



COL864: Special Topics in AI
Semester II, 2021-22

Reinforcement Learning II – Model-free Methods

Rohan Paul

Outline

- Last Class
 - Reinforcement Learning
 - Model-Based Reinforcement Learning
- This Class
 - Model-Free Reinforcement Learning
- Reference Material
 - Please follow the notes as the primary reference on this topic.

Acknowledgements

These slides are intended for teaching purposes only. Some material has been used/adapted from web sources and from slides by Nicholas Roy, Wolfram Burgard, Dieter Fox, Sebastian Thrun, Siddharth Srinivasa, Dan Klein, Pieter Abbeel, Max Likhachev and others.

Model-Based RL vs. Model-Free RL

- Model-Based RL
 - Used data to infer a model and compute a policy
 - Problem
 - Storing the model (and computing the policy) can be difficult.
- Can we compute a policy in a way that “avoids” storing the model?
 - Note that we still need to store “some” statistics over our experience
 - Hence, maintain an estimate of the “value” function.
- Model-Free RL
 - Bypass explicit learning of the intermediate model.
 - Can sample trajectories from the world directly and estimate a value function without the model.

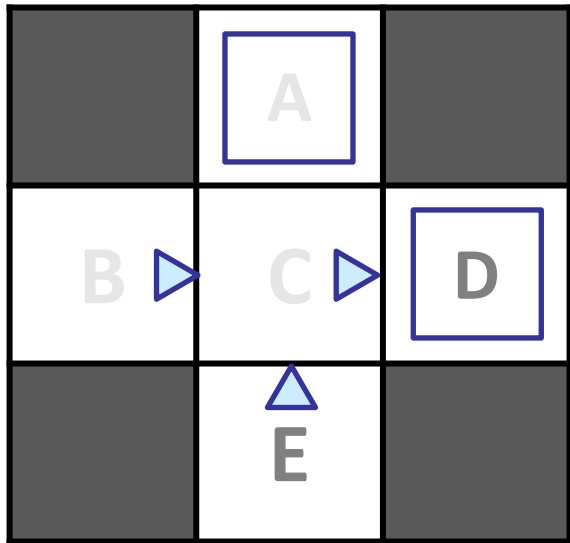
Monte Carlo Methods

- Learning the state-value function for a given policy
 - What is the value of a state?
 - Expected return – expected cumulative future discounted reward
- Key Idea
 - Sample trajectories from the world directly and estimate a value function without a model
 - Simply average the returns observed after visits to that state.
 - As more returns are observed the average should converge to the expected value.
 - Each occurrence of a state in an episode is called a visit to the state.

Toy Example: Monte Carlo Method

Input Policy π

Observed Episodes (Training)



Assume: $\gamma = 1$

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, x, +10

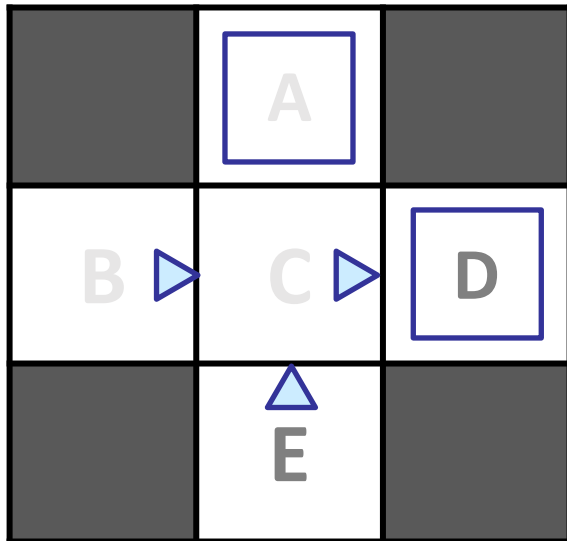
Episode 4

E, north, C, -1
C, east, A, -1
A, exit, x, -10

Value/utility of state c
 $V^\pi(C) = ((9 + 9 + 9 + (-11))/4)$
 $= 4$

Toy Example: Monte Carlo Method

Input Policy π



Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, x, +10

Episode 4

E, north, C, -1
C, east, A, -1
A, exit, x, -10

Value/utility of state c
 $V^\pi(C) = ((9 + 9 + 9 + (-11))/4)$
 $= 4$

Output Values

	-10	
+8	+4	+10
B	C	D
	-2	
	E	

Monte Carlo Methods

- Advantage
 - Do not require the MDP dynamics or rewards
- Disadvantage
 - Can only be applied to episodic MDPs
 - Averaging over the returns from a complete episode
 - Requires each episode to terminate

