

# CS234 Notes - Lecture 3

## Model Free Policy Evaluation: Policy Evaluation Without Knowing How the World Works

Youkow Homma, Emma Brunskill

March 20, 2018

### 4 Model-Free Policy Evaluation

*{* *Markov process*  
*Temporal Difference*

In the previous lecture, we began by discussing three problem formulations of increasing complexity, which we recap below.

1. A Markov process (MP) is a stochastic process augmented with the Markov property.
2. A Markov reward process (MRP) is a Markov process with rewards at each time step and the accumulation of discounted rewards, called values.
3. A Markov decision process (MDP) is a Markov reward process augmented with choices, or actions, at each state.

In the second half of last lecture, we discussed two methods for evaluating a given policy in an MDP and three methods for finding the optimal policy of an MDP. The two methods for policy evaluation were directly solving via a linear system of equations and dynamic programming. The three methods for control were brute force policy search, policy iteration and value iteration.

Implicit in all of these methods was the assumption that we know both the rewards and probabilities for every transition. However, in many cases, such information is not readily available to us, which necessitates **model-free algorithms**. In this set of lectures notes, we will be discussing **model-free policy evaluation**. That is, we will be given a policy and will attempt to learn the value of that policy without leveraging knowledge of the rewards or transition probabilities. In particular, we will not be discussing how to improve our policies in the model-free case until next lecture.

#### 4.1 Notation Recap

Before diving into some methods for model-free policy evaluation, we'll first recap some of the notation surrounding MDP's from last lecture that we'll need in this lecture.

Recall that we define the return of an MRP as the discounted sum of rewards starting from time step  $t$  and ending at horizon  $H$ , where  $H$  may be infinite. Mathematically, this takes the form

$$G_t = \sum_{i=t}^{H-1} \gamma^{i-t} r_i, \quad (1)$$

for  $0 \leq t \leq H - 1$ , where  $0 < \gamma \leq 1$  is the discount factor and  $r_i$  is the reward at time step  $i$ . For an MDP, the return  $G_t$  is defined identically, and the rewards  $r_i$  are generated by following policy  $\pi(a|s)$ .

The state value function  $V^\pi(s)$  is the expected return from starting at state  $s$  under stationary policy  $\pi$ . Mathematically, we can express this as

$$V^\pi(s) = \mathbb{E}_\pi[G_t | s_t = s] \quad (2)$$

$$= \mathbb{E}_\pi \left[ \sum_{i=t}^{H-1} \gamma^{i-t} r_i | s_t = s \right]. \quad (3)$$

The state-action value function  $Q^\pi(s, a)$  is the expected return from starting in state  $s$ , taking action  $a$ , and then following stationary policy  $\pi$  for all transitions thereafter. Mathematically, we can express this as

$$Q^\pi(s, a) = \mathbb{E}_\pi[G_t | s_t = s, a_t = a] \quad (4)$$

$$= \mathbb{E}_\pi \left[ \sum_{i=t}^{H-1} \gamma^{i-t} r_i | s_t = s, a_t = a \right]. \quad (5)$$

Throughout this lecture, we will assume an infinite horizon as well as stationary rewards, transition probabilities and policies. This allows us to have time-independent state and state-action value functions, as derived last lecture.

There is one new definition that we will use in this lecture, which is the definition of a history.

**Definition 4.1.** *The **history** is the ordered tuple of states, actions and rewards that an agent experiences. In episodic domains, we will use the word *episode* interchangeably with *history*. When considering many interactions, we will index the histories in the following manner: the  $j$ th history is*

$$h_j = (s_{j,1}, a_{j,1}, r_{j,1}, s_{j,2}, a_{j,2}, r_{j,2}, \dots, s_{j,L_j}),$$

where  $L_j$  is the length of the interaction, and  $s_{j,t}, a_{j,t}, r_{j,t}$  are the state, action and reward at time step  $t$  in history  $j$ , respectively.

$L_j$  is basically  $H_j$

## 4.2 Dynamic Programming

Recall also from last lecture the dynamic programming algorithm to calculate the value of an infinite horizon MDP with  $\gamma < 1$  under policy  $\pi$ , which we rewrite here for convenience as Algorithm 1.

