

# Model Predictive Control and Reinforcement Learning – Monte Carlo RL, Temporal Difference and Q-Learning –

Joschka Boedecker and Moritz Diehl

University Freiburg

October 6, 2023

**universität freiburg**



- 1 Markov Decision Processes
- 2 Policies and Value Functions
- 3 Policy and Value Iteration
- 4 Monte Carlo Reinforcement Learning
- 5 Monte Carlo Prediction
- 6 Monte Carlo Control
- 7 TD Prediction
- 8 TD Control (SARSA, Q-Learning)

# Acknowledgement

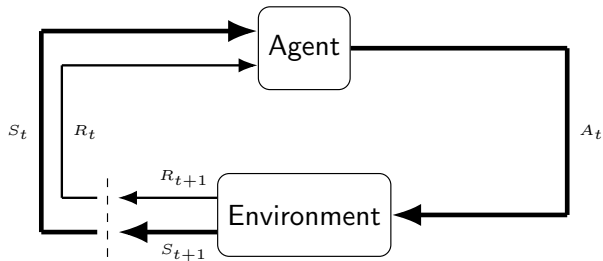


Slide contents are partially based on *Reinforcement Learning: An Introduction* by Sutton and Barto and the Reinforcement Learning lecture by David Silver.



- 1 Markov Decision Processes
- 2 Policies and Value Functions
- 3 Policy and Value Iteration
- 4 Monte Carlo Reinforcement Learning
- 5 Monte Carlo Prediction
- 6 Monte Carlo Control
- 7 TD Prediction
- 8 TD Control (SARSA, Q-Learning)

# Agent and Environment



Time steps  $t$ :  $0, 1, 2, \dots$

States:  $S_0, S_1, S_2, \dots$

Actions:  $A_0, A_1, A_2, \dots$

Rewards:  $R_1, R_2, R_3, \dots$



A finite Markov Decision Process (MDP) is a 4-tuple  $\langle \mathcal{S}, \mathcal{A}, p, \mathcal{R} \rangle$ , where

- ▶  $\mathcal{S}$  is a finite set of states,
- ▶  $\mathcal{A}$  is a finite set of actions,
- ▶  $p$  is the transition probability function  $p : \mathcal{S} \times \mathcal{R} \times \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ ,
- ▶ and  $\mathcal{R}$  is a finite set of scalar rewards. We can then define expected reward  $r(s, a) = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$  and  $r(s, a, s') = \mathbb{E}[R_{t+1} | S_t = s, A_t = a, S_{t+1} = s']$ .

## Markov Property

A state-reward pair  $(S_{t+1}, R_{t+1})$  has the Markov property iff:

$$\Pr\{S_{t+1}, R_{t+1} | S_t, A_t\} = \Pr\{S_{t+1}, R_{t+1} | S_t, A_t, \dots, S_0, A_0\}.$$

*The future is independent of the past given the present.*



A finite Markov Decision Process (MDP) is a 4-tuple  $\langle \mathcal{S}, \mathcal{A}, p, \mathcal{R} \rangle$ , where

- ▶  $\mathcal{S}$  is a finite set of states,
- ▶  $\mathcal{A}$  is a finite set of actions,
- ▶  $p$  is the transition probability function  $p : \mathcal{S} \times \mathcal{R} \times \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ ,
- ▶ and  $\mathcal{R}$  is a finite set of scalar rewards. We can then define expected reward  $r(s, a) = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$  and  $r(s, a, s') = \mathbb{E}[R_{t+1} | S_t = s, A_t = a, S_{t+1} = s']$ .

A deterministic system is a special case of an MDP:

$$p(s_{t+1} | s_t, u_t) = \begin{cases} 1 & s_{t+1} = f(s_t, a_t) \\ 0 & \text{otherwise} \end{cases}$$



- ▶ A reward  $R_t$  in time step  $t$  is a **scalar** feedback signal.
- ▶  $R_t$  indicates how well an agent is performing **at single time step  $t$** .