First-Visit Monte Carlo (FVMC) Policy Evaluation

```
First-visit MC prediction, for estimating  $V \approx v_\pi$ 

Input: a policy  $\pi$  to be evaluated
Initialize:
   $V(s) \in \mathbb{R}$ , arbitrarily, for all  $s \in \mathcal{S}$ 
   $Returns(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$ 

Loop forever (for each episode):
  Generate an episode following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 
   $G \leftarrow 0$ 
  Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :
     $G \leftarrow \gamma G + R_{t+1}$ 
    Unless  $S_t$  appears in  $S_0, S_1, \dots, S_{t-1}$ :
      Append  $G$  to  $Returns(S_t)$ 
       $V(S_t) \leftarrow \text{average}(Returns(S_t))$ 
```

Initialize the value function arbitrarily.

Loop over the episode from the end till the start.

Account for the discounting of the reward

Except if (unless) the state has been visited from time 0 till (t-1) append and average out the results.

I.e., only update the value estimate if this is the first visit to the state.

First-Visit Monte Carlo (FVMC)

- First-Visit Monte Carlo (FVMC)
 - Averages the returns following the first visit to a state s in the episode.
- Every-visit Monte Carlo (EVMC)
 - Averages returns following all the visits to s .
- Convergence
 - FVMC – error falls as $1/N(s)$. Needs lots of data
 - EVMC - error falls quadratically, slightly better data efficiency.

So, now we know how to evaluate a policy.
What about improving the policy now?

Policy Improvement

Two problems in the last pseudo-code in doing policy improvement.

1. How to obtain the policy?

- Note: we only stored the state value function. In the absence of a model, we cannot compute the policy.
- *Solution: Store the state-action values.*

2. How to ensure the coverage of states?

- *Solution: Use epsilon-greedy policies*
 - Most of the time select an action that has maximal estimated action value.
 - But with probability epsilon, instead, select an action at random.

Policy Improvement (Monte Carlo Control)

Estimate the Q-function.
Derive the policy from there.