---

**Algorithm 1** Iterative algorithm to calculate MDP value function for a stationary policy  $\pi$

- 1: **procedure** POLICY EVALUATION( $M, \pi, \epsilon$ )
  - 2:   For all states  $s \in S$ , define  $R^\pi(s) = \sum_{a \in A} \pi(a|s)R(s, a)$
  - 3:   For all states  $s, s' \in S$ , define  $P^\pi(s'|s) = \sum_{a \in A} \pi(a|s)P(s'|s, a)$
  - 4:   For all states  $s \in S$ ,  $V_0(s) \leftarrow 0$
  - 5:    $k \leftarrow 0$
  - 6:   **while**  $k = 0$  or  $\|V_k - V_{k-1}\|_\infty > \epsilon$  **do**
  - 7:      $k \leftarrow k + 1$
  - 8:     For all states  $s \in S$ ,  $V_k(s) = R^\pi(s) + \gamma \sum_{s' \in S} P^\pi(s'|s)V_{k-1}(s')$
  - 9: **return**  $V_k$
- 

But we don't know  $R$  and  $P$  for model-free policy evaluation so DP is not a model-free algorithm.

Written in this form, we can think of  $V_k$  in a few ways. First,  $V_k(s)$  is the exact value of following policy  $\pi$  for  $k$  additional transitions, starting at state  $s$ . Second, for large  $k$ , or when Algorithm 1 terminates,  $V_k(s)$  is an estimate of the true, infinite horizon value  $V^\pi(s)$  of state  $s$ .

We can additionally express the behavior of this algorithm via a **backup diagram**, which is shown in Figure 1. This backup diagram is read top-down with white circles representing states, black circles

*policy is deterministic in  
this graph.*

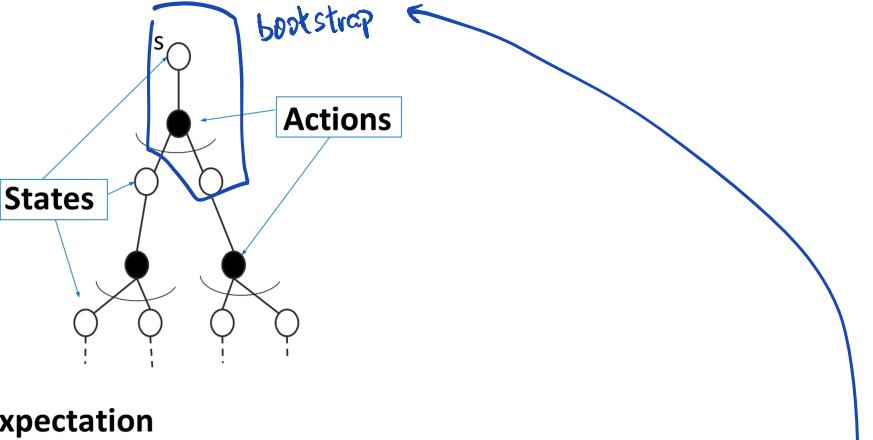


Figure 1: Backup diagram for the dynamic programming policy evaluation algorithm.

representing actions and arcs representing taking the expectation. The diagram shows the branching effect of starting at state  $s$  at the top and transitioning two time steps as we move down the diagram. It furthermore shows how after starting at state  $s$  and taking an action under policy  $\pi$ , we take the expectation over the value of the next state. In dynamic programming, we **bootstrap**, or estimate, the value of the next state using our current estimate,  $V_{k-1}(s')$ .

### 4.3 Monte Carlo On Policy Evaluation

We now describe our first model-free policy evaluation algorithm which uses a popular computational method called the Monte Carlo method. We first walk through an example of the Monte Carlo method outside the context of reinforcement learning, then discuss the method more generally, and finally apply Monte Carlo to reinforcement learning. We emphasize here that this method only works in episodic environments, and we'll see why this is as we examine the algorithm more carefully in this section. H $\infty$

Suppose we want to estimate how long the commute from your house to Stanford's campus will take today. Suppose we also have access to a commute simulator which models our uncertainty of how bad the traffic will be, the weather, construction delays, and other variables, as well as how these variables interact with each other. One way to estimate the expected commute time is to simulate our commute many times on the simulator and then take an average over the simulated commute times. This is called a Monte Carlo estimate of our commute time.

In general, we get the Monte Carlo estimate of some quantity by observing many iterations of how that quantity is generated either in real life or via simulation and then averaging over the observed quantities. By the law of large numbers, this average converges to the expectation of the quantity.

