	Right
52(3,10)	-178.926	-200.541	-197.892	-198.042	Up
38 (2,10)	-160.128	-185.670	-920.7	-176.306	Up
24 (1,10)	-160.117	-163.168	-1011.96	-139.927	Right
25(1,11)	-142.189	-171.418	-174.618	-118.293	Right
26(1,12)	-117.947	-1015.56	-114.821	-95.197	Right
27(1,13)	-96.008	-70.609	-100.867	-102.131	Down
41 (2,13)	-94.385	-46.54	-921.6	-69.859	Down
55(3,13)	-49.779	-23	-1013.21	-46.3788	Down
69(4,13)			GOAL		

From the V table, we can see that all the coordinates that the agent move are the box location. It moves pass (4,1) then (3,1) and eventually reach the goal. If the agent follow the same box location does it mean this method work? But I did not include this as my solution. This is just an experiment.

### Policy and Value Table for Q learning without Box coordinate

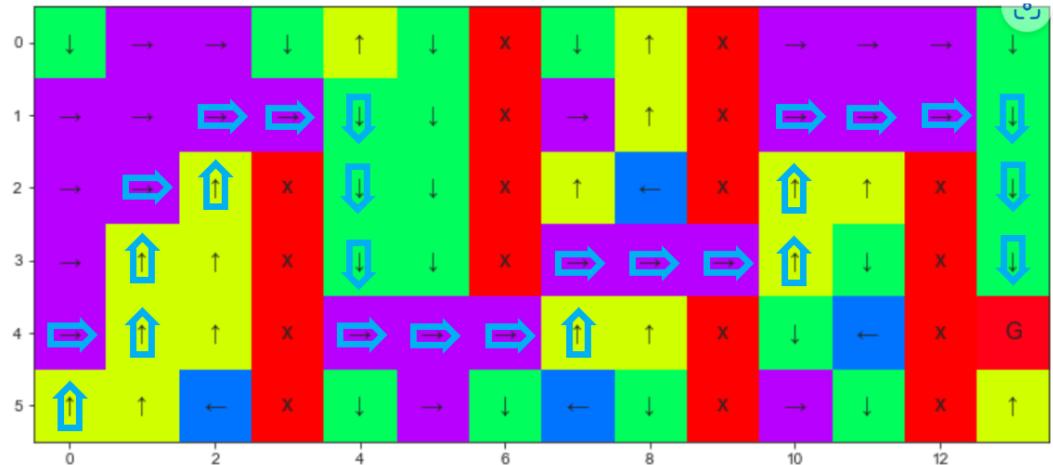


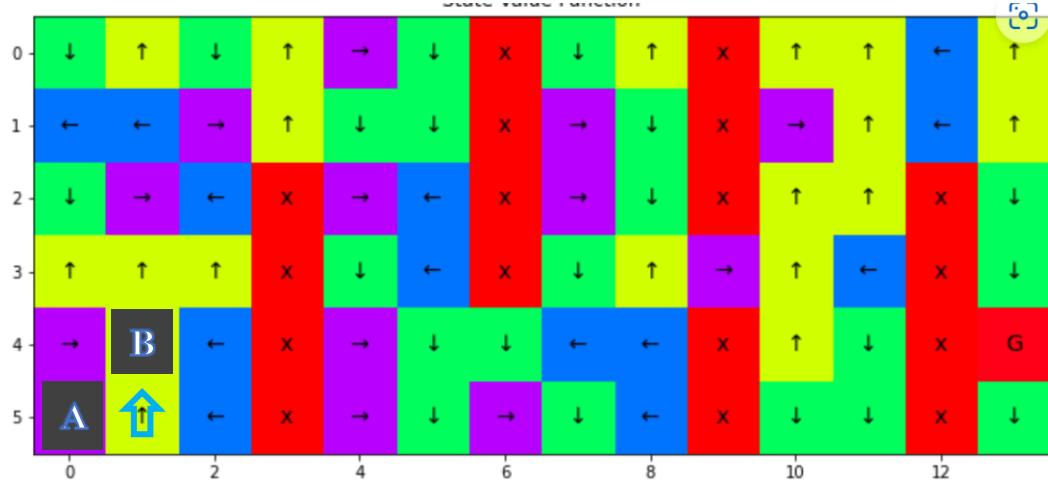
Figure 11: Policy for Q learning (with no box) with arrows from (5,0) to Goal (4,13)

0	-367.7	-364.9	-357.1	-350.8	-344.4	-337.0	0.0	-247.6	-241.6	0.0	-143.0	-118.9	-120.5	-95.3
1	-370.1	-366.6	-322.3	-312.5	-301.5	-323.9	0.0	-236.4	-243.3	0.0	-139.9	-118.3	-95.2	-70.6
2	-377.8	-338.6	-331.0	0.0	-291.4	-319.7	0.0	-242.0	-244.2	0.0	-160.1	-157.6	0.0	-46.5
3	-383.7	-347.6	-371.3	0.0	-282.0	-308.4	0.0	-227.2	-212.4	-196.3	-178.9	-174.6	0.0	-23.0
4	-363.8	-356.0	-380.6	0.0	-273.5	-263.8	-252.8	-240.6	-251.8	0.0	-190.3	-175.9	0.0	0.0
5	-372.7	-395.8	-389.0	0.0	-305.4	-299.5	-283.1	-271.4	-265.1	0.0	-199.1	-184.7	0.0	0.0

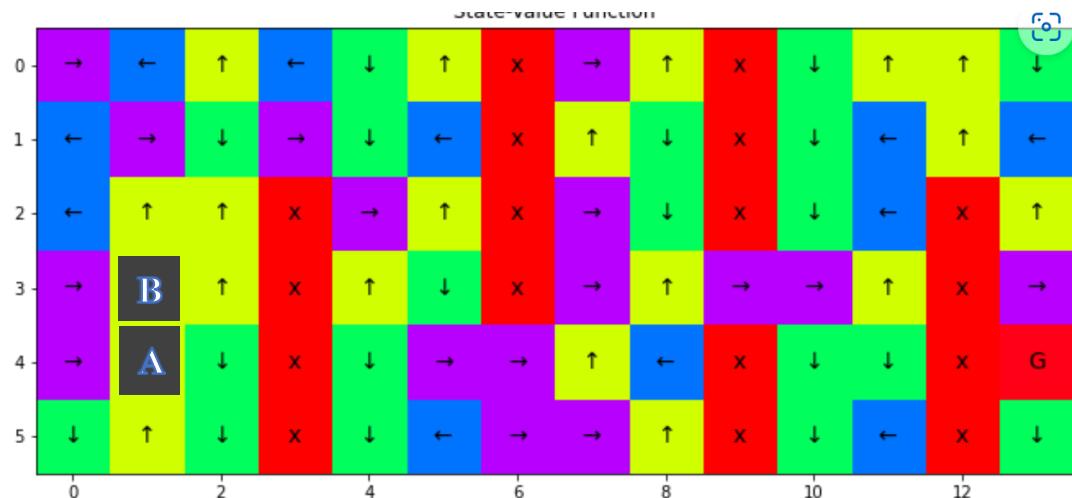
Figure 12: Policy for Q learning (with no box) with values for 6 x 14

## ANNEX A (Q learning policy step by step)

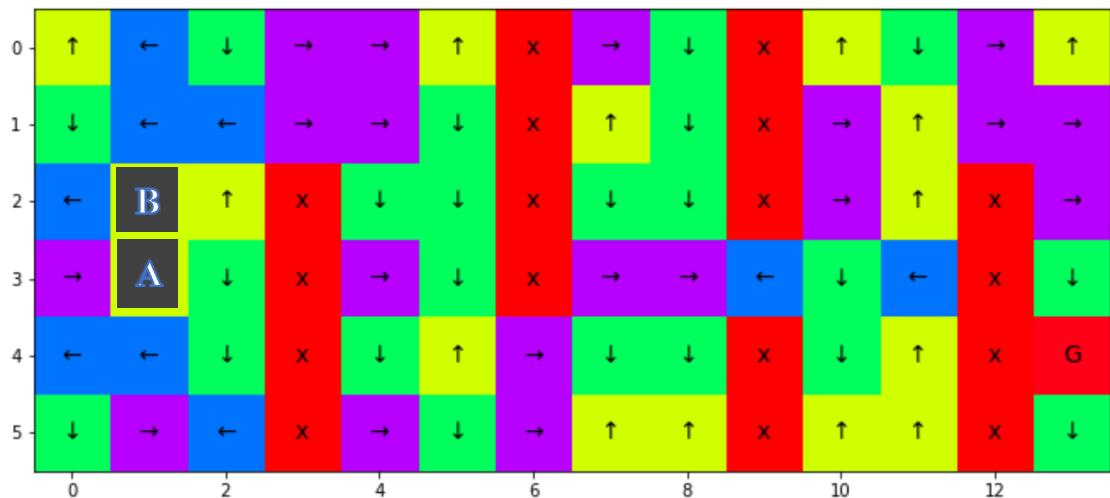
Agent (5,0) move → to (5,1) then move ↑ to Box (4,1)



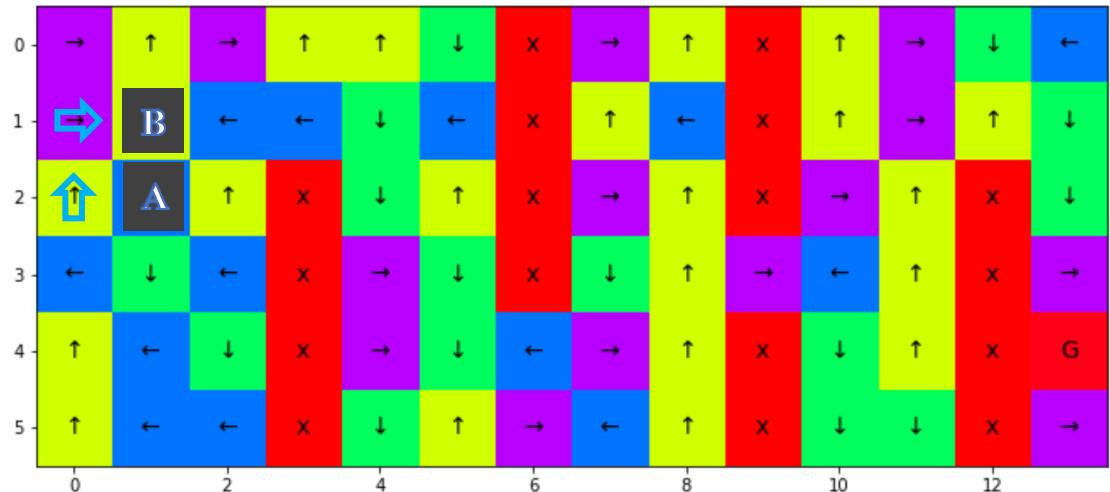
Agent (4,1) move ↑ to Box (3,1)