All non-greedy actions have a small probability of being selected. The bulk of the likelihood is given to the greedy action.

```
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  $Returns(S_t, A_t)$ 
       $Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$ 
       $A^* \leftarrow \text{argmax}_a Q(S_t, a)$  (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 literature, a policy is called *soft* if $\pi(a|s) > 0$ for all $s \in \mathcal{S}$ and all $a \in \mathcal{A}(s)$.

An ϵ -*soft* policy is one for which $\pi(a|s) \geq \frac{\epsilon}{|\mathcal{A}(s)|}$ for all states and actions

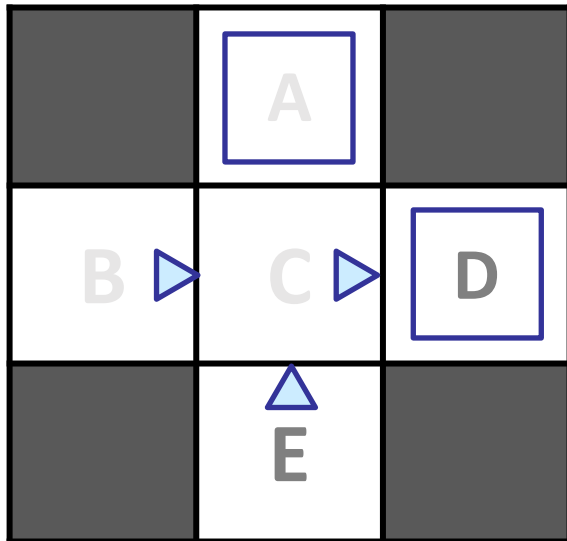
for some ϵ

Problem with Monte Carlo

- Limitations
 - Each state must be learned separately, loses the state connection information.
 - Estimate of one state is not taking advantage of the estimates of the other states.
 - Note: Bellman equations tell us that value function for states has a recursive relationship.
 - Could only be used in an episodic setting.

Recall our example

Input Policy π



Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, x, +10

Episode 4

E, north, C, -1
C, east, A, -1
A, exit, x, -10

Output Values

	-10	
	A	
+8	+4	+10
B	C	D
	-2	
	E	

Problem: we have lost the connection between states.
We go through state C, 4 times.
We use only 2 estimates for E.
But, C and E are adjacent!

Temporal Difference (TD) Learning

- Model-Free combination of
 - Monte Carlo (learning from sample trajectories/experience) and
 - Dynamic programming (via Bellman Equations)
- Incorporate Bootstrapping
 - Update value function estimates of a state based on others
 - Adjust the value function estimate using the Bellman Equation relationship between the value function of successor states.
 - More data-efficient than a Monte Carlo method (discussed previously)
- Setting
 - Can be used in an *episodic* or *infinite-horizon non-episodic* settings
 - Immediately updates the estimate of $V(s)$ after each (s, a, s', r) tuple.

“If one had to identify one idea that is central and novel to reinforcement learning, it would undoubtedly be temporal-difference (TD) learning”, Sutton and Barto 2017

TD Learning

- Aim: estimate $V^\pi(s)$ given episodes generated under policy π
