

## Outline and References

- Outline
  - ▶ TD(0)
  - ▶ SARSA
  - ▶ Q-Learning

## References

- ▶ Reinforcement Learning: An Introduction [RL Chapter 6 and 7]
- ▶ Sample-based Learning Methods (Coursera)



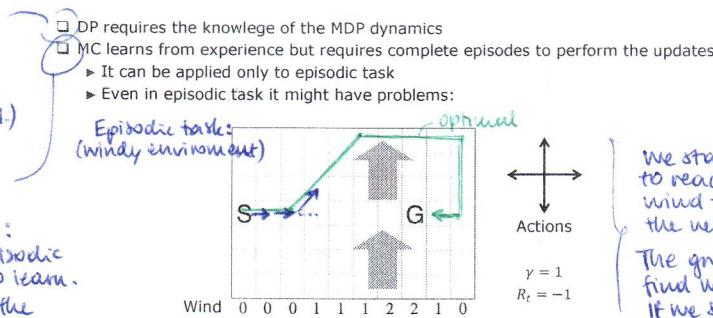
Machine Learning - Daniele Loiacono

2

## Why temporal-difference?

Dynamic Programming was not enough because it required us to know the one step dynamics (and we typically don't have this knowl.)

Monte Carlo can learn from experience but there are some issues/limitations:  
It can be applied only to episodic task  $\Rightarrow$  it needs episodes to learn.  
The point is that MC needs the whole episode to START learning.



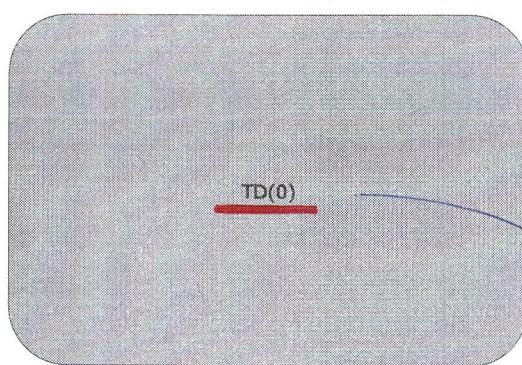
We start from the start  $S$  and we have to reach the goal  $G$ . There is some wind that changes the transition to the next state.

The green policy is extremely difficult to find with an explorative policy, for example. If we start with the uniform random actions ( $\uparrow 25\%$ ,  $\downarrow 25\%$ ,  $\rightarrow 25\%$ ,  $\leftarrow 25\%$ )

we cannot really reach the goal by chance easily. Given that, in order to START MonteCarlo, we need to reach the goal (so that the episode ends, we compute the returns and we start updating), this becomes a big issue.

Machine Learning - Daniele Loiacono

3



It combines together MonteCarlo AND Dynamic Programming (it combines: SAMPLING (from MC): we don't need the knowledge of the dynamics, we need experience, and BOOTSTRAPPING: we exploit the recursive structures of the problem to learn/update the value of the value function in a state  $S$  based on the values of the value function in the states that are next to  $S$ )

Machine Learning - Daniele Loiacono

4

## Temporal-Difference Policy Evaluation - TD(0)

- Temporal-Difference combines MC (model-free) with DP (bootstrapping):

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

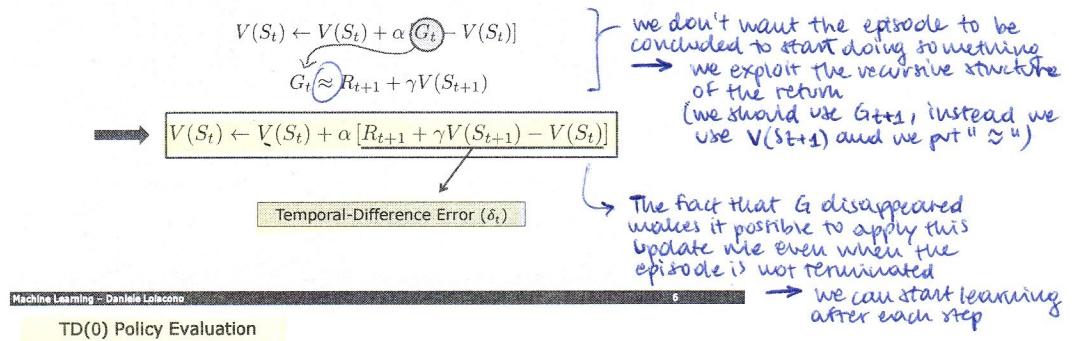
To update this, we need the return  $G_t$  (which means that the episode must be completed)