## Reward Hypothesis

All of what we mean by goals and purposes can be well thought of as the maximization/minimization of the expected value of the cumulative sum of a received scalar signal (called reward/cost).

Examples:

- ▶ Chess: +1 for winning, -1 for losing
- ▶ Walking: +1 for every time step not falling over
- ▶ Investment Portfolio: difference in value between two time steps





- ▶ The agent aims at maximizing the expected **cumulative reward**
- ▶ Non-discounted:  $G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$
- ▶ Discounted:  $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$
- ▶ Discounting with  $\gamma \in [0, 1]$  to prevent from infinite returns (e.g. in infinite horizon control problems)
- ▶ Returns at successive time steps are related to each other:

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$



## Description

Imagine a house cleaning robot. It can have three charge levels: *high*, *low* and *none*. At every point in time, the robot can decide to *recharge* or to *explore* unless it has no battery. When exploring, the charge level can reduce with probability  $\rho$ . Exploring is preferable to recharging, however it has to avoid running out of battery.

Formalize the above problem as an MDP.



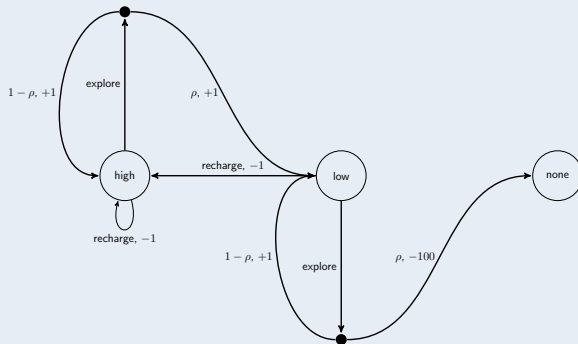
## Solution

For the given problem, we set:

- ▶  $\mathcal{S} = \{\text{high}, \text{low}, \text{none}\}$
- ▶  $\mathcal{A} = \{\text{explore}, \text{recharge}\}$
- ▶  $\mathcal{R} = \{+1, -1, -100\}$  for exploring, recharging, and transitions leading to none, respectively.
- ▶  $p$  has entries with value 1 for transitions (high,  $-1$ , high, recharge), (low,  $-1$ , high, recharge) and (none, 0, none,  $\cdot$ ). It further has entries with value  $\rho$  for transitions (high,  $+1$ , low, explore) and (low,  $-100$ , none, explore) and entries with value  $1 - \rho$  for transitions (high,  $+1$ , high, explore) and (low,  $+1$ , low, explore).

## Solution

The transition graph therefore is:





- 1 Markov Decision Processes
- 2 Policies and Value Functions**
- 3 Policy and Value Iteration
- 4 Monte Carlo Reinforcement Learning
- 5 Monte Carlo Prediction
- 6 Monte Carlo Control
- 7 TD Prediction
- 8 TD Control (SARSA, Q-Learning)



- ▶ The policy defines the behaviour of the agent:
  - ▶ can be stochastic:  $\pi(a|s) = \mathbb{P}[A_t = a|S_t = s]$
  - ▶ or deterministic:  $\pi(s) = a$
- ▶ Due to the Markov property, knowledge of the current state  $s$  is sufficient to make an informed decision.



- ▶ Value Function  $v_\pi(s)$  is the expected return when starting in  $s$  and following  $\pi$ :

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s \right]$$

- ▶ Action-Value Function  $q_\pi$  is the expected return when starting in  $s$ , taking action  $a$  and following  $\pi$  thereafter:

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s, A_t = a \right]$$

- ▶ Simple connection:

$$v_\pi(s) = \mathbb{E}_\pi[q_\pi(s, \pi(s))] \tag{1}$$



# Bellman Equation

- ▶ The Bellman Equation expresses a relationship between the value of a state and the values of its successor states
- ▶ The value function  $v_\pi$  is the unique solution to its Bellman Equation

$$\begin{aligned}v_\pi(s) &= \mathbb{E}_\pi[G_t | S_t = s] \\&= \mathbb{E}_\pi[R_t + \gamma G_{t+1} | S_t = s] \\&= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s']] \\&= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')]\end{aligned}$$