In the context of reinforcement learning, the quantity we want to estimate is  $V^\pi(s)$ , which is the average of returns  $G_t$  under policy  $\pi$  starting at state  $s$ . We can thus get a Monte Carlo estimate of  $V^\pi(s)$  through three steps:

1. Execute a **rollout** of policy  $\pi$  until termination many times
2. Record the returns  $G_t$  that we observe when starting at state  $s$
3. Take an average of the values we get for  $G_t$  to estimate  $V^\pi(s)$ .

The backup diagram for Monte Carlo policy evaluation can be seen in Figure 2. The new blue line indicates that we sample an entire episode until termination starting at state  $s$ .

There are two forms of Monte Carlo on policy evaluation, which are differentiated by whether we take an average over just the first time we visit a state in each rollout or every time we visit the state in each

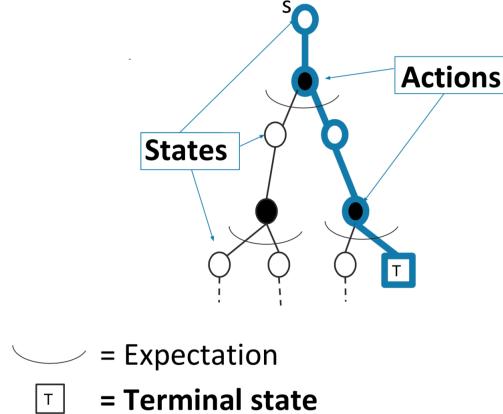


Figure 2: Backup diagram for the Monte Carlo policy evaluation algorithm.

rollout. These are called First-Visit Monte Carlo and Every-Visit Monte Carlo On Policy Evaluation, respectively.

More formally, we describe the First-Visit Monte Carlo in Algorithm 2 and the Every-Visit Monte Carlo in Algorithm 3.

---

**Algorithm 2** First-Visit Monte Carlo Policy Evaluation

---

```

1: procedure FIRST-VISIT-MONTE-CARLO( $h_1, \dots, h_j$ )
2:   For all states  $s$ ,  $N(s) \leftarrow 0$ ,  $S(s) \leftarrow 0$ ,  $V(s) \leftarrow 0$ 
3:   for each episode  $h_j$  do
4:     for  $t = 1, \dots, L_j$  do
5:       if  $s_{j,t} \neq s_{j,u}$  for  $u < t$  then - only difference compared to Every-Visit
6:          $N(s_{j,t}) \leftarrow N(s_{j,t}) + 1$ 
7:          $S(s_{j,t}) \leftarrow S(s_{j,t}) + G_{j,t}$ 
8:          $V^\pi(s_{j,t}) \leftarrow S(s_{j,t})/N(s_{j,t})$ 
9:   return  $V^\pi$ 
```

---



---

**Algorithm 3** Every-Visit Monte Carlo Policy Evaluation

---

```

1: procedure EVERY-VISIT-MONTE-CARLO( $h_1, \dots, h_j$ )
2:   For all states  $s$ ,  $N(s) \leftarrow 0$ ,  $S(s) \leftarrow 0$ ,  $V(s) \leftarrow 0$ 
3:   for each episode  $h_j$  do
4:     for  $t = 1, \dots, L_j$  do
5:        $N(s_{j,t}) \leftarrow N(s_{j,t}) + 1$ 
6:        $S(s_{j,t}) \leftarrow S(s_{j,t}) + G_{j,t}$ 
7:        $V^\pi(s_{j,t}) \leftarrow S(s_{j,t})/N(s_{j,t})$ 
8:   return  $V^\pi$ 
```

---

Note that in the body of the for loop in Algorithms 2 and 3, we can remove vector  $S$  and replace the update for  $V^\pi(s_{j,t})$  with

$$V^\pi(s_{j,t}) \leftarrow V^\pi(s_{j,t}) + \frac{1}{N(s_{j,t})}(G_{j,t} - V^\pi(s_{j,t})). \quad (6)$$

This is because the new average is the average of  $N(s_{j,t}) - 1$  of the old values  $V^\pi(s_{j,t})$  and the new return  $G_{j,t}$ , giving us

$$V^\pi(s_{j,t}) \leftarrow \frac{V^\pi(s_{j,t}) \times (N(s_{j,t}) - 1) + G_{j,t}}{N(s_{j,t})} = V^\pi(s_{j,t}) + \frac{1}{N(s_{j,t})}(G_{j,t} - V^\pi(s_{j,t})), \quad (7)$$

First-Visit is unbiased  
Every-visit is biased but often better MSE. If actually Markovian, every-visit is more data efficient.

which is precisely the new form of the update.