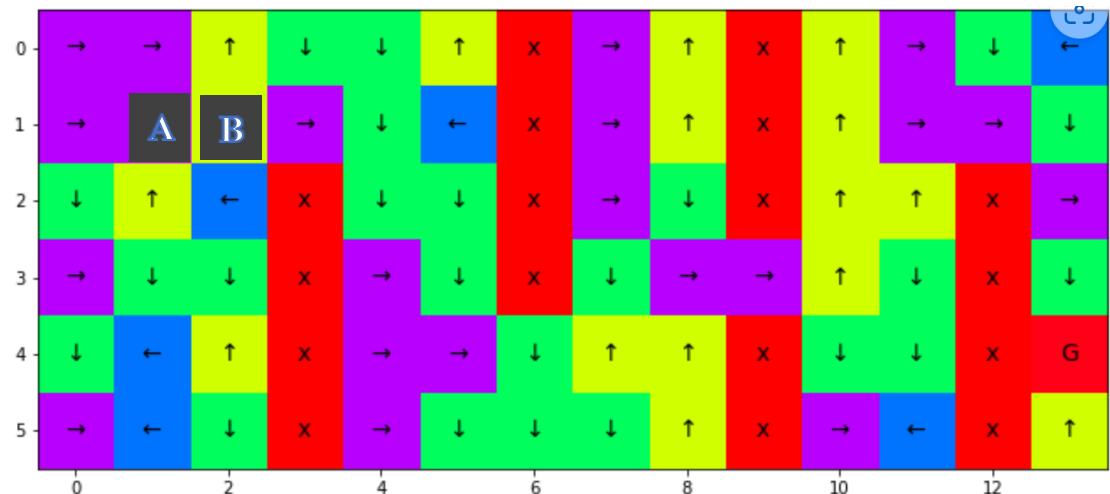
Agent (3,1) move ↑ to Box (2,1)



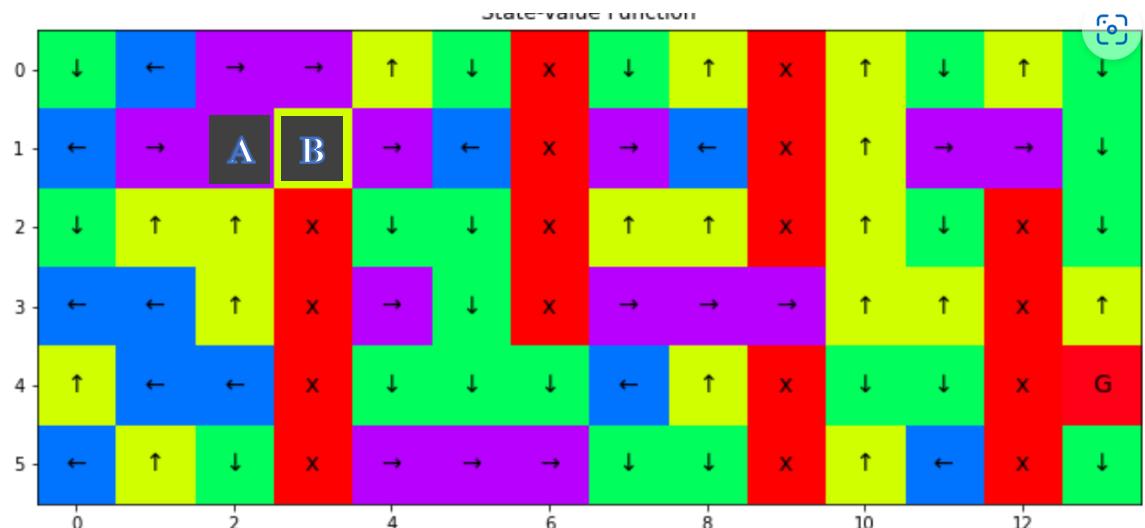
Agent (2,1) move  $\leftarrow$  to (2,0) then  $\uparrow$  to (1,0) then  $\rightarrow$  to Box (1,1)



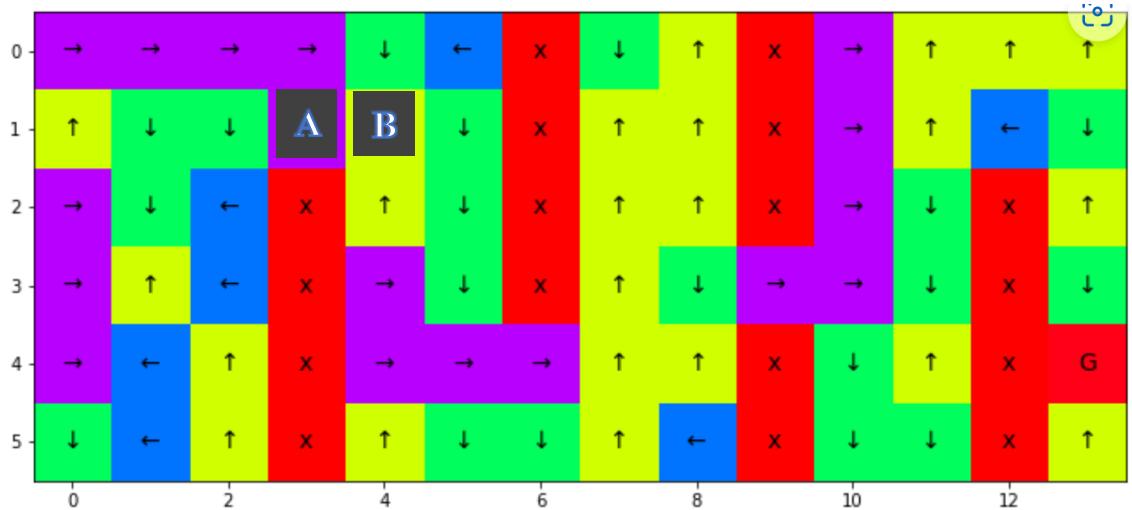
Agent (1,1) move  $\rightarrow$  to Box (1,2)



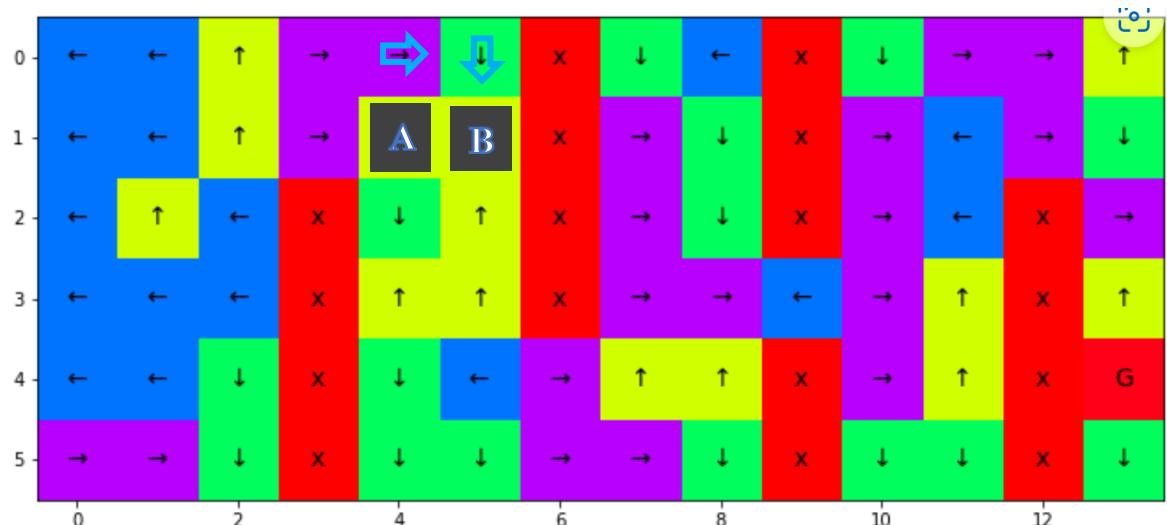
Agent (1,2) move  $\rightarrow$  to Box (1,3)



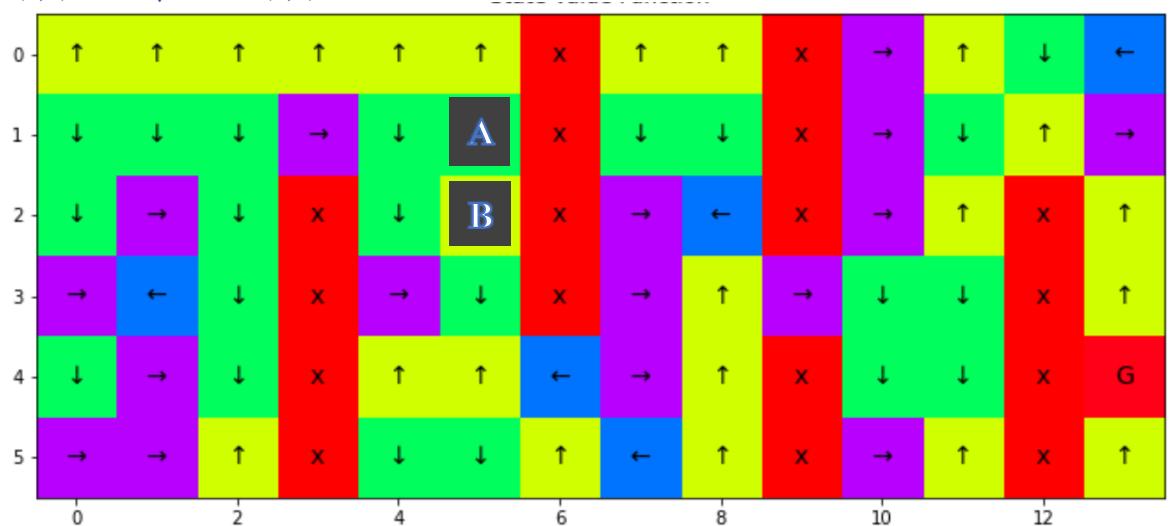
Agent (1,3) move → to Box (1,4)



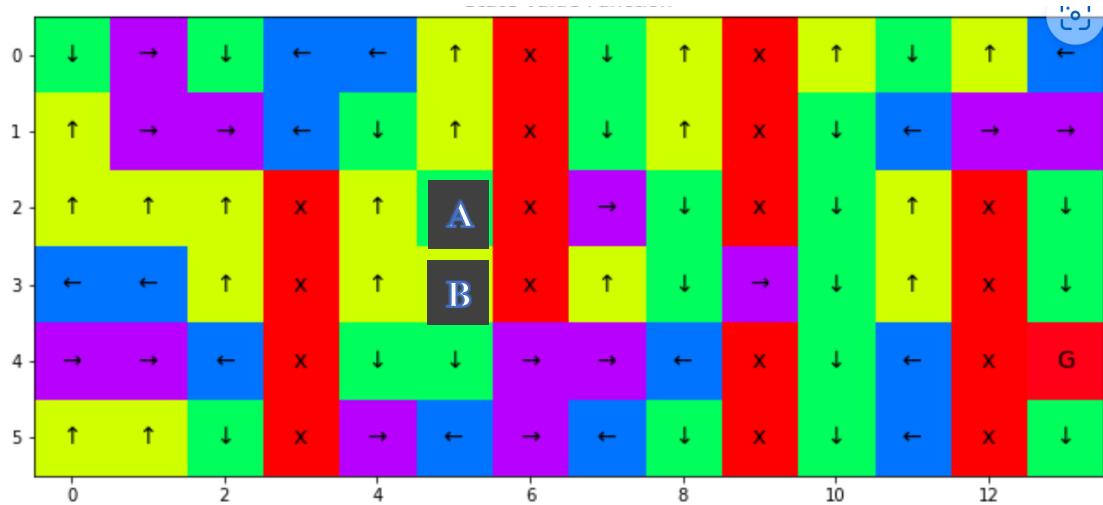
Agent (1,4) move ↑ to (0,4) then → to (0,5) then move ↓ to Box (1,5)