## Bellman Equation for $v_\pi$

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





For a deterministic system and a deterministic policy, the Bellman Equation simplifies to:

Bellman equation for value-function  $v_\pi$  for a deterministic system and policy

$$v_\pi(s) = r + \gamma v_\pi(f(s, \pi(s))).$$

We equivalently obtain a corresponding system of equations for the Q-function:

Bellman Equation for action-value function  $q_\pi$

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



We consider a policy as optimal if the value (i.e. its expected return under the policy) in every state is at least as high as for any other policy:

## Optimality of a policy $\pi_*$

A policy  $\pi_*$  is called *optimal*  $:\Leftrightarrow$

For all  $s \in \mathcal{S}$  :

$$v_{\pi_*}(s) \geq v_{\pi}(s) \text{ for all } \pi \quad (2)$$

The corresponding *optimal value function* is denoted by  $v_*$ .

- ▶ This requires a search among all, possibly infinitely many, policies. This seems to be rather impractical.
- ▶ Is there an easier way to check if a policy  $\pi$  and corresponding value function  $v_{\pi}$  is actually optimal?



# Bellman Optimality Equation

Intuitively, the Bellman Optimality Equation expresses the fact that the value of a state under an optimal policy must equal the expected return for the best action from that state:

$$\begin{aligned}v_*(s) &= \max_a q_{\pi_*}(s, a) \\&= \max_a \mathbb{E}_{\pi_*} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\&= \max_a \mathbb{E}_{\pi_*} [R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\&= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]\end{aligned}$$

## Bellman Optimality Equation for $v_*$

The Bellman Equation for the optimal value function  $v_*$  is defined as:

$$v_*(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')].$$

# Bellman Optimality Equation



For a deterministic system and a deterministic policy, the Bellman Optimality Equation simplifies to:

Bellman equation for the optimal value-function  $v_*$  for a deterministic system and policy

$$v_*(s) = \max_a [r + \gamma v_*(f(s, a))].$$

Equivalently, there exists a Bellman optimality equation for Q-functions:

Bellman equation for the optimal action-value function  $q_*$

$$q_*(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')].$$

How can we turn these equations into practical algorithms to find optimal policies  $\pi_*$ ?



- 1 Markov Decision Processes
- 2 Policies and Value Functions
- 3 Policy and Value Iteration**
- 4 Monte Carlo Reinforcement Learning
- 5 Monte Carlo Prediction
- 6 Monte Carlo Control
- 7 TD Prediction
- 8 TD Control (SARSA, Q-Learning)



Idea: Alternate **evaluating** the value function  $v_\pi$  and **improving** the policy  $\pi$  to convergence.

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi_* \xrightarrow{E} v_*$$



Compute the state-value function  $v_\pi$  for an arbitrary policy  $\pi$ .

$\forall s \in \mathcal{S}$  :

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

If the environments dynamics are completely known, this is a system of  $|\mathcal{S}|$  simultaneous linear equations in  $|\mathcal{S}|$  unknowns. With the Bellman equation, we can iteratively update an initial approximation  $v_0$ :

$$\begin{aligned} v_{k+1}(s) &\doteq \mathbb{E}_\pi [R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_k(s')] \end{aligned}$$



## 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  $\Delta < \theta$





Once we have the value function for a policy, we consider which action  $a$  to select in a state  $s$  when we follow our old policy  $\pi$  afterwards. To decide this, we look at the Bellman equation of the state-action value function:

$$\begin{aligned}q_{\pi}(s, a) &\doteq \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]\end{aligned}$$

## Policy improvement theorem

Let  $\pi$  and  $\pi'$  be any pair of deterministic policies. If,  $\forall s \in S$ ,

$$q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s),$$

then the policy  $\pi'$  must be as good as, or better than,  $\pi$ . It follows that,  $\forall s \in S$ :

$$v_{\pi'}(s) \geq v_{\pi}(s)$$



To implement this, we compute  $q_\pi(s, a)$  for *all* states and *all* actions, and consider the greedy policy:

$$\begin{aligned}\pi'(s) &\doteq \arg \max_a q_\pi(s, a) \\ &= \arg \max_a \mathbb{E} [R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = a] \\ &= \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')]\end{aligned}$$