 - $G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots$ in MDP M under policy π
 - $V^\pi(s) = E_\pi[G_t]$
- Bellman Operator (if we know MDP models)

$$B^\pi V(s) = r(s, \pi(s)) + \gamma \sum_{s'} p(s'|s, \pi(s)) V(s')$$

- In incremental every-visit MC, update estimate using 1 sample of return (for the current i^{th} episode)

$$V^\pi(s) = V^\pi(s) + \alpha(G_{i,t} - V^\pi(s))$$

- Insight: have an estimate of V^π , use to estimate expected return

$$V^\pi(s_t) = V^\pi(s_t) + \alpha \underbrace{([r_t + \gamma V^\pi(s_{t+1})])}_{\text{TD target}} - V^\pi(s_t)$$

- Can immediately update the value estimate after each (s, a, r, s') tuple. Do not need the episodic setting.

TD Methods

- The updates are based on the difference in value functions at each time step, the TD error,

$$\delta_t = r_t + \gamma V^\pi(s_{t+1}) - V^\pi(s_t)$$

hence the name *temporal difference learning*.

- α is the learning rate.
- TD can be generalised to n -step returns:

$$R_t^{(n)} = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots + \gamma^n V_t(s_{t+n})$$

Note: earlier we used the 1-step return in the TD update. Now, we can generalize to how many steps in the future to update from.

$$V^\pi(s_t) = V^\pi(s_t) + \alpha \underbrace{([r_t + \gamma V^\pi(s_{t+1})])}_{\text{TD target}} - V^\pi(s_t)$$

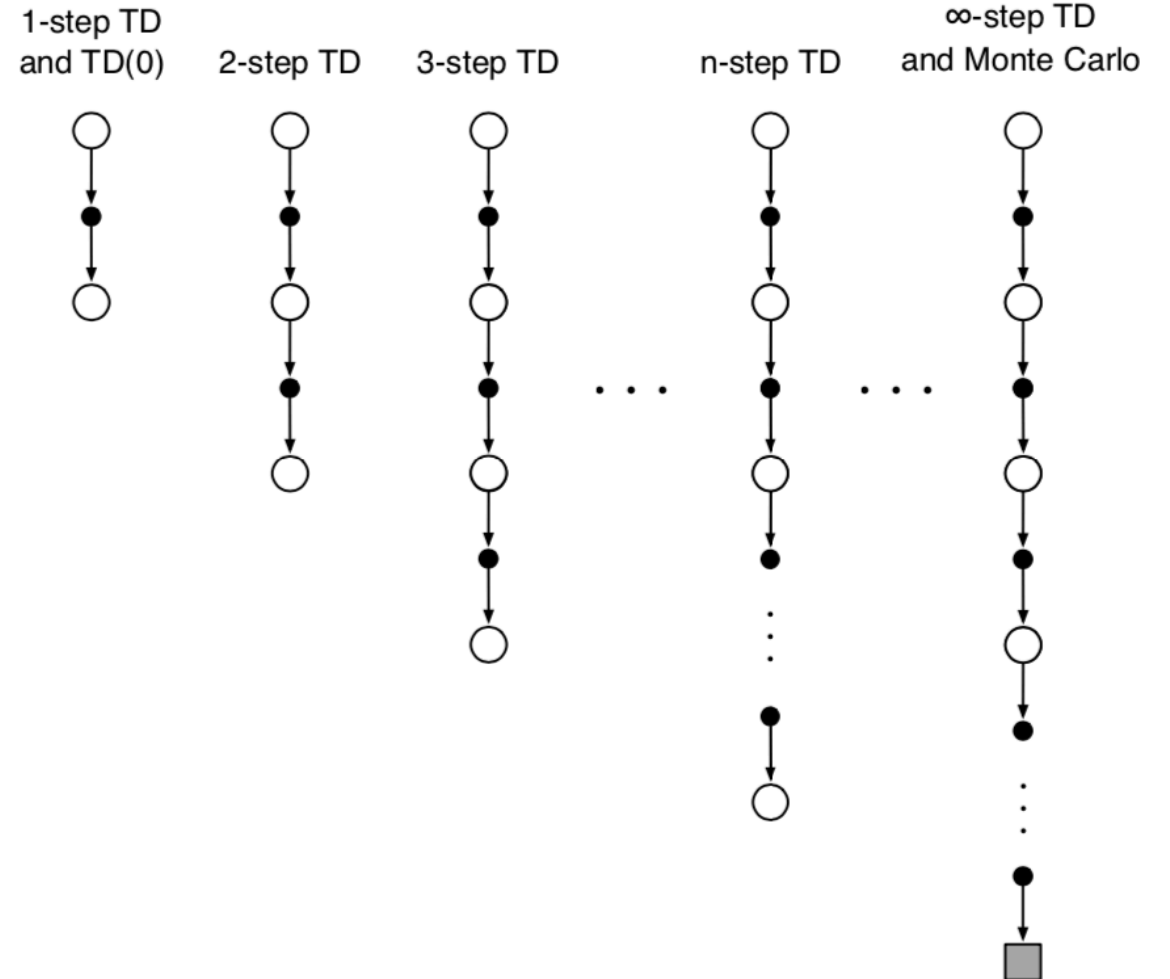
TD Methods

TD can be generalised to n -step returns:

$$R_t^{(n)} = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots + \gamma^n V_t(s_{t+n})$$

TD Learning (intuitively)

- Nudge our prior estimate of the value function for a state using the given experience.
- Shift the estimate based on the error in what we are experiencing and what estimate we had before
- Weighted by the learning rate.



Temporal Difference Learning: Example

States

	A	
B	C	D
	E	

Observed Transitions

B, east, C, -2

	0	
0	0	8
	0	

C, east, D, -2

	0	
-1	0	8
	0	

	0	
-1	3	8
	0	

Assume: $\gamma = 1$, $\alpha = 1/2$

$$V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + \alpha [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

SARSA

- The TD algorithm describes how to evaluate a policy, but does not describe how to improve the policy.