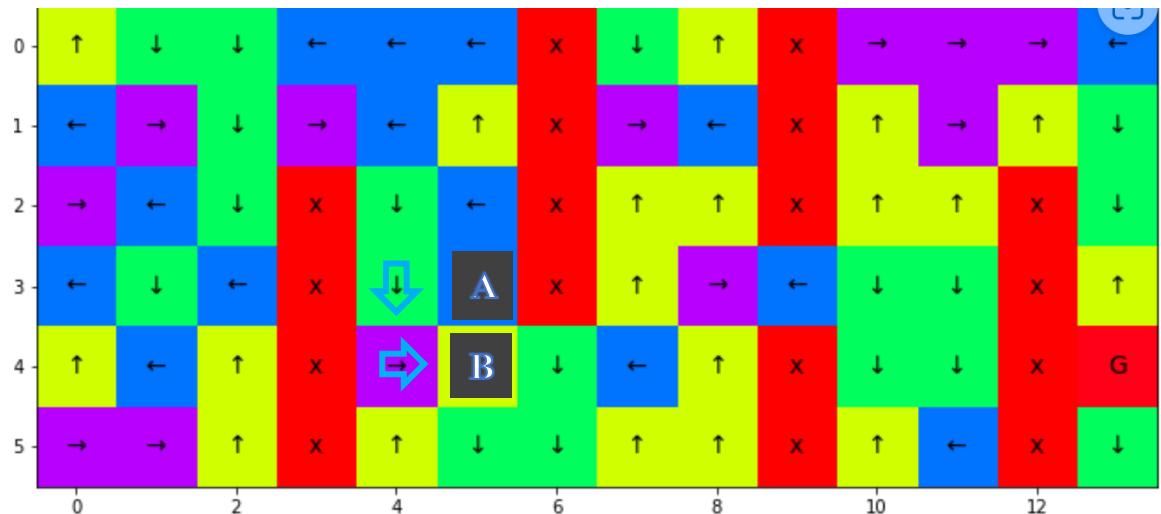
Agent (1,5) move ↓ to Box (2,5)



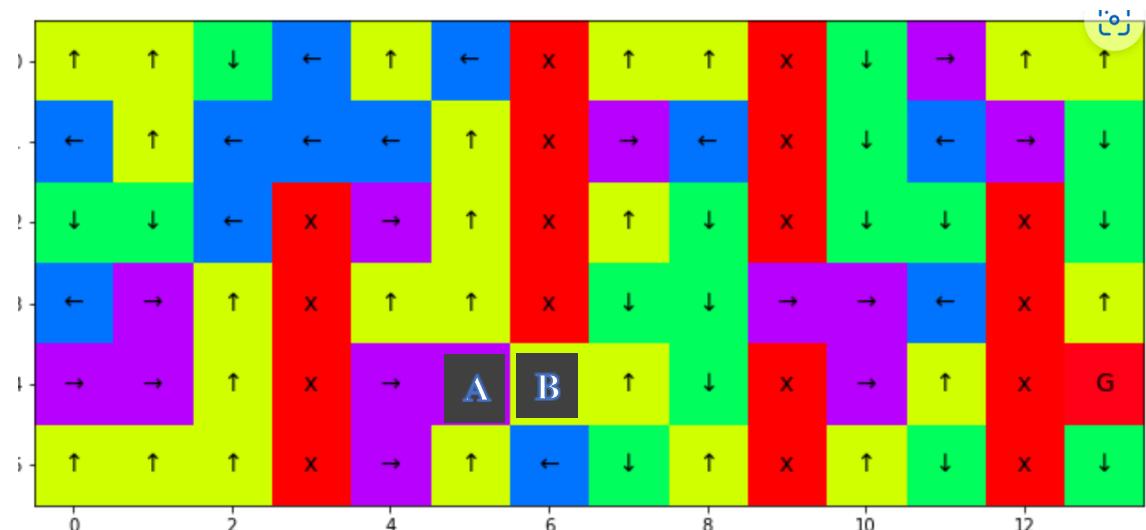
**Agent (2,5) move ↓ to Box (3,5)**



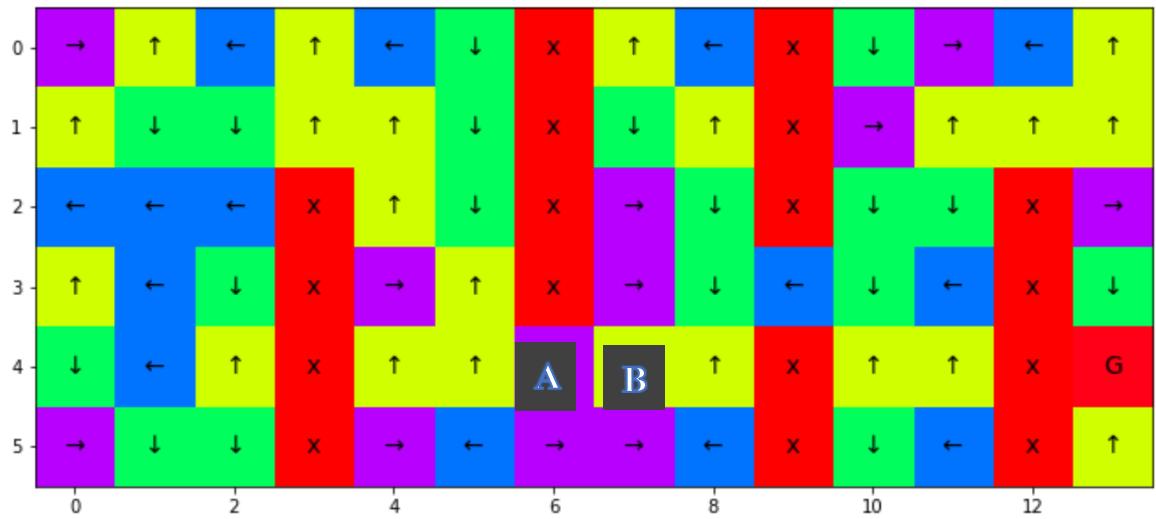
**Agent (3,5) move ← to (3,4) the move ↓ to (4,4) then move → to Box (4,5)**



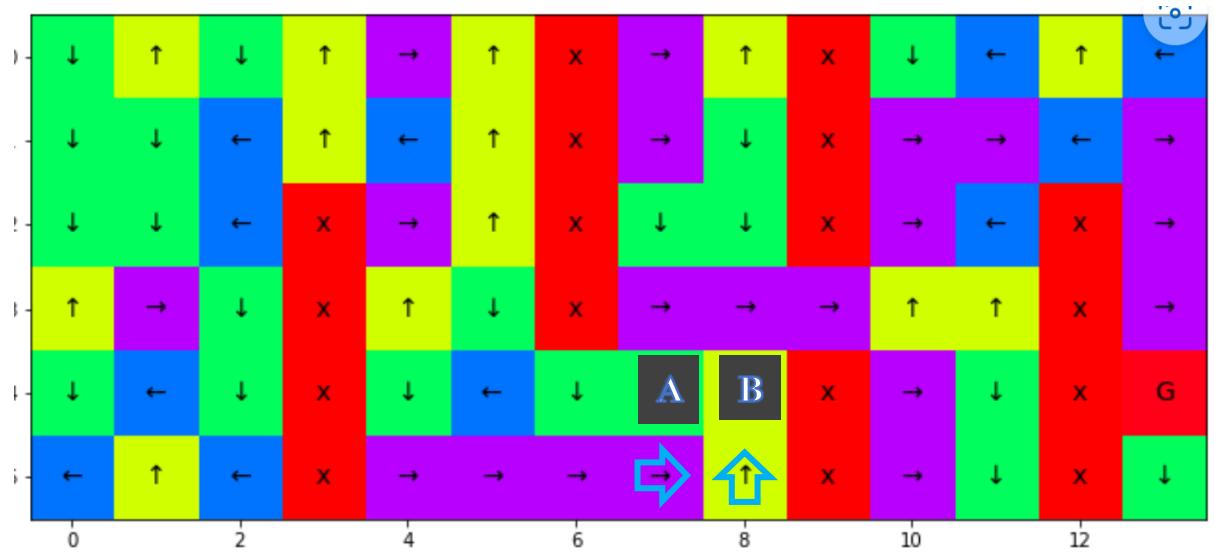
**Agent (4,5) move → to Box (4,6)**



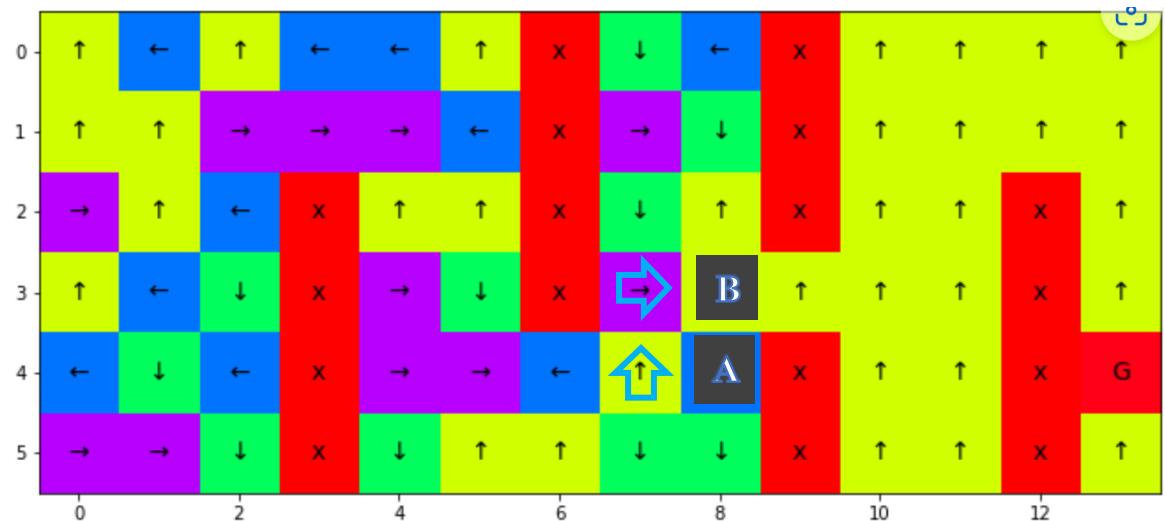
Agent (4,6) move → to Box (4,7)



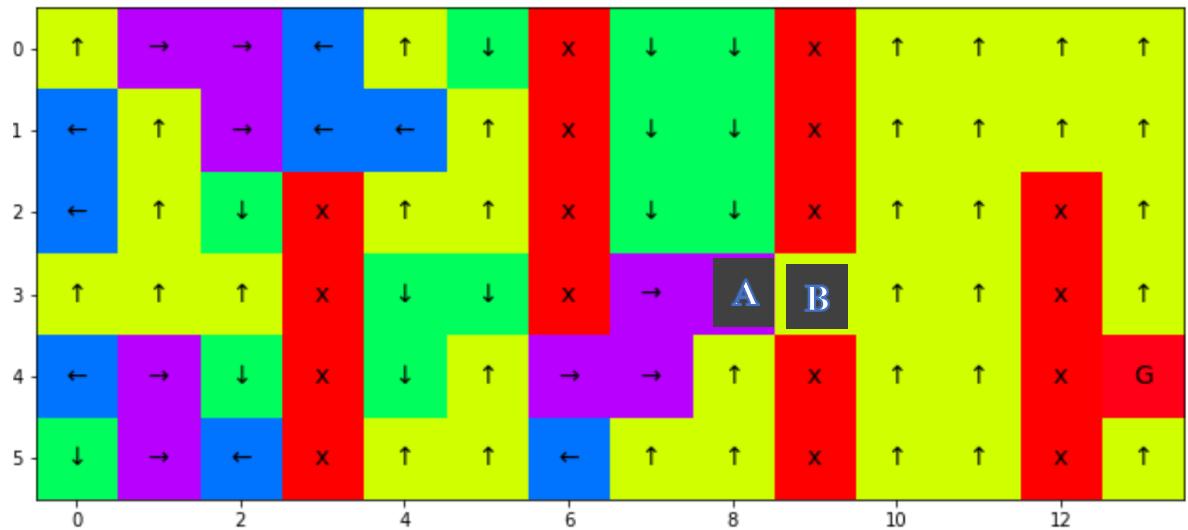
Agent (4,7) move ↓ to (5,7) then move → to (5,8) then move ↑ to Box (4,8)