## 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



Performing policy evaluation to convergence *in every iteration* is costly and often not necessary. A special case is to evaluate just once and combine it with the policy improvement step:

$$\begin{aligned}v_{k+1}(s) &\doteq \max_a \mathbb{E} [R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_k(s')]\end{aligned}$$



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

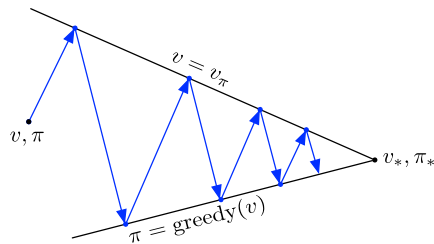
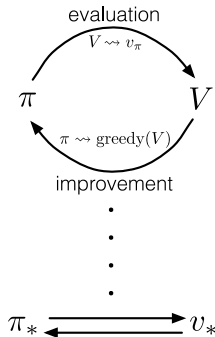
Algorithm parameter: a small threshold  $\theta > 0$  determining accuracy of estimation  
Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

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

# Generalized Policy Iteration



- ▶ Policy Evaluation: estimate  $v_\pi$
- ▶ Policy Improvement: greedy



- ▶ MDPs allow us to formalize RL (and more generally, stochastic optimal control) problems, 4-tuple  $\langle \mathcal{S}, \mathcal{A}, p, \mathcal{R} \rangle$ , assume Markov Property holds
- ▶ Bellman Equations express a relationship between the value of a state and the values of its successor states, provide structure to search for an optimal policy *intelligently*
- ▶ Policy Iteration and Value iteration use the structure of the Bellman Equations and turn them into iterative algorithms for finding optimal policies given an MDP (with and without explicit representation of the policy)



- 1 Markov Decision Processes
- 2 Policies and Value Functions
- 3 Policy and Value Iteration
- 4 Monte Carlo Reinforcement Learning**
- 5 Monte Carlo Prediction
- 6 Monte Carlo Control
- 7 TD Prediction
- 8 TD Control (SARSA, Q-Learning)





Estimate / Optimize the value function of an *unknown* MDP.

- ▶ MC methods learn from episodes of *experiences*  
*experiences* = sequences of states, actions, and rewards
- ▶ MC is model-free: no knowledge required about MDP dynamics
- ▶ MC learns from complete episodes (no bootstrapping), based on averaging sample returns



- 1 Markov Decision Processes
- 2 Policies and Value Functions
- 3 Policy and Value Iteration
- 4 Monte Carlo Reinforcement Learning
- 5 Monte Carlo Prediction**
- 6 Monte Carlo Control
- 7 TD Prediction
- 8 TD Control (SARSA, Q-Learning)



- ▶ Goal: learn the state-value function  $v_\pi$  for a given policy  $\pi$

$$S_0, A_0, R_1, \dots, S_T \sim \pi$$

- ▶ Idea: estimate it from experience by average the returns observed after visits to that state
- ▶ Recall: the *return* is the total discounted reward

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

- ▶ Recall: the value function is the expected return

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$$

- ▶ Monte-Carlo policy prediction uses the empirical mean return instead of expected return



- ▶ We can compute the mean of a sequence  $x_1, x_2, \dots$  incrementally:

$$\begin{aligned}\mu_k &= \frac{1}{k} \sum_{j=1}^k x_j \\ &= \frac{1}{k} \left( x_k + \sum_{j=1}^{k-1} x_j \right) \\ &= \frac{1}{k} (x_k + (k-1)\mu_{k-1}) \\ &= \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})\end{aligned}$$