We start from this: Monte Carlo policy evaluation (version which uses rolling means instead of incremental mean)

Machine Learning - Daniele Loiacono

5

- Temporal-Difference combines MC (model-free) with DP (bootstrapping):



```

Input: the policy  $\pi$  to be evaluated
Algorithm parameter: step size  $\alpha \in (0, 1]$ 
Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$ 
Loop for each episode:
    Initialize  $S$ 
    Loop for each step of episode:
         $A \leftarrow$  action given by  $\pi$  for  $S$ 
        Take action  $A$ , observe  $R, S'$ 
         $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$ 
         $S \leftarrow S'$ 
    until  $S$  is terminal
  
```

Machine Learning - Daniela Loloceno

7

TD(0) Policy Evaluation

```

Input: the policy  $\pi$  to be evaluated
Algorithm parameter: step size  $\alpha \in (0, 1]$ 
Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$ 
Loop for each episode:
    Initialize  $S$ 
    Loop for each step of episode:
         $A \leftarrow$  action given by  $\pi$  for  $S$ 
        Take action  $A$ , observe  $R, S'$ 
         $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$ 
         $S \leftarrow S'$ 
    until  $S$  is terminal
  
```

Machine Learning - Daniela Loloceno

8

TD(0) Policy Evaluation

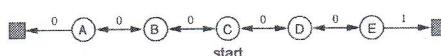
```

Input: the policy  $\pi$  to be evaluated
Algorithm parameter: step size  $\alpha \in (0, 1]$ 
Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$ 
Loop for each episode:
    Initialize  $S$ 
    Loop for each step of episode:
         $A \leftarrow$  action given by  $\pi$  for  $S$ 
        Take action  $A$ , observe  $R, S'$ 
         $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$ 
         $S \leftarrow S'$ 
    until  $S$  is terminal
  
```

the update is still delayed in time, but now is not delayed at the end of the episode but just of 1 step (HUGE advantage)

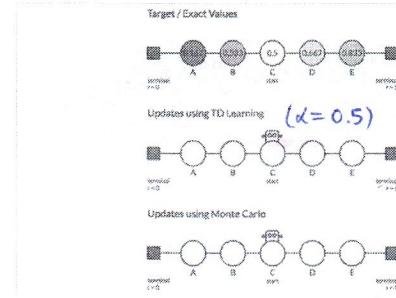
Machine Learning - Daniela Loloceno

Random Walk Example – Comparison between MC and TD(0)



- Let  $C$  be the starting state
- Let  $\pi$  be left or right with equal probability in all states (random policy)
- Reward +1 on right termination, 0 otherwise
- Assuming  $\gamma = 1$ ,  $V_\pi(s)$  represents the probability of ending on the right side

## Random Walk Example - Simulation

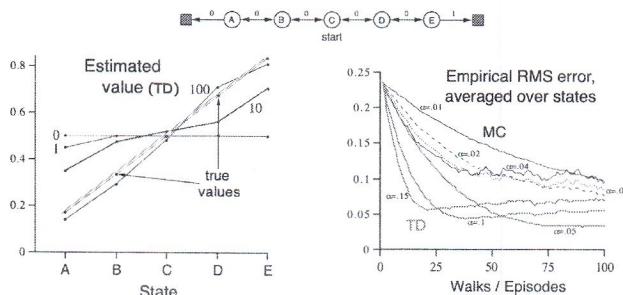


Shown in animation is that in both cases we converge to the target, however the two methods behave very differently when updating. Typically, TD is much faster than MC.