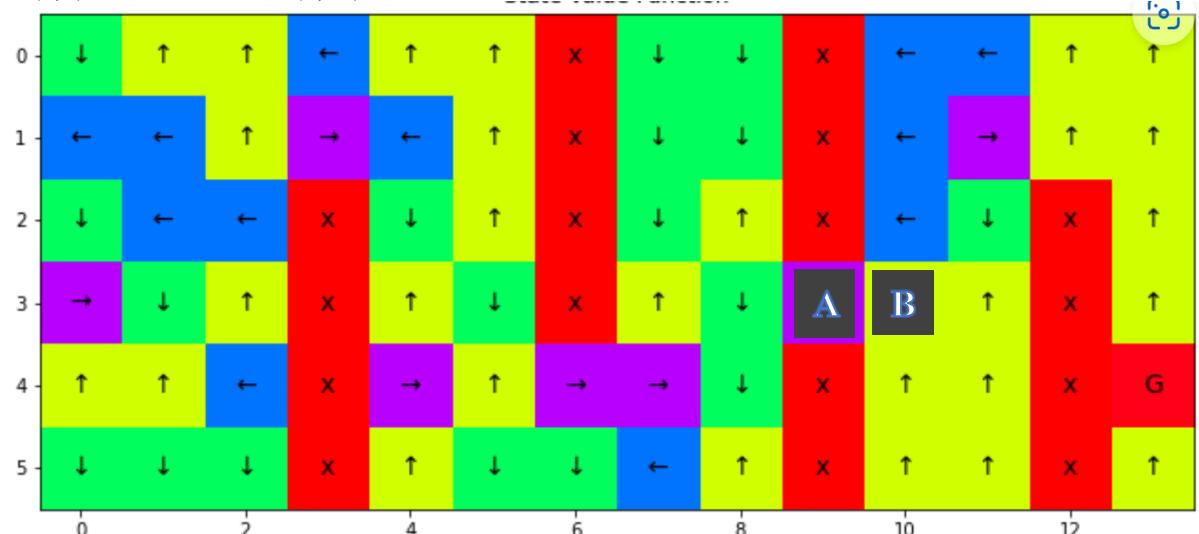
Agent (4,8) move ← to (4,7) then ↑ to (3,7) then → to Box (3, 8)



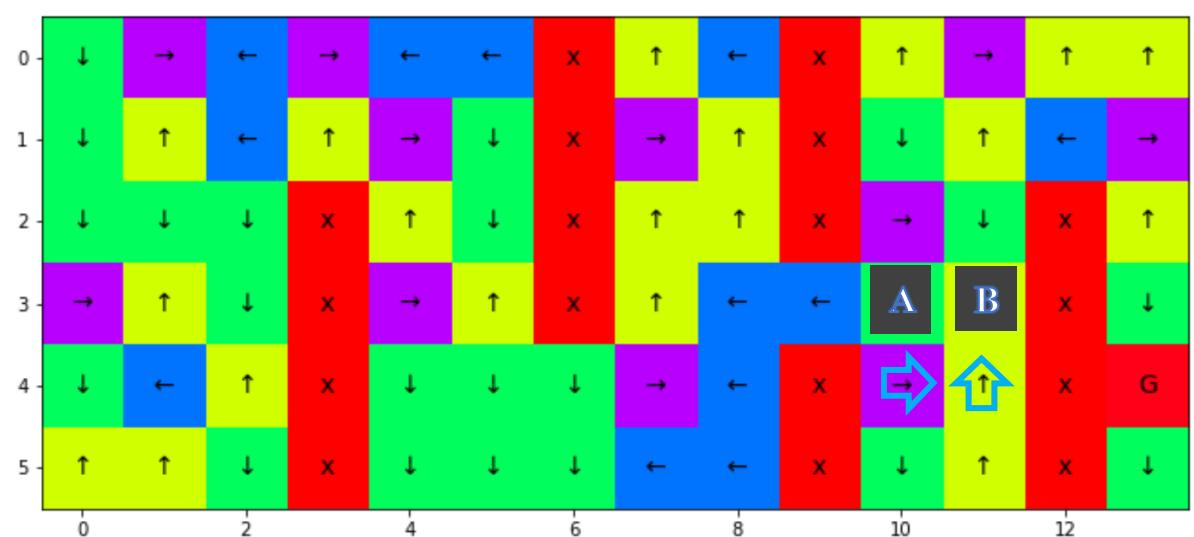
Agent (3,8) move → to Box (3,9)



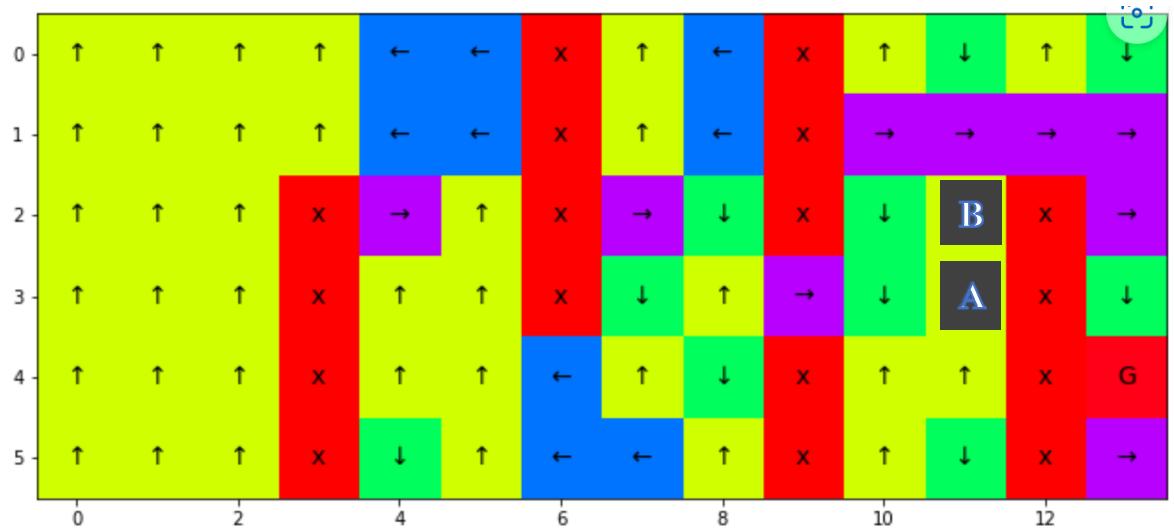
Agent (3,9) move → to Box (3,10)



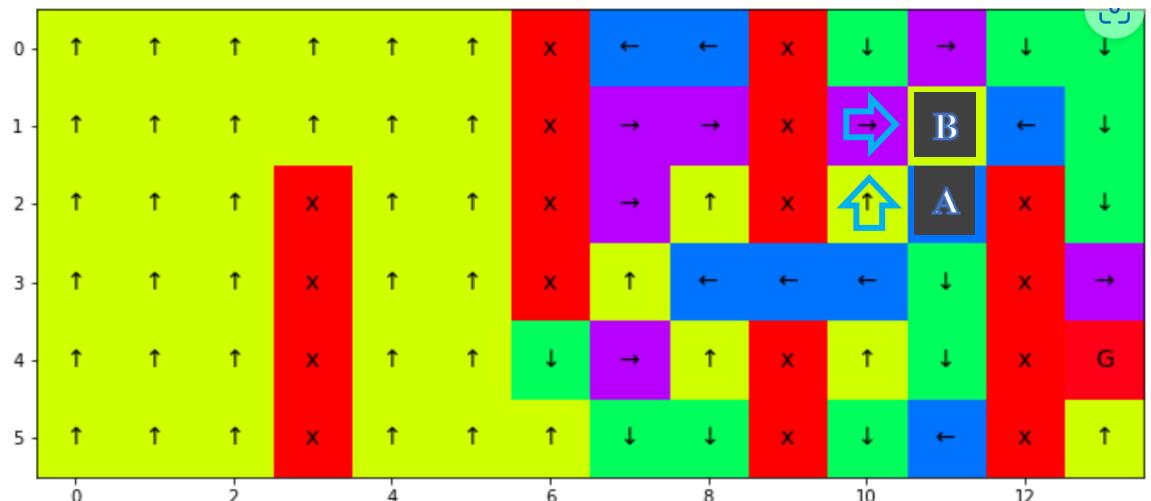
Agent (3,10) move ↓ to (4,10) then move → to (4,11) then move ↑ to Box (3,11)



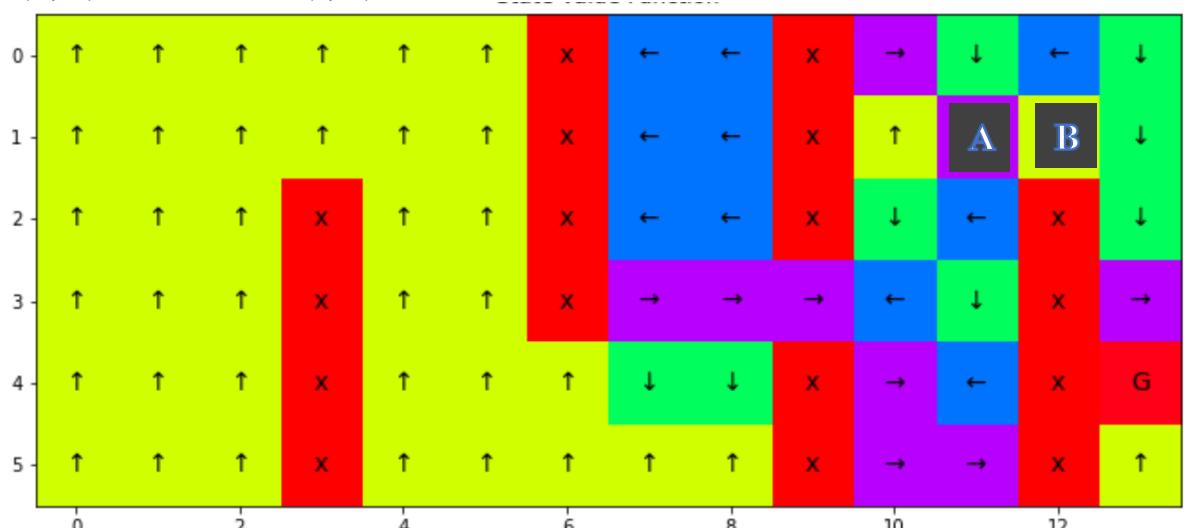
**Agent (3,11) move  $\uparrow$  to Box (2,11)**



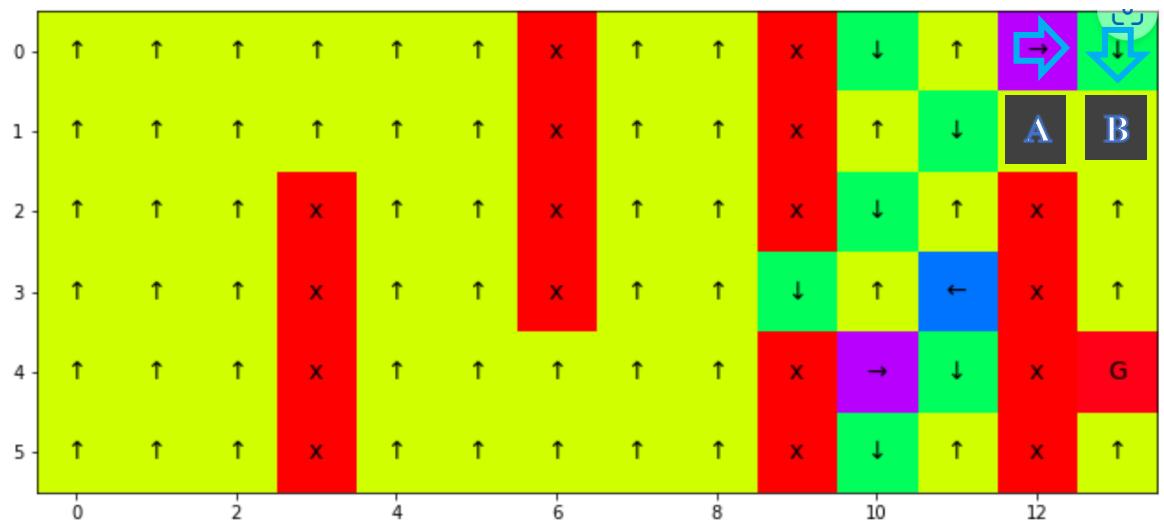
**Agent (2,11) move  $\leftarrow$  to (2,10) then  $\uparrow$  to (1,10) then  $\rightarrow$  to Box (1 ,11)**