Replacing  $\frac{1}{N(s_{j,t})}$  with  $\alpha$  in this new update gives us the more general **Incremental Monte Carlo**

**On Policy Evaluation.** Algorithms 4 and 5 detail this procedure in the First-Visit and Every-Visit cases, respectively.

---

**Algorithm 4** Incremental First-Visit Monte Carlo Policy Evaluation

```

1: procedure INCREMENTAL-FIRST-VISIT-MONTE-CARLO( $\alpha, h_1, \dots, h_j$ )
2:   For all states  $s$ ,  $N(s) \leftarrow 0$ ,  $V(s) \leftarrow 0$ 
3:   for each episode  $h_j$  do
4:     for  $t = 1, \dots, \text{terminal}$  do
5:       if  $s_{j,t} \neq s_{j,u}$  for  $u < t$  then
6:          $N(s_{j,t}) \leftarrow N(s_{j,t}) + 1$ 
7:          $V^\pi(s_{j,t}) \leftarrow V^\pi(s) + \alpha(G_{j,t} - V^\pi(s))$ 
8:   return  $V^\pi$ 

```

---

Analogous to algorithm 2

**Algorithm 5** Incremental Every-Visit Monte Carlo Policy Evaluation

Analogous to algorithm 3

```

1: procedure INCREMENTAL-EVERY-VISIT-MONTE-CARLO( $\alpha, h_1, \dots, h_j$ )
2:   For all states  $s$ ,  $N(s) \leftarrow 0$ ,  $V(s) \leftarrow 0$ 
3:   for each episode  $h_j$  do
4:     for  $t = 1, \dots, \text{terminal}$  do
5:        $N(s_{j,t}) \leftarrow N(s_{j,t}) + 1$ 
6:        $V^\pi(s_{j,t}) \leftarrow V^\pi(s) + \alpha(G_{j,t} - V^\pi(s))$ 
7:   return  $V^\pi$ 

```

---

Setting  $\alpha = \frac{1}{N(s_{j,t})}$  recovers the original Monte Carlo On Policy Evaluation algorithms given in Algorithms 2 and 3, while setting  $\alpha > \frac{1}{N(s)}$  gives a higher weight to newer data, which can help learning in non-stationary domains. If we are in a truly Markovian-domain, Every-Visit Monte Carlo will be more data efficient because we update our average return for a state every time we visit the state.

**Exercise 4.1.** Recall our Mars Rover MDP from last lecture, shown in Figure 3 below. Suppose that our estimate for the value of each state is currently 0. If we experience the history

$$h = (S3, TL, +0, S2, TL, +0, S1, TL, +1, \text{terminal}),$$

then:

1. What is the first-visit Monte Carlo estimate of  $V$  at each state? }  $V(S_3) = Y^2$ ,  $V(S_2) = Y$ ,  $V(S_1) = 1$
  2. What is the every-visit Monte Carlo estimate of each state?
  3. What is the incremental first-visit Monte Carlo estimate of  $V$  with  $\alpha = \frac{2}{3}$ ? }  $V(S_3) = \frac{2}{3}Y^2$ ,  $V(S_2) = \frac{2}{3}Y$ ,  $V(S_1) = \frac{2}{3}$
  4. What is the incremental every-visit Monte Carlo estimate of  $V$  with  $\alpha = \frac{2}{3}$ ? }
- first-visit and every-visit yield the same result because we didn't visit any state more than once.

#### 4.4 Monte Carlo Off Policy Evaluation

In the section above, we discussed the case where we are able to obtain many realizations of  $G_t$  under the policy  $\pi$  that we want to evaluate. However, in many costly or high stakes situations, we aren't able to obtain rollouts of  $G_t$  under the policy that we wish to evaluate. For example, we may have data associated with one medical policy, but want to determine the value of a different medical policy. In this section, we describe Monte Carlo off policy evaluation, which is a method for using data taken from one policy to evaluate a different policy.

Reward for both policies are the same.

Action needs to exist in both policies  $\rightarrow$  e.g.  $\pi_{b'}(s) = \text{use drug A}$  } off-policy doesn't work because  $q_{\pi}(s) = 0$ .