- Can be done very easily with SARSA.
- Rather than learning the value function $V(s)$, we will learn the state-action value $Q(s, a)$, learning from the tuple $\{s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}\}$ hence the name SARSA or “state action reward state action”.
- If we generalize the TD rule to state-action functions, such that

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha [(r_{t+1} + \gamma Q(s_{t+1}, a_{t+1})) - Q(s_t, a_t)],$$

can then choose a new policy by maximizing over the next action we might take, as in

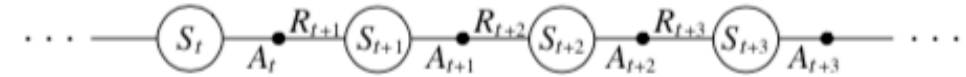
$$\pi(s) = \max_a Q(s, a)$$

- We modify this policy slightly, using a policy $\pi_\epsilon(s)$ such that

$$\pi_\epsilon(s) = \begin{cases} \max_a Q(s, a) & \text{with probability}(1 - \epsilon) \\ \text{RAND}(A) & \text{otherwise} \end{cases}$$

in order to ensure that we visit all state-action pairs sufficiently often.

- The parameters ϵ can be lowered as $\epsilon_t = 1/t$ so that the agent eventually stops exploring and converges to the optimal policy.



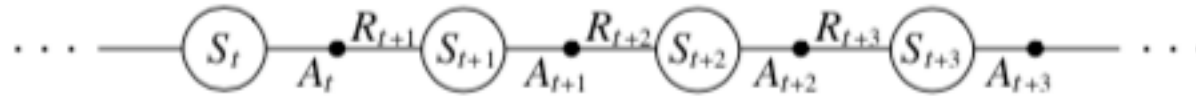
Update using {state, action, reward, state, action}.

SARSA is an on-policy algorithm

- Updates are using the transitions actually taken by the agent

- Note: that policy may or may not be good.

SARSA



Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

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

Loop for each episode:

Initialize S

Choose A from S using policy derived from Q (e.g., ε -greedy)

Loop for each step of episode:

Take action A , observe R, S'

Choose A' from S' using policy derived from Q (e.g., ε -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

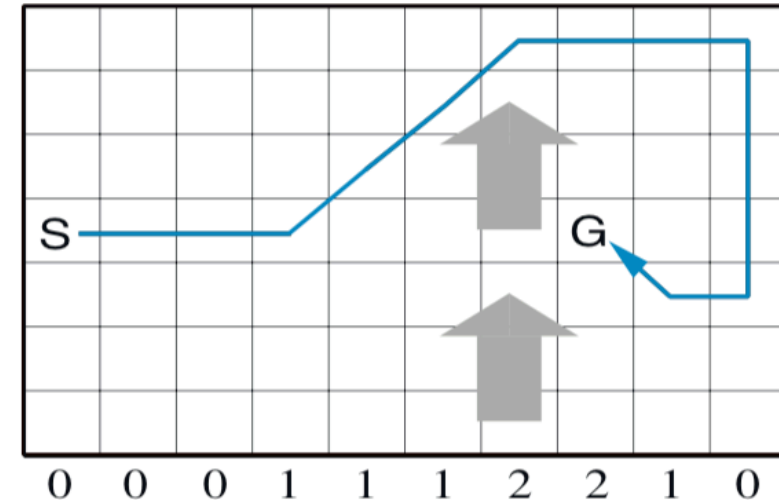
$S \leftarrow S'; A \leftarrow A';$

until S is terminal

State-action values
updated using the SARSA
Learning Rule

Example: Windy Grid World

- Setup
 - Standard grid world with start and goal state.
 - Crosswind running upward through the middle of the grid.
 - Actions: Up, down, left, right
 - Wind strength varies from column to column (written below) in number of grid cells shifted upwards.
 - E.g., if you are one cell to the right of the goal, then the action left takes you to the cell just above the goal.