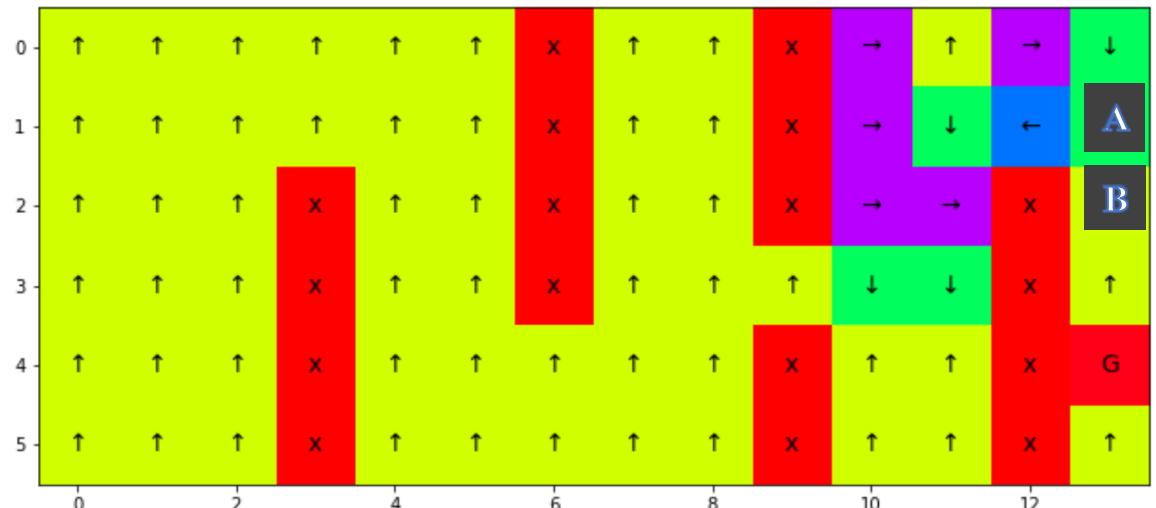
**Agent (1 ,11) move  $\rightarrow$  to Box (1,12)**



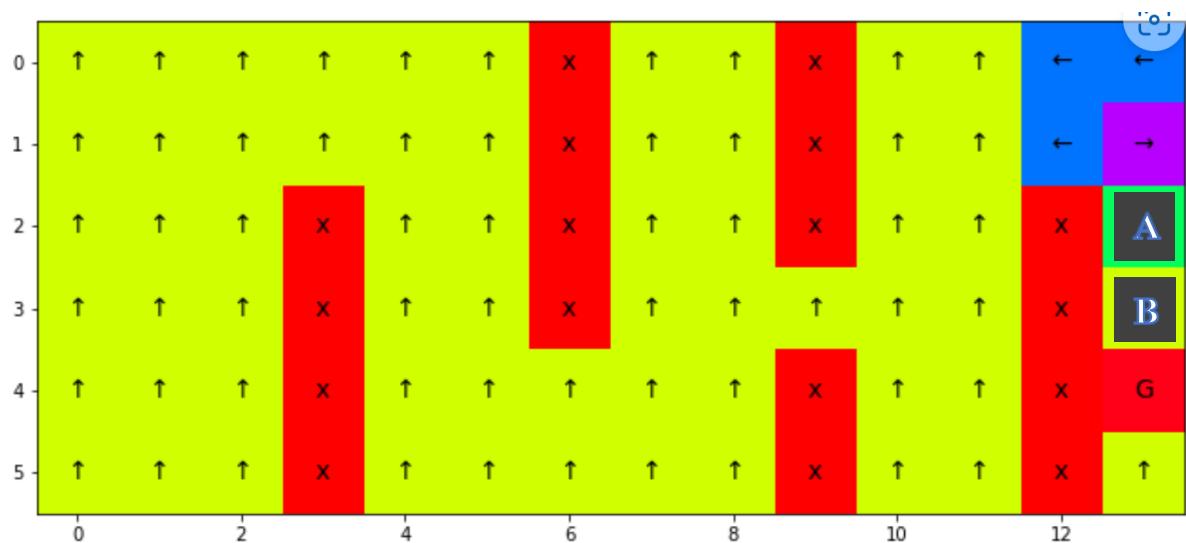
**Agent (1 ,12) move  $\uparrow$  to (0,12) then  $\rightarrow$  to (0,13) then  $\downarrow$  to Box (1,13)**



**Agent (1 ,13) move  $\downarrow$  to Box (2,13)**

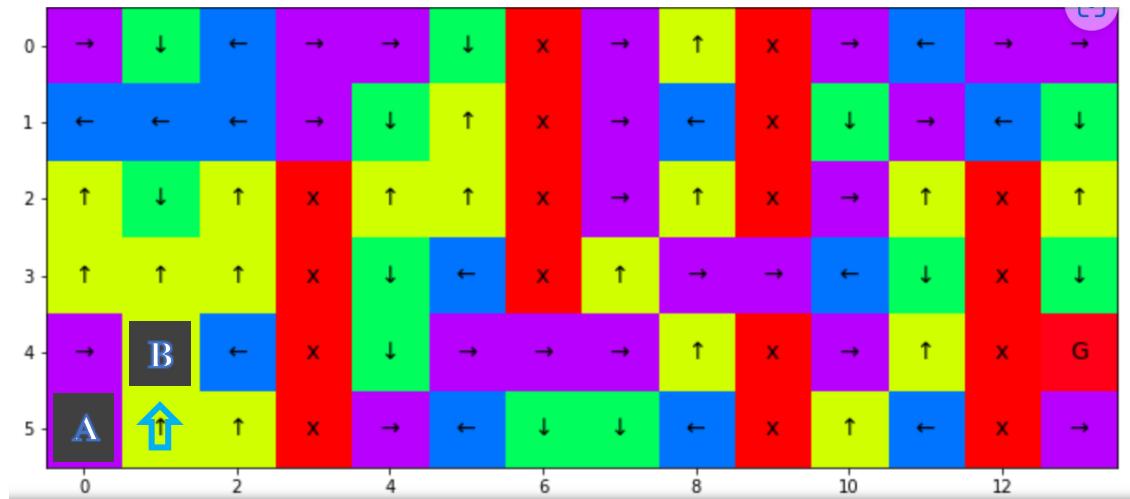


**Agent (2 ,13) move  $\downarrow$  to Box (3,13)**

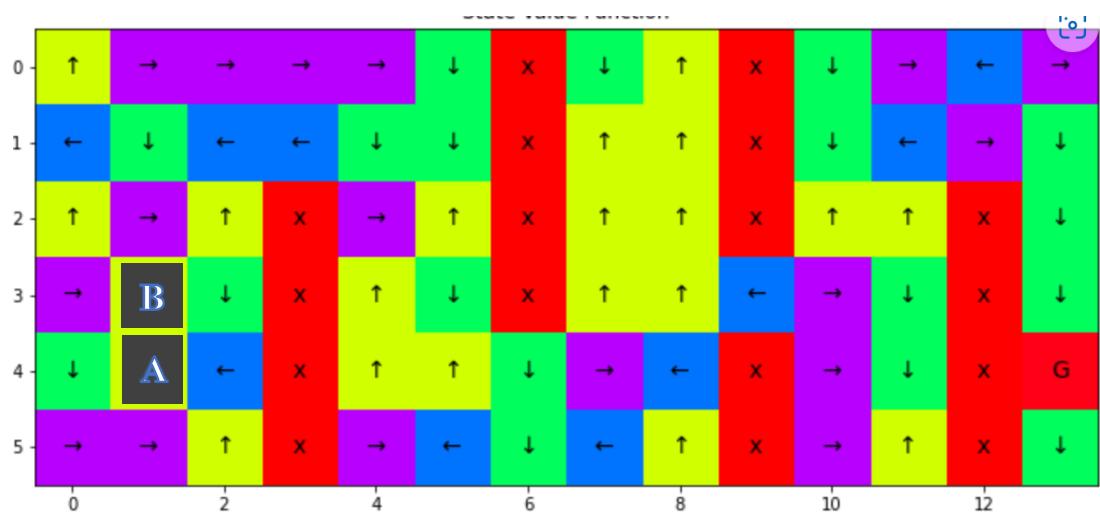


## ANNEX B SARSA (Step-by-step)

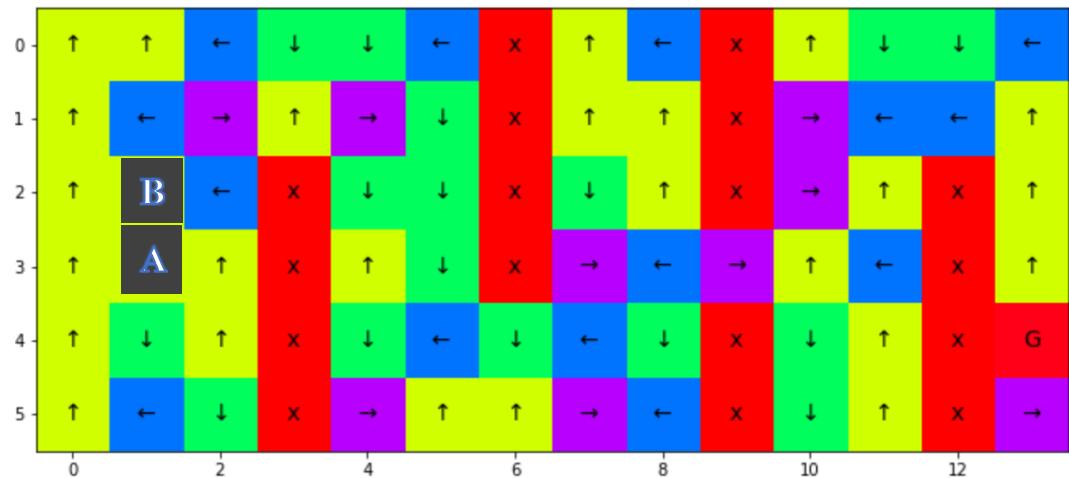
Agent (5,0) move → to (5,1) then move ↑ to Box (4,1)