$\pi_{b'}(s) = \text{drug A for time drug B 10\%}$

off policy works }  
 Tie (b) = drug A 60%      drug B 40%

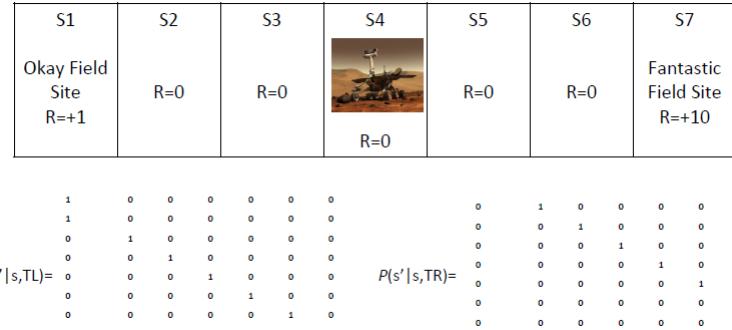


Figure 3: Mars Rover Markov Decision Process with actions Try Left (TL) and Try Right (TR)

#### 4.4.1 Importance Sampling

The key ingredient of off policy evaluation is a method called importance sampling. The goal of importance sampling is to estimate the expected value of a function  $f(x)$  when  $x$  is drawn from distribution  $q$  using only the data  $f(x_1), \dots, f(x_n)$ , where  $x_i$  are drawn from a different distribution  $p$ . In summary, given  $q(x_i), p(x_i), f(x_i)$  for  $1 \leq i \leq n$ , we would like an estimate for  $\mathbb{E}_{x \sim q}[f(x)]$ . We can do this via the following approximation:

$$\mathbb{E}_{x \sim q}[f(x)] = \int_x q(x)f(x)dx \quad (8)$$

$$= \int_x p(x) \left[ \frac{q(x)}{p(x)} f(x) \right] dx \quad (9)$$

$$= \mathbb{E}_{x \sim p} \left[ \frac{q(x)}{p(x)} f(x) \right] \quad (10)$$

$$\approx \sum_{i=1}^n \left[ \frac{q(x_i)}{p(x_i)} f(x_i) \right]. \quad (11)$$

The last equation gives us the **importance sampling estimate** of  $f$  under distribution  $q$  using samples of  $f$  under distribution  $p$ . Note that the first step only holds if  $q(x)f(x) > 0$  implies  $p(x) > 0$  for all  $x$ .

#### 4.4.2 Importance Sampling for Off Policy Evaluation

We now apply the general result of importance sampling estimates to reinforcement learning. In this instance, we want to approximate the value of state  $s$  under policy  $\pi_1$ , given by  $V^{\pi_1}(s) = \mathbb{E}[G_t | s_t = s]$ , using  $n$  histories  $h_1, \dots, h_n$  generated under policy  $\pi_2$ . Using the importance sampling estimate result gives us that

$$V^{\pi_1}(s) \approx \frac{1}{n} \sum_{j=1}^n \frac{p(h_j | \pi_1, s)}{p(h_j | \pi_2, s)} G(h_j), \quad (12)$$

where  $G(h_j) = \sum_{t=1}^{L_j-1} \gamma^{t-1} r_{j,t}$  is the total discounted sum of rewards for history  $h_j$ .

Now, for a general policy  $\pi$ , we have that the probability of experiencing history  $h_j$  under policy  $\pi$  is

$$p(h_j|\pi, s = s_{j,1}) = \prod_{t=1}^{L_j-1} p(a_{j,t}|s_{j,t})p(r_{j,t}|s_{j,t}, a_{j,t})p(s_{j,t+1}|s_{j,t}, a_{j,t}) \quad (13)$$

$$= \boxed{\prod_{t=1}^{L_j-1} \pi(a_{j,t}|s_{j,t})p(r_{j,t}|s_{j,t}, a_{j,t})p(s_{j,t+1}|s_{j,t}, a_{j,t})}, \quad (14)$$

where  $L_j$  is the length of the  $j$ th episode. The first line follows from looking at the three components of each transition. The components are:

1.  $p(a_{j,t}|s_{j,t})$  - probability we take action  $a_{j,t}$  at state  $s_{j,t}$
2.  $p(r_{j,t}|s_{j,t}, a_{j,t})$  - probability we experience reward  $r_{j,t}$  after taking action  $a_{j,t}$  in state  $s_{j,t}$
3.  $p(s_{j,t+1}|s_{j,t}, a_{j,t})$  - probability we transition to state  $s_{j,t+1}$  after taking action  $a_{j,t}$  in state  $s_{j,t}$