## Machine Learning - Daniele Lolli

11

### Random Walk Example – Value Function



## Machine Learning - Daniele Lolli

12

### MC vs TD

- TD can learn before knowing the final outcome
  - TD can learn after every step
  - MC must wait until end of episode before return is known
- TD can learn without the final outcome
  - TD can learn from incomplete sequences
  - MC can only learn from complete sequences
- TD works in continuing (non-terminating) environments MC only works for episodic (terminating) environments

## Machine Learning - Daniele Lolli

13

### MC vs TD: Bias-Variance

- MC target has lower bias
  - MC return  $R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_{t+T}$  is an **unbiased estimate** of  $V_\pi(S_t)$
  - TD target  $R_{t+1} + \gamma V(S_{t+1})$  is a **biased estimate** of  $V_\pi(S_t)$ , as  $V(S_{t+1}) \neq V_\pi(S_{t+1})$
- TD target has lower variance
  - Return depends on many random actions, transitions, rewards
  - TD target depends on one random action, transition, reward
- Overall
  - MC works well with function approximation and it is not very sensitive to initial values
  - TD has problem with function approximation and it is more sensitive to initial values

In MC we have multiple sources of stochasticity (since we consider  $R_{t+1}, R_{t+2}, \dots, R_{t+T}$ ), instead, in TD we only have 2 sources of stochasticity ( $R_{t+1}$  and  $S_{t+1}$ ).

We've used to represent the value function as a table:

$$V = \begin{matrix} s_0 & s_1 & & & s_N \\ | & | & | & \cdots & | \\ V(s_0) & & & & \end{matrix}$$

This is not feasible in a huge problem → in real problem we replace the tabular representation with a parametrized function (that can be anything: SVM, regression,..)

## Machine Learning - Daniele Lolli

14

### Sampling and Bootstrapping: overview

- Bootstrapping: update involves an estimate
  - MC does not bootstrap
  - DP bootstraps
  - TD bootstraps
- Sampling: update does not involve an expected value
  - MC samples
  - DP does not sample
  - TD samples

## Machine Learning - Daniele Lolli

15



Let's go back to MC Policy Iteration

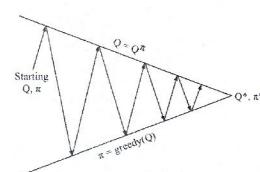
□ Evaluation

What we do in SARSA is to replace the return with an approx.:  $G_t \approx R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$

$$Q_\pi(s, a) \approx \text{average}[G_t | S_t = s, A_t = a]$$

□ Improvement

$$\pi'(s) = \arg \max_a Q_\pi(s, a)$$



On-Policy Control with SARSA

because we're learning all action value function by following a policy; an the action value function we're learning is the one of the policy we're following

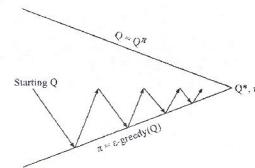
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

□ Evaluation

► SARSA

□ Improvement

►  $\varepsilon$ -Greedy Policy Improvement



On-Policy Control with SARSA

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.,  $\varepsilon$ -greedy)  
   Loop for each step of episode:  
     Take action  $A$ , observe  $R, S'$   
     Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)  

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

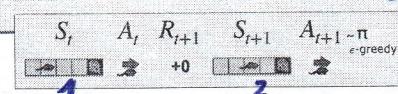
$$S \leftarrow S'; A \leftarrow A';$$
  
   until  $S$  is terminal

On-Policy Control with SARSA

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.,  $\varepsilon$ -greedy)  
   Loop for each step of episode:  
     Take action  $A$ , observe  $R, S'$   
     Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -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



the mouse want to reach the cheese, he goes right, he doesn't get the reward, he goes again right (accordingly with the policy that he follows, that is an  $\varepsilon$ -greedy policy based on the current known of the problem) and then we use the new information to update the previous pair. The information of 2 update the information of 1.

\* From here we get the name:

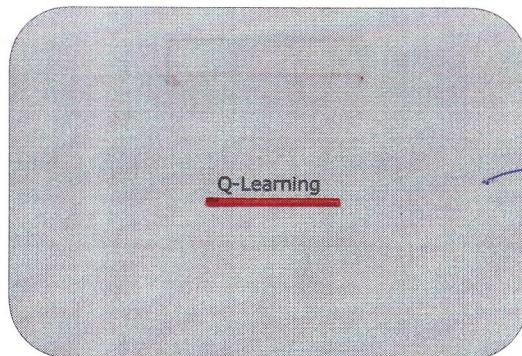
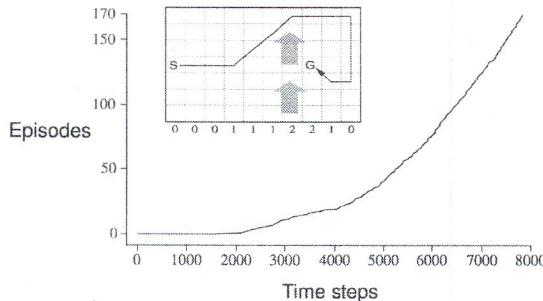
$S_t$  = current state