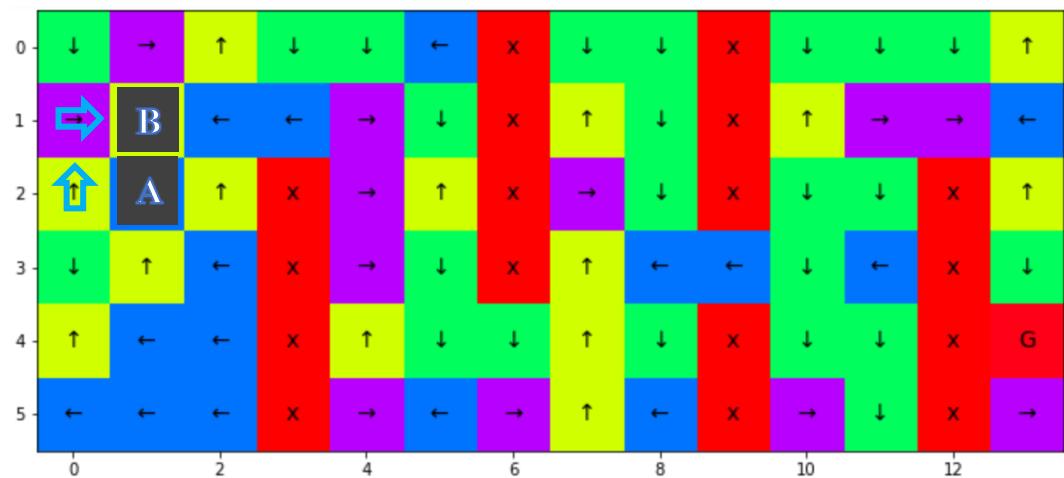
Agent (4,1) move ↑ to Box (3,1)



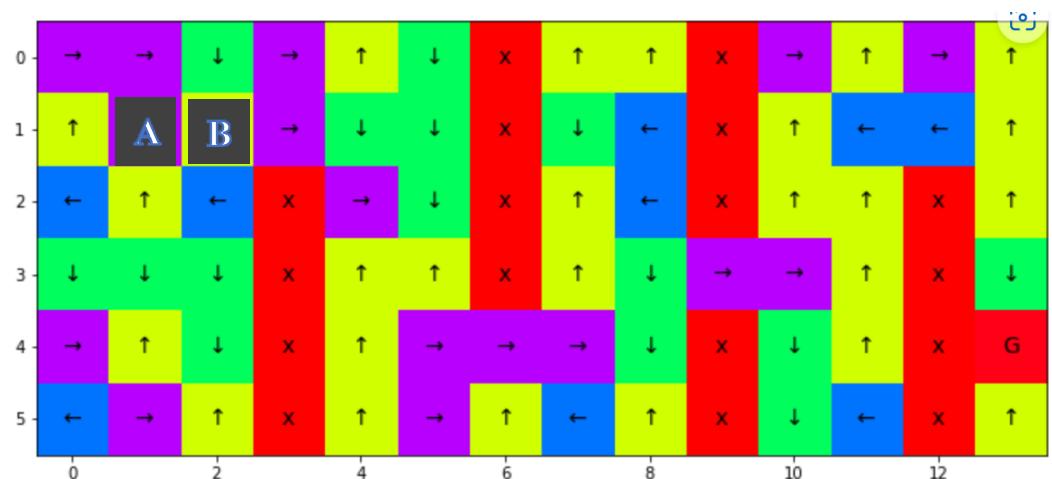
**Agent (3,1) move  $\uparrow$  to Box (2,1)**



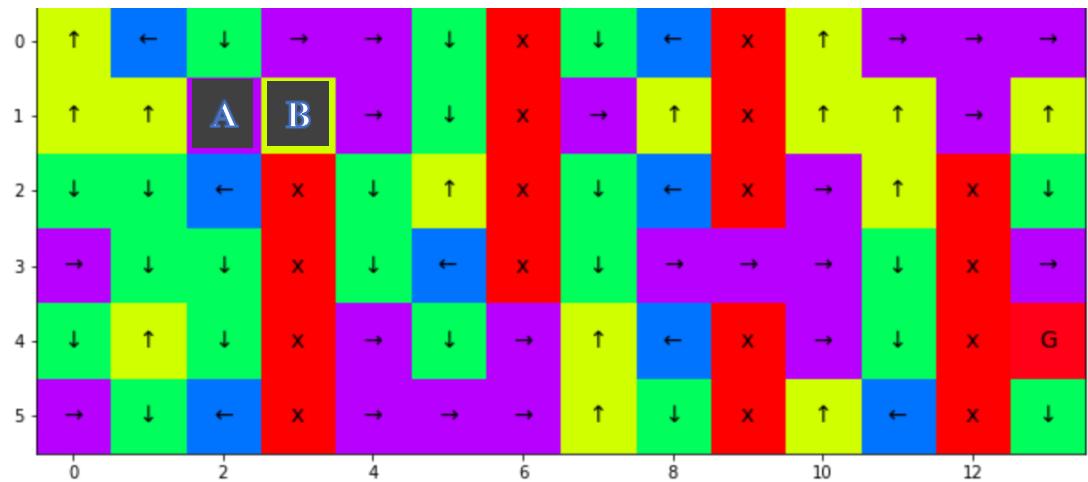
**Agent (2,1) move  $\leftarrow$  to (2,0) then  $\uparrow$  to (1,0) then move  $\rightarrow$  to Box (1,1)**



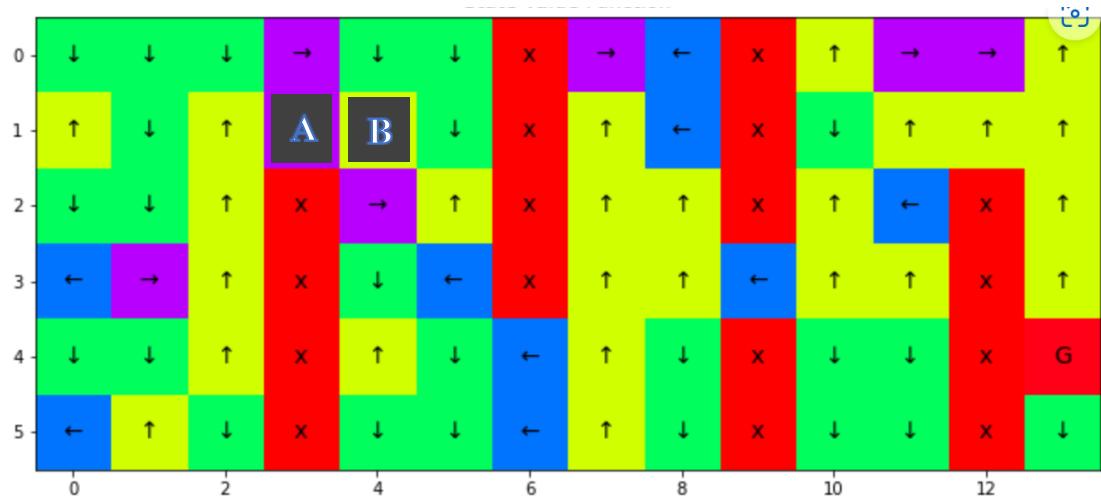
**Agent (1,1) move  $\rightarrow$  to Box (1,2)**



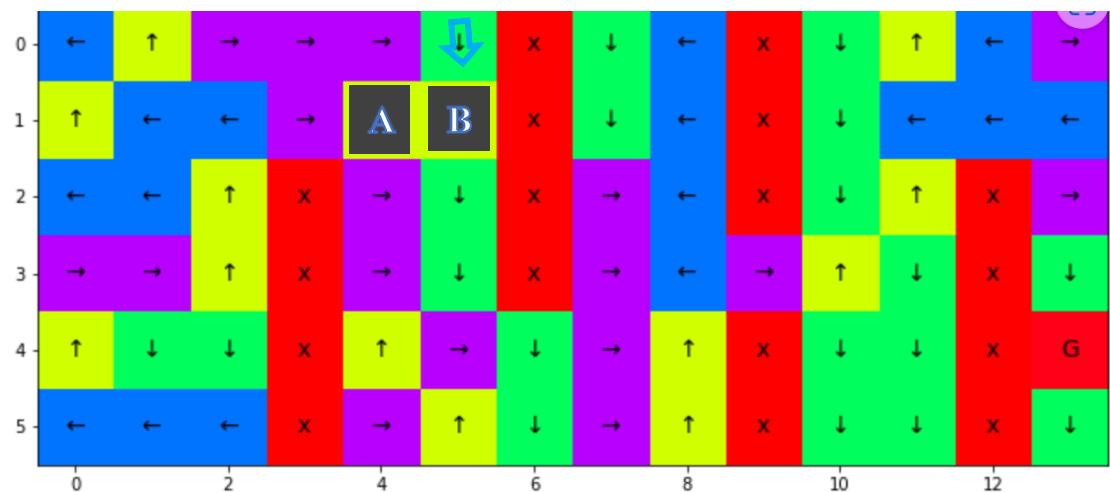
**Agent (1,2) move → to Box (1,3)**



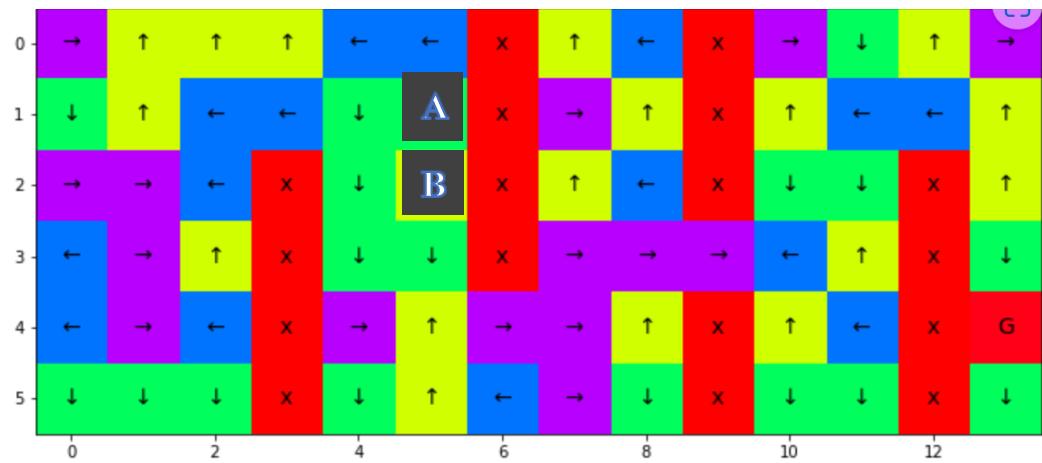
**Agent (1,3) move → to Box (1,4)**



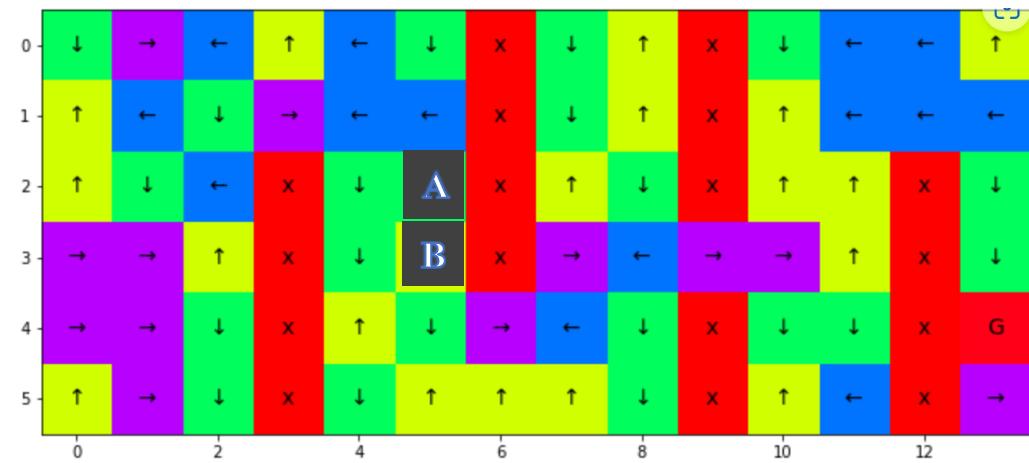
**Agent (1,4) move ↑ to (0,4) then move → to (0,5) then move ↓ to Box (1,5)**