- ▶ Thus, we can update  $V(s)$  incrementally by:

$$V(s) \leftarrow V(s) + \frac{1}{N(s)}(G_t - V(s)),$$

where  $\frac{1}{N(s)}$  is the state-visitation counter

- ▶ Instead  $\frac{1}{k}$ , we can use step size  $\alpha$  to calculate a running mean:

$$V(s) \leftarrow V(s) + \alpha(G_t - V(s))$$



## 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  $S_0, S_1, \dots, S_{t-1}$ :

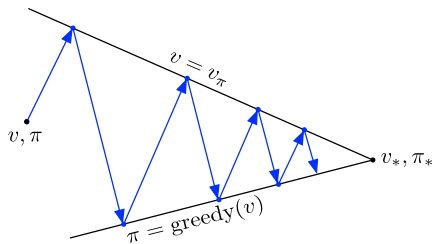
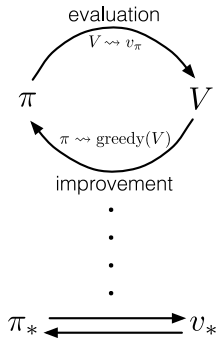
Append  $G$  to  $Returns(S_t)$

$V(S_t) \leftarrow \text{average}(Returns(S_t))$



- 1 Markov Decision Processes
- 2 Policies and Value Functions
- 3 Policy and Value Iteration
- 4 Monte Carlo Reinforcement Learning
- 5 Monte Carlo Prediction
- 6 Monte Carlo Control**
- 7 TD Prediction
- 8 TD Control (SARSA, Q-Learning)

# Generalized Policy Iteration with MC Evaluation



- ▶ Monte Carlo Policy Evaluation:  $V \approx v_\pi$
- ▶ Policy Improvement: greedy?



# Monte Carlo Estimation of Action Values

- ▶ Greedy policy improvement over  $V(s)$  requires a model of the MDP

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

- ▶ Greedy policy improvement over  $Q(s, a)$  is model-free

$$\pi(s) = \arg \max_{a \in \mathcal{A}} Q(s, a)$$

Generalized Policy Iteration with action-value function:

- ▶ Monte Carlo Policy Evaluation:  $Q \approx q_\pi$
- ▶ Policy Improvement: greedy?



- ▶ We have to ensure that each state-action pair is visited a sufficient (infinite) number of times
- ▶ Simple idea:  $\epsilon$ -greedy
- ▶ All actions have non-zero probability
- ▶ With probability  $\epsilon$  choose a random action, with probability  $1 - \epsilon$  take the greedy action.

$$\pi(a|s) = \begin{cases} \frac{\epsilon}{|\mathcal{A}|} + 1 - \epsilon & \text{if } a = \arg \max_{a' \in \mathcal{A}} Q(s, a') \\ \frac{\epsilon}{|\mathcal{A}|} & \text{otherwise} \end{cases}$$



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

Algorithm parameter: small  $\varepsilon > 0$

Initialize:

$\pi \leftarrow$  an arbitrary  $\varepsilon$ -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  $Returns(S_t, A_t)$

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

$A^* \leftarrow \operatorname{argmax}_a Q(S_t, a)$  (with ties broken arbitrarily)

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

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



Estimate/ optimize the value function of an *unknown* MDP using Temporal Difference Learning.

- ▶ TD is a combination of Monte Carlo and dynamic programming ideas
- ▶ Similar to MC methods, TD methods learn directly raw experiences without a dynamic model
- ▶ TD learns from *incomplete* episodes by bootstrapping
- ▶ Bootstrapping: update estimated based on other estimates without waiting for a final outcome (update a guess towards a guess)



- 1 Markov Decision Processes
- 2 Policies and Value Functions
- 3 Policy and Value Iteration
- 4 Monte Carlo Reinforcement Learning
- 5 Monte Carlo Prediction
- 6 Monte Carlo Control
- 7 TD Prediction**
- 8 TD Control (SARSA, Q-Learning)



## Monte Carlo Update

Update value  $V(S_t)$  towards the *actual* return  $G_t$ .

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

$\alpha$  is a step-size parameter.