$A_t$  = current action

$R_t$  = current reward (immediate)

$S_{t+1}$  = next state

$A_{t+1}$  = next action



## Q-Learning: Off-Policy TD Control

- SARSA, as Policy Iteration in DP, is based on Bellman Expectation Equation

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

$$Q_\pi(s, a) = \sum_{s', r} p(s', r|s, a) \left( r + \gamma \sum_{a'} \pi(a'|s') Q_\pi(s', a') \right)$$

- Q-Learning, as Value Iteration in DP, is based on Bellman Optimality Equation

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t))$$

$$Q^*(s, a) = \sum_{s', r} p(s', r|s, a) \left( r + \gamma \max_{a'} Q^*(s', a') \right)$$

] we can see a parallel between the SARSA structure of the update rule and the update rule used in policy iteration in dynamic programming

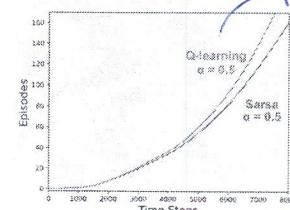
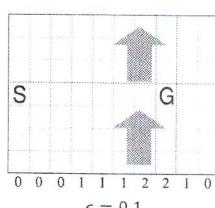
] similarly as above, the Q-learning takes inspiration from the value iteration in dynamic programming (which exploit the Bellman Optimality Equation instead of the Bellman Expectation Equation)

## Q-Learning: Off-Policy TD Control

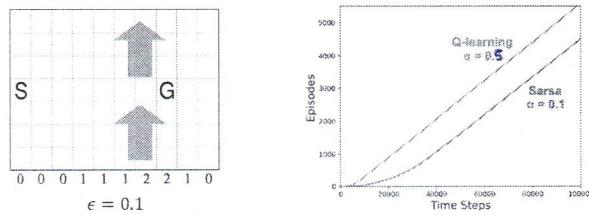
```

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$ 
Initialize  $Q(s, a)$ , for all  $s \in S^+$ ,  $a \in \mathcal{A}(s)$ , 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.,  $\varepsilon$ -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
  
```

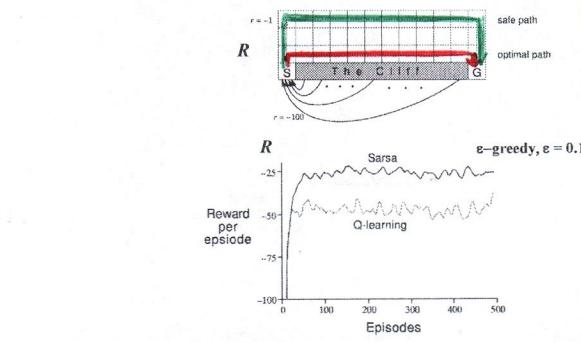
## Q-Learning vs SARSA: Windy Gridworld



Q-learning learns slightly faster than SARSA



## Q-Learning vs SARSA: Cliff Walking Environment



We have a grid: we can move everywhere, we get  $-1$  as reward everywhere, we have to reach a goal starting from state  $S$ , we have special states ("the cliff") and when we move into these states we have a deterministic transition back to the start and a reward that is  $-100$ .

The optimal path is the red (and Q-learn is able to learn it). SARSA is not able to learn that path, it's going to learn the safest one (the green).

Why? Because SARSA goes with an  $\epsilon$ -greedy policy → it won't choose a path which has a probability (small, but)  $> 0$  to end up in the cliff. Q-learning doesn't have this problem because in the maximization step it doesn't take into account the worst possible scenario (only the optimal one).

→ Depending on the problem settings we are into, we may want to prefer one or the other. (depending on how much we care about the safety w.r.t. the optimality)