 - Undiscounted constant rewards of -1 till the goal state is reached.
- SARSA with epsilon-greedy does well on this control task.
 - Can learn that the blue path is good.
- Monte Carlo methods cannot easily be used for this task because termination is not guaranteed for all policies.
 - If a policy was found that caused the agent to stay in the same state, then the next episode would never end.
- Step-by-step learning methods such as SARSA do not have this problem because they quickly learn *during the episode*.



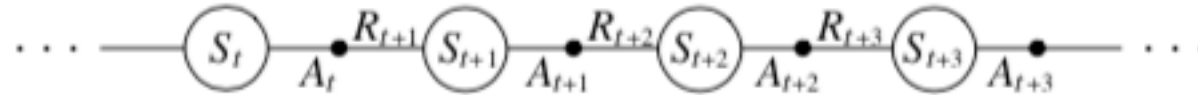
Q-Learning

- One disadvantage to SARSA is that it is an on-policy algorithm
 - That is, we only get an estimate of the Q function for tuples that are directly experienced by executing the policy
 - What if our initial policy is really poor?
 - In that case could take a long time to find the optimal policy.
- The update rule can be modified to improve the *best* state-action tuple, rather than the *experienced* tuple by putting the max operator directly inside the update rule, as in

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha \left[(r_t + \gamma \max_{a'} Q(s_{t+1}, a')) - Q(s_t, a_t) \right].$$

- Q-LEARNING is one of the core algorithms of reinforcement learning.
- It is an off-policy algorithm, in that the Q function can be shown to converge, regardless of the underlying policy, so long as the underlying policy is guaranteed to visit all state-action pairs infinitely often.

Q-Learning



Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

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

Loop for each episode:

Initialize S

Loop for each step of episode:

Choose A from S using policy derived from Q (e.g., ε -greedy)

Take action A , observe R, S'

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

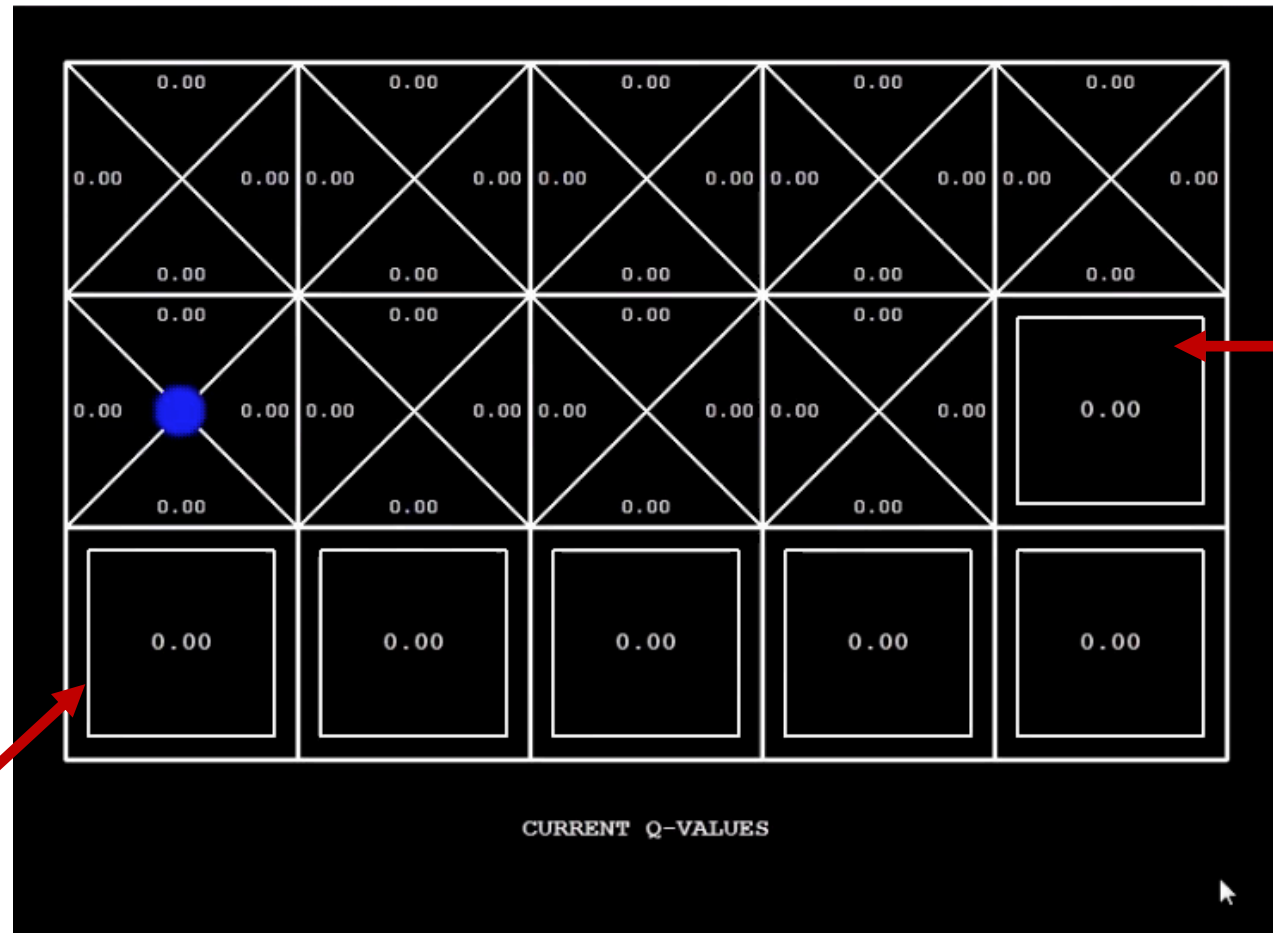
$S \leftarrow S'$

until S is terminal

Q-Learning update rule that uses the max of the next state-action tuples to update.

Example: Cliff Walking

- The agent does not know the rewards a-priori.
 - Learns the effect of the east action over time.
 - Only the actions taken by the agent contribute to updates.