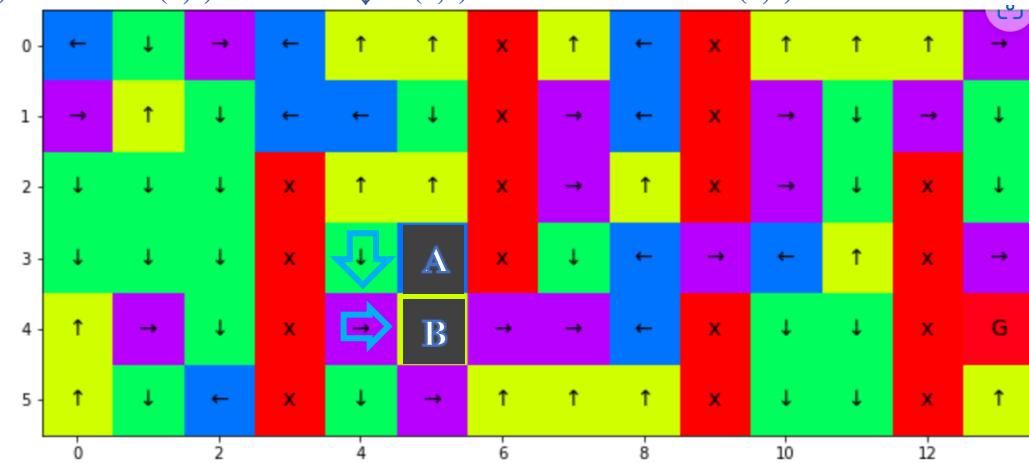
Agent (1,5) move ↓ to Box (2,5)



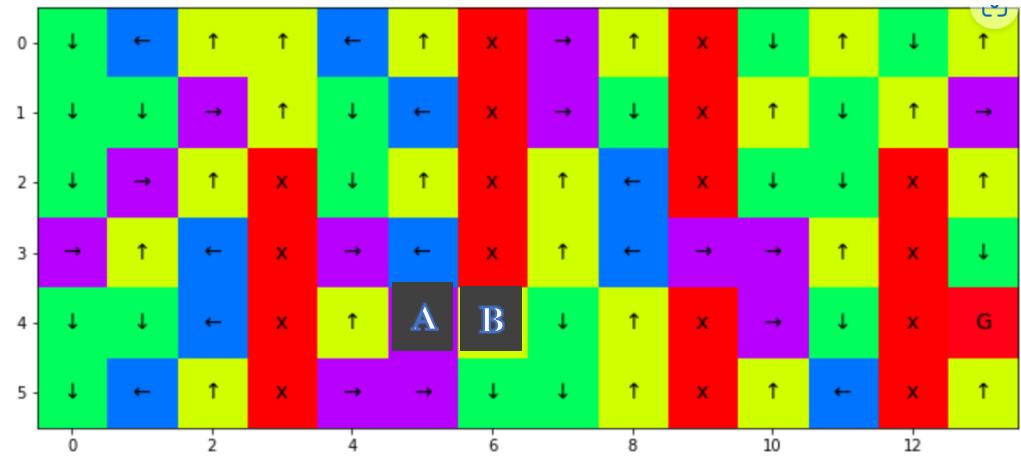
Agent (2,5) move ↓ to Box (3,5)



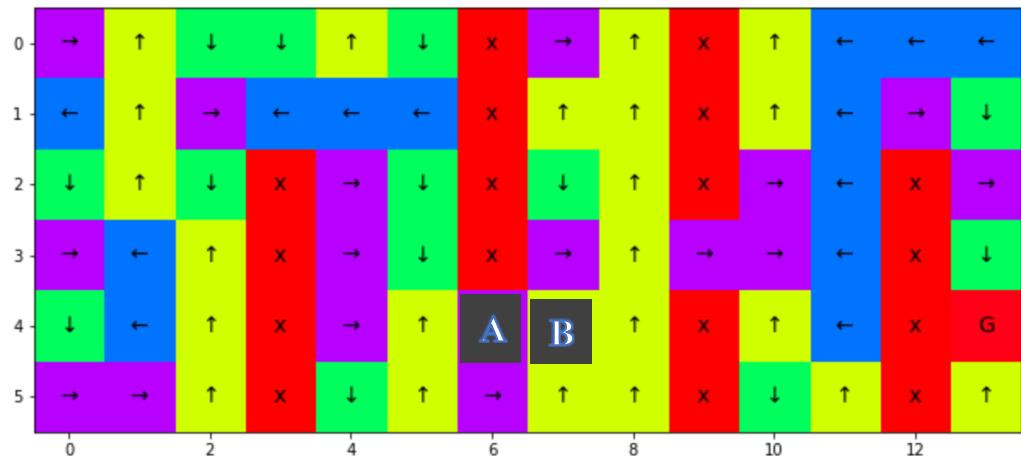
Agent (3,5) move ← to (3,4) then move ↓ to (4,4) then move → to Box (4,5)



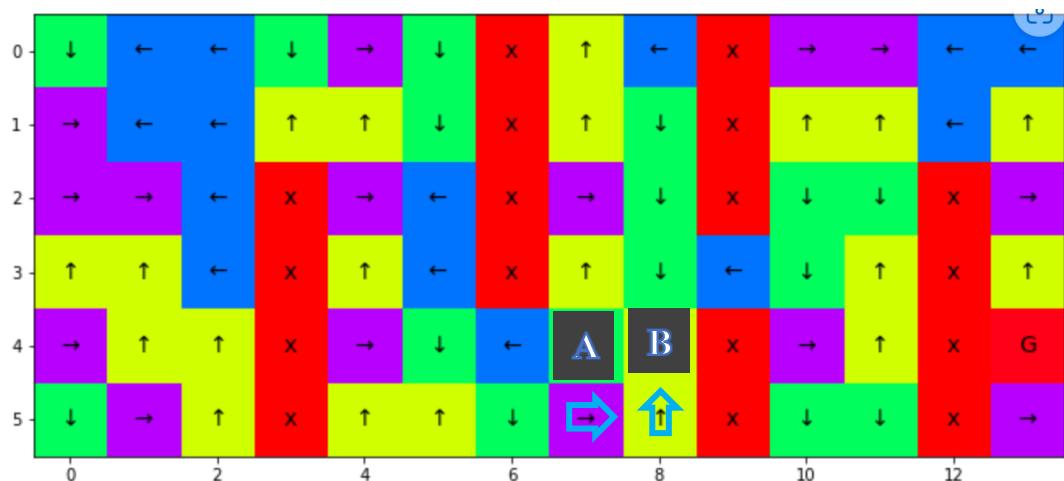
**Agent (4,5) move → to Box (4,6)**



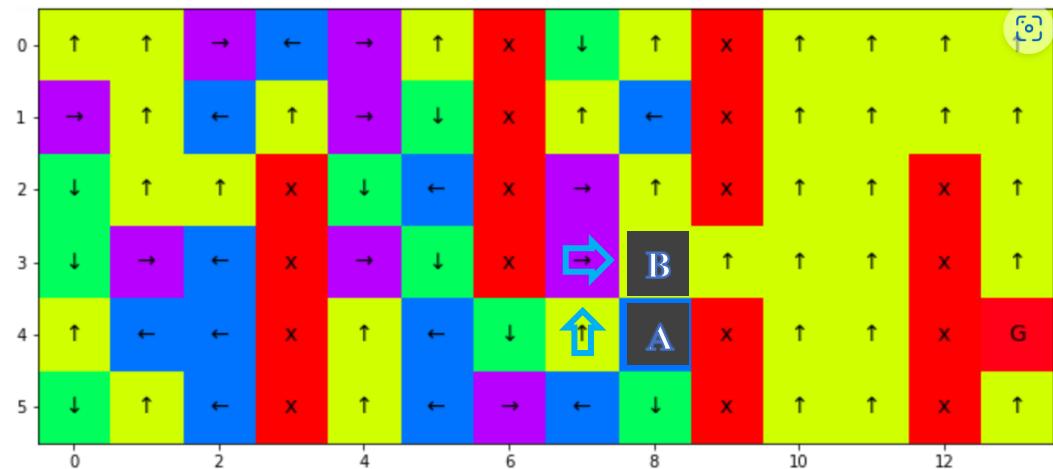
**Agent (4,6) move → to Box (4,7)**



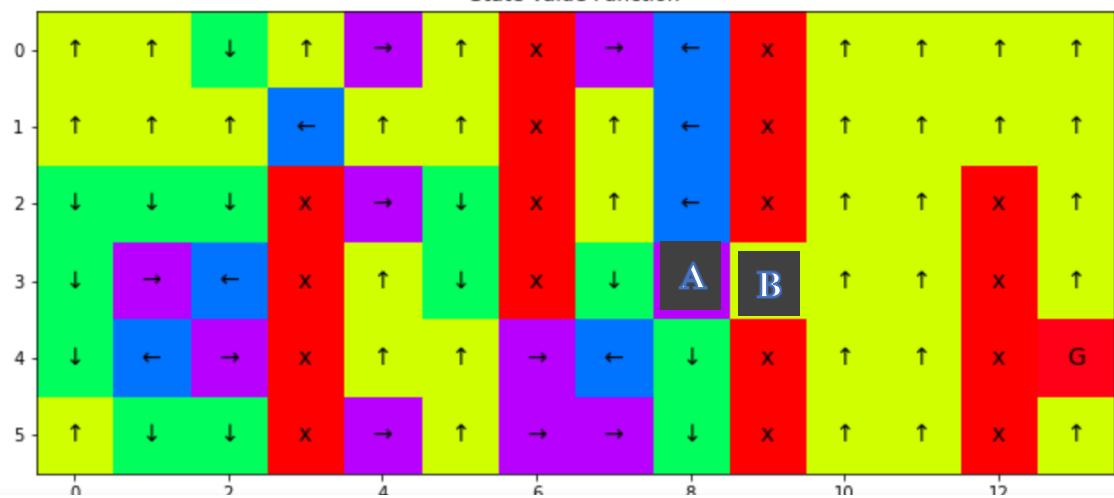
**Agent (4,7) move ↓ to (5,7) then move → to (5,8) the move ↑ to Box (4,8)**



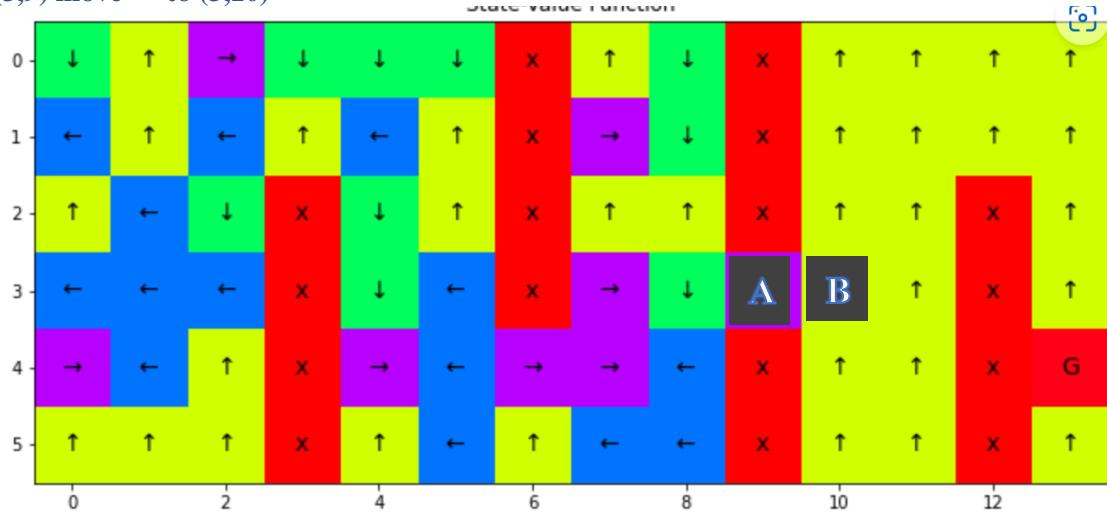
**Agent (4,8) move  $\leftarrow$  to (4,7) then move  $\uparrow$  to (3,7) then move  $\rightarrow$  to Box (3,8)**