Now, combining our importance sampling estimate for  $V^{\pi_1}(s)$  with our decomposition of the history probabilities,  $p(h_j|\pi, s = s_{j,1})$ , we get that

$$V^{\pi_1}(s) \approx \frac{1}{n} \sum_{j=1}^n \frac{p(h_j|\pi_1, s)}{p(h_j|\pi_2, s)} G(h_j) \quad (15)$$

$$= \frac{1}{n} \sum_{j=1}^n \frac{\prod_{t=1}^{L_j-1} \pi_1(a_{j,t}|s_{j,t})p(r_{j,t}|s_{j,t}, a_{j,t})p(s_{j,t+1}|s_{j,t}, a_{j,t})}{\prod_{t=1}^{L_j-1} \pi_2(a_{j,t}|s_{j,t})p(r_{j,t}|s_{j,t}, a_{j,t})p(s_{j,t+1}|s_{j,t}, a_{j,t})} G(h_j) \quad (16)$$

$$= \boxed{\frac{1}{n} \sum_{j=1}^n G(h_j) \prod_{t=1}^{L_j-1} \frac{\pi_1(a_{j,t}|s_{j,t})}{\pi_2(a_{j,t}|s_{j,t})}}. \quad (17)$$

Notice we can now explicitly evaluate the expression without the transition probabilities or rewards since all of the terms involving model dynamics canceled out in the second step of the equation. In particular, we are given the histories  $h_j$ , so we can calculate  $G(h_j) = \sum_{t=1}^{L_j-1} \gamma^{t-1} r_{j,t}$ , and we know the two policies  $\pi_1$  and  $\pi_2$ , so we can also evaluate the second term.

## 4.5 Temporal Difference (TD) Learning

So far, we have two methods for policy evaluation: dynamic programming and Monte Carlo. Dynamic programming leverages bootstrapping to help us get value estimates with only one backup. On the other hand, Monte Carlo samples many histories for many trajectories which frees us from using a model. Now, we introduce a new algorithm that combines bootstrapping with sampling to give us a second model-free policy evaluation algorithm.

To see how to combine sampling with bootstrapping, let's go back to our incremental Monte Carlo update:

$$V^\pi(s_t) \leftarrow V^\pi(s_t) + \alpha(G_t - V^\pi(s_t)). \quad (18)$$

Recall that  $G_t$  is the return after rolling out the policy from time step  $t$  to termination starting at state  $s_t$ . Let's now replace  $G_t$  with a Bellman backup like in dynamic programming. That is, let's replace  $G_t$  with  $r_t + \gamma V^\pi(s_{t+1})$ , where  $r_t$  is a sample of the reward at time step  $t$  and  $V^\pi(s_{t+1})$  is our current estimate of the value at the next state. Making this substitution gives us the TD-learning update

$$V^\pi(s_t) \leftarrow V^\pi(s_t) + \alpha(r_t + \gamma V^\pi(s_{t+1}) - V^\pi(s_t)). \quad (19)$$

The difference

$$\delta_t = r_t + \gamma V^\pi(s_{t+1}) - V^\pi(s_t) \quad (20)$$

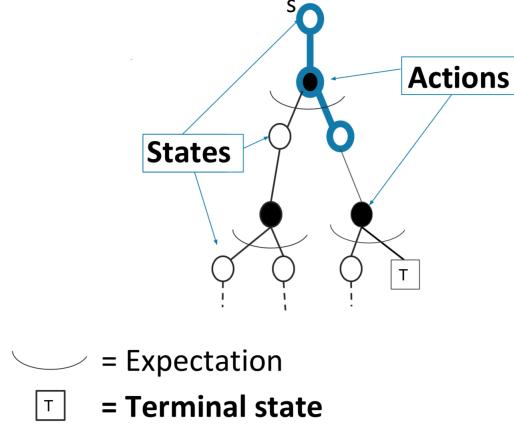


Figure 4: Backup diagram for the TD Learning policy evaluation algorithm.

is commonly referred to as the **TD error**, and the sampled reward combined with the bootstrap estimate of the next state value,

$$r_t + \gamma V^\pi(s_{t+1}), \quad (21)$$

is referred to as the **TD target**. The full TD learning algorithm is given in Algorithm 6. We can see that using this method, we update our value for  $V^\pi(s_t)$  directly after witnessing the transition  $(s_t, a_t, r_t, s_{t+1})$ . In particular, we don't need to wait for the episode to terminate like in Monte Carlo.