- Occasionally falls in the cliff and gets the negative reward.
- Note that the max of the Q values is propagated (green values) to other states as it is approximating the optimal Q value.



Cliff at the bottom,
negative reward here.

Exit with
+10 reward

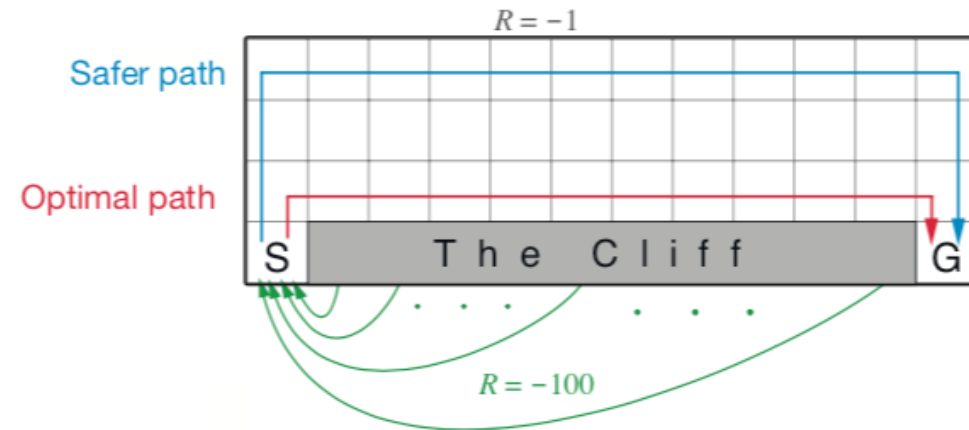
The Off-Policy Nature of Q-Learning

- Off-policy methods
 - Learns the optimal state-action value function, independent of the policy being followed.
- Q-learning converges to optimal policy
 - Even if the agent is acting *sub-optimally*.
 - Under certain conditions such as when exploration is enough, the learning rate becomes small enough.
- The algorithm learns about the value of the optimal policy without knowing or using the optimal policy.
 - Example: sporadic falls in the cliff does not affect the knowledge that going right leads to high reward.
- Note that SARSA was an on-policy algorithm

Example: Another Cliff Walking Domain

Task: Undiscounted, episodic task with start and goal states. Reward of (-1) for living. Falling in the cliff leads to a reward of (-100) and sends back to the start state.

- Q-learning learns to walk along the cliff (red path) occasionally falling into the cliff due to epsilon greedy action selection and gets lower rewards.
- SARSA takes the action selection into account. It leans the longer but safe path (blue) through the upper part of the grid.
- Although Q-learning actually learns the values of the optimal policy, its online performance is worse in comparison to SARSA.
- If epsilon is reduced then both methods perform similarly.



Expected SARSA

- Expected SARSA just like Q-learning except takes the “expectation” instead of the max.
- Considers how likely the action is under the current policy.
- Expected SARSA is more complex computationally but eliminates the variance due to the random selection of actions.
- Expected SARSA performs better than both SARSA and Q-learning on the cliff walking task.

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \mathbb{E}_{\pi} [Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_t, A_t) \right]$$
$$\leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \sum_a \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right],$$