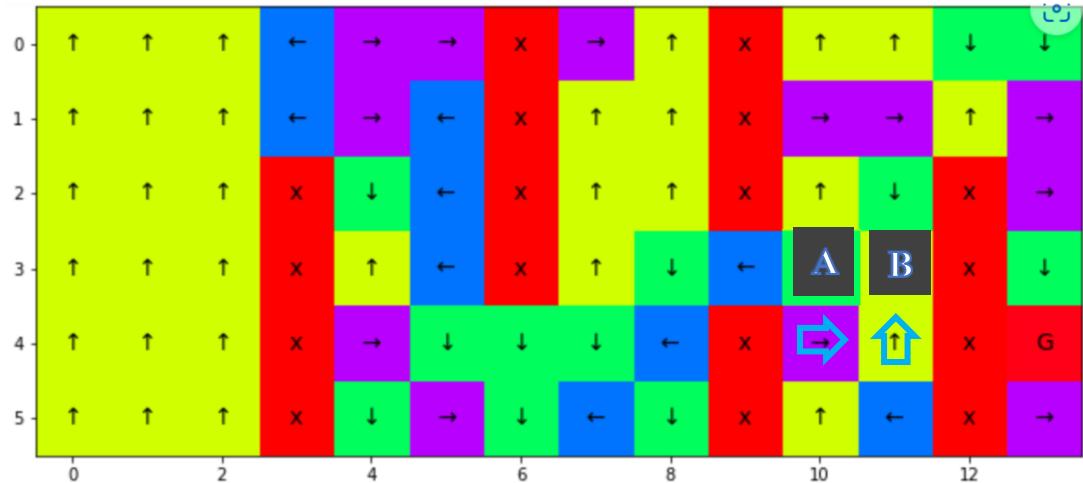
Agent (3,8) move → to (3,9)



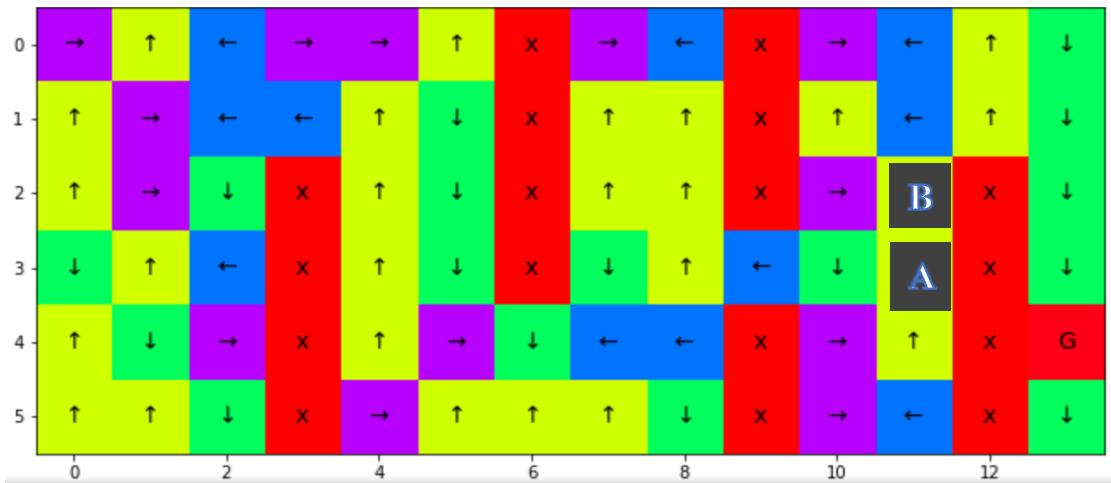
Agent (3,9) move → to (3,10)



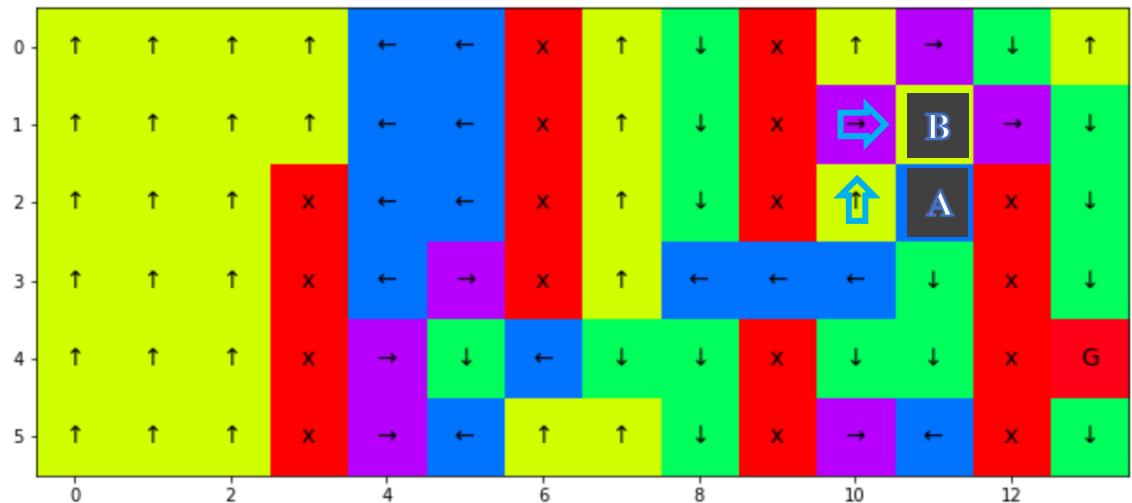
Agent (3,10) move ↓ to (4,7) then move → to (4,11) then move ↑ to Box (3,11)



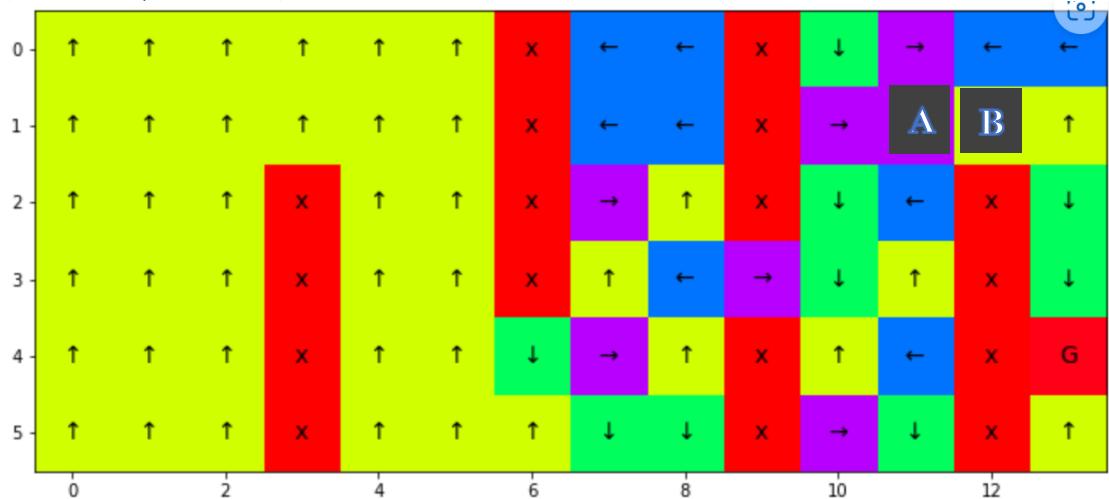
Agent (3,11) move ↑ to Box (2,11)



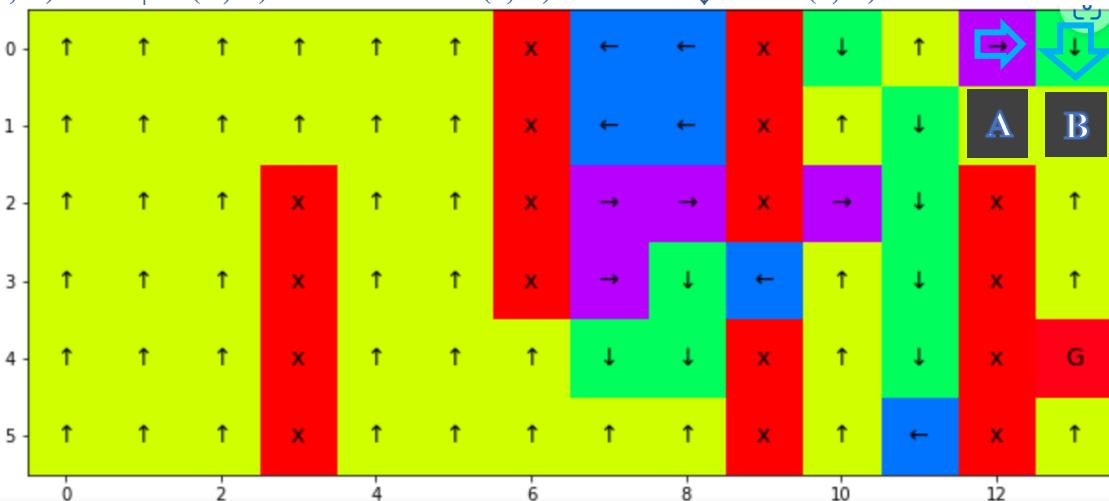
Agent (2,11) move ← to (2,10) then move ↑ to (1,10) then move → to Box (1,11)



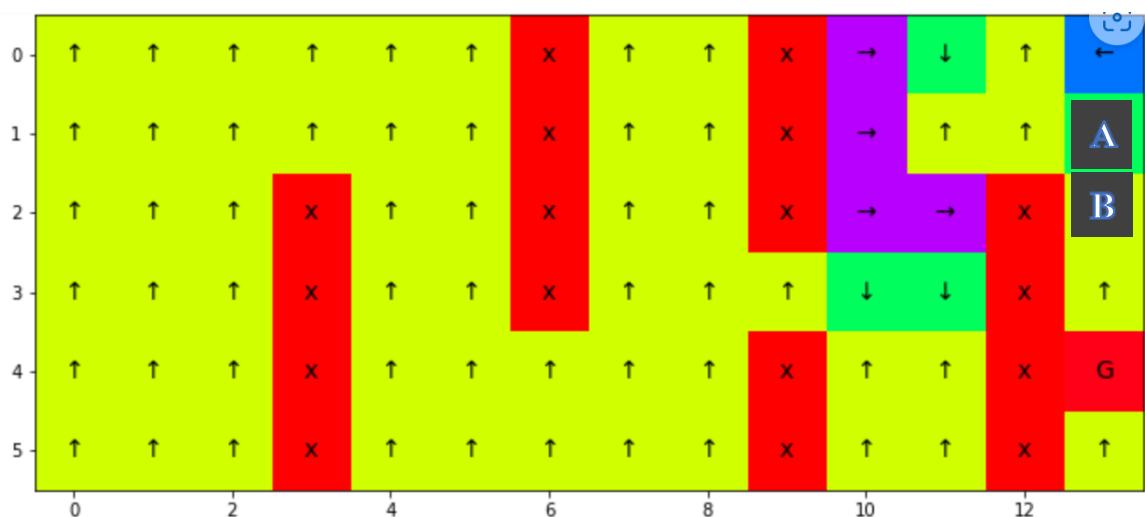
**Agent (1,11) move  $\uparrow$  to Box (1 ,12)**



**Agent (1,12) move  $\uparrow$  to (0 ,12) then move  $\rightarrow$  to (0,13) then move  $\downarrow$  to Box (1,13)**



**Agent (1,13) move  $\downarrow$  to Box (2,13)**



**Agent (2,13) move ↓ to Box (3,13)**