## Simplest temporal-difference learning algorithm: $TD(0)$

Update value  $V(S_t)$  towards the *estimated* return  $R_{t+1} + \gamma V(S_{t+1})$ .

$$V(s_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

- ▶  $R_{t+1} + \gamma V(S_{t+1})$  is called the *TD target*
- ▶  $\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$  is called the *TD error*



## Tabular TD(0) for estimating $v_\pi$

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

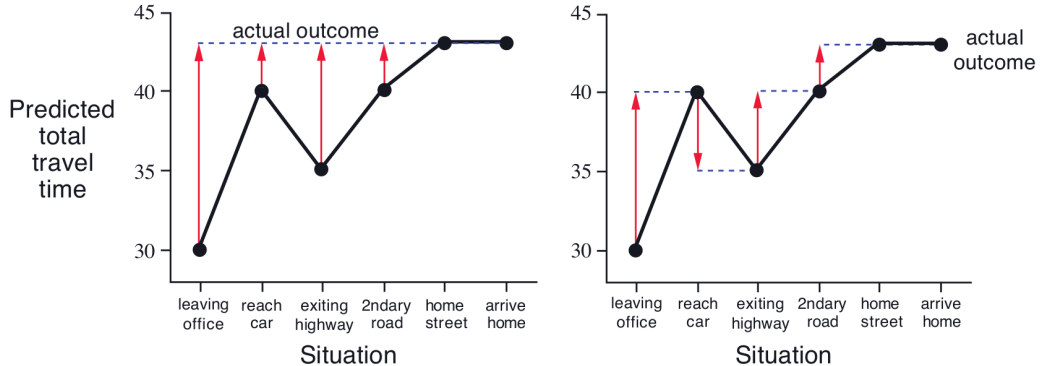
# Driving Home Example



State	Elapsed Time (minutes)	Predicted Time to Go	Predicted Total Time
leaving office, friday at 6	0	30	30
reach car, raining	5	35	40
exiting highway	20	15	35
2ndary road, behind truck	30	10	40
entering home street	40	3	43
arrive home	43	0	43



# TD Prediction



**Figure 6.1:** Changes recommended in the driving home example by Monte Carlo methods (left) and TD methods (right).



- ▶ TD can learn online after every step, MC has to wait for the final outcome/return
- ▶ TD can even learn without ever getting a *final* outcome, which is especially important for infinite horizon tasks
- ▶ The return  $G_t$  depends on many random actions, transitions and rewards, the TD-target depends on one random action, transition and reward
- ▶ Therefore, the TD-target has lower *variance* than the return
- ▶ But the TD-target is a *biased* estimate of  $v_\pi$
- ▶ This is known as the bias/variance trade-off



- ▶ MC and TD converge if every state and every action are visited an infinite number of times
- ▶ What about finite experience?

Imagine two states,  $A$  and  $B$ , and the following transitions:

$A,0,B,0$	$B,1$
$B,1$	$B,1$
$B,1$	$B,1$
$B,1$	$B,0$

What are the values of  $A$  and  $B$  given this data?

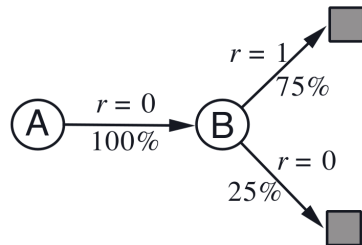


- ▶ W.r.t.  $B$ , the process terminated immediately  $6/8$  times with a return of 1, 0 otherwise
- ▶ Thus, it is reasonable to assume a value of 0.75
- ▶ What about  $A$ ?



- ▶ W.r.t.  $B$ , the process terminated immediately  $6/8$  times with a return of 1, 0 otherwise
- ▶ Thus, it is reasonable to assume a value of 0.75
- ▶ What about  $A$ ?

$A$  led to  $B$  in all cases. Thus,  $A$  could have a value of 0.75 as well. This answer is based on first modelling the Markov Process and then computing the values given the model. TD is leading to this value. MC gives a value of 0 – which is also the solution with 0 MSE on the given data. One can assume, however, that the former gives lower error on *future* data.





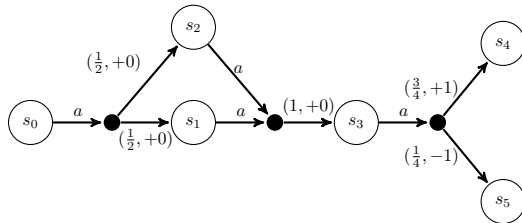
- ▶ Batch MC converges to the solution with minimum MSE on the observed returns
- ▶ Batch TD converges to the solution of the maximum-likelihood Markov model
- ▶ Given this model, we can compute the estimate of the value-function that would be exactly correct, if the model were exactly correct
- ▶ This is called the *Certainty Equivalence*

# MC vs TD: Example



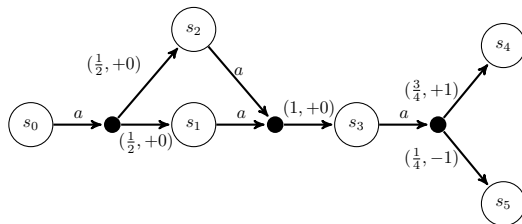
Assume that the agent encounters the following set of trajectories at every iteration  $i$  (where  $i \bmod 4 = 0$ ):

$t_i$	:	$s_0$	$\rightarrow$	$s_1$	$\rightarrow$	$s_3$	$\rightarrow$	$s_4$
$t_{i+1}$	:	$s_0$	$\rightarrow$	$s_1$	$\rightarrow$	$s_3$	$\rightarrow$	$s_4$
$t_{i+2}$	:	$s_0$	$\rightarrow$	$s_1$	$\rightarrow$	$s_3$	$\rightarrow$	$s_4$
$t_{i+3}$	:	$s_0$	$\rightarrow$	$s_2$	$\rightarrow$	$s_3$	$\rightarrow$	$s_5$



## Description

Given these trajectories, explain why TD-learning is better fitted to estimate the value function compared to MC. Assume no discount and that the value function is initialized with zeros. To which value function is MC going to converge, given a suitable learning rate  $\alpha$ ? What about TD?



## Solution

MC always takes the full return to update its values. Therefore,  $s_1$  only updates on return  $+1$ , whereas  $s_2$  only updates on return  $-1$ . TD takes this into account due to bootstrapping. MC converges to:  $v(s_0) = v(s_3) = 0.5$ ,  $v(s_1) = +1$  and  $v(s_2) = -1$ . TD converges to the true value function  $v(s_0) = v(s_1) = v(s_2) = v(s_3) = 0.5$  and  $v(s_4) = v(s_5) = 0$ , since  $\frac{3}{4}$  trajectories end with a return of  $+1$  and  $\frac{1}{4}$  with a return of  $-1$  – which corresponds to the true distribution of the MDP.





- 1 Markov Decision Processes
- 2 Policies and Value Functions
- 3 Policy and Value Iteration
- 4 Monte Carlo Reinforcement Learning
- 5 Monte Carlo Prediction
- 6 Monte Carlo Control
- 7 TD Prediction
- 8 TD Control (SARSA, Q-Learning)



- ▶ SARSA: State, Action, Reward, State, Action
- ▶ Why is it considered an on-policy method?

## 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.,  $\varepsilon$ -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



- ▶ Why is it considered an off-policy method?

## 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.,  $\varepsilon$ -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



- ▶ Monte Carlo RL methods average sample returns from episodes of experience interacting with the environment, making it possible to learn without a given transition model
- ▶ Temporal Difference methods learn from incomplete episodes by bootstrapping, combining ideas from Monte Carlo and dynamic programming
- ▶ TD can learn online after every step, MC has to wait for the final outcome/return
- ▶ TD target has lower variance than the MC return, but is biased due to bootstrapping with wrong initial values
- ▶ Q-Learning can learn to approximate  $q_*$  even while gathering data with a different policy (off-policy learning)