---

### TD(0) Algorithm 6 TD Learning to evaluate policy $\pi$

---

```

1: procedure TDLEARNING(step size  $\alpha$ , number of trajectories  $n$ )
2:   For all states  $s$ ,  $V^\pi(s) \leftarrow 0$ 
3:   while  $n > 0$  do
4:     Begin episode  $E$  at state  $s$ 
5:     while  $n > 0$  and episode  $E$  has not terminated do
6:        $a \leftarrow$  action at state  $s$  under policy  $\pi$ 
7:       Take action  $a$  in  $E$  and observe reward  $r$ , next state  $s'$ 
8:        $V^\pi(s) \leftarrow V^\pi(s) + \alpha(R + \gamma V^\pi(s') - V^\pi(s))$ 
9:        $s \leftarrow s'$ 
10:    return  $V^\pi$ 

```

---

We can again examine this algorithm via a backup diagram as shown in Figure 4. Here, we see via the blue line that we sample one transition starting at  $s$ , then we estimate the value of the next state via our current estimate of the next state to construct a full Bellman backup estimate.

There is actually an entire spectrum of ways we can blend Monte Carlo and dynamic programming using a method called  $\text{TD}(\lambda)$ . When  $\lambda = 0$ , we get the TD-learning formulation above, hence giving us the alias  $\text{TD}(0)$ . When  $\lambda = 1$ , we recover Monte Carlo policy evaluation, depending on the formulation used. When  $0 < \lambda < 1$ , we get a blend of these two methods. For a more thorough treatment of  $\text{TD}(\lambda)$ , we refer the interested reader to Sections 7.1 and 12.1-12.5 of Sutton and Barto [1] which detail  $n$ -step TD learning and  $\text{TD}(\lambda)$ /eligibility traces, respectively.

**Exercise 4.2.** Consider again the Mars Rover example in Figure 3. Suppose that our estimate for the value of each state is currently 0. If we experience the history

$$h = (S3, TL, +0, S2, TL, +0, S2, TL, +0, S1, TL, +1, \text{terminal}),$$

then:

$$\begin{aligned}
 1. \text{ What is the } TD(0) \text{ estimate of } V \text{ with } \alpha = 1? & \quad V(s_3) = 0 \quad V(s_2) = 0 \quad V(s_1) = 1 \\
 2. \text{ What is the } TD(0) \text{ estimate of } V \text{ with } \alpha = \frac{2}{3}?
 & \quad V(s_3) = 0 \quad V(s_2) = 0 \quad V(s_1) = \frac{2}{3}
 \end{aligned}$$

## 4.6 Summary of Methods Discussed

### Comparison between MC and TD(0)

In this lecture, we re-examined policy evaluation using dynamic programming from last lecture, and we introduced two new methods for policy evaluation, namely Monte Carlo evaluation and Temporal Difference (TD) learning.

First, we motivated the introduction of Monte Carlo and TD-Learning by noting that dynamic programming relied on a model of the world. That is, we needed to feed our dynamic programming policy evaluation algorithm with the rewards and transition probabilities of the domain. Monte Carlo and TD-Learning are both free from this constraint, making them **model-free** methods.

In Monte Carlo policy evaluation, we generate many histories and then average the returns over the states that we encounter. In order for us to generate these histories in a finite amount of time, we require the domain to be **episodic** - that is, we need to ensure that each history that we observe terminates. In both dynamic programming and temporal difference learning, we only backup over one transition (we only look one step ahead in the future), so termination of histories is not a concern, and we can apply these algorithms to non-episodic domains.

On the flip side, the reason we are able to backup over just one transition in dynamic programming and TD learning is because we leverage the **Markovian assumption** of the domain. Furthermore, **Incremental Monte Carlo policy evaluation**, described in Algorithms 4 and 5 can be utilized in **non-Markovian domains**.

In all three methods, we converge to the true value function. Last lecture, we proved this result for dynamic programming by using the fact that the Bellman backup operator is a contraction. We saw in today's lecture that Monte Carlo policy evaluation converges to the policy's value function due to the law of large numbers. TD(0) converges to the true value as well, which we will look at more closely in the next section on batch learning.

Because we are taking an average over the true distribution of returns in Monte Carlo, we obtain an **unbiased estimator** of the value at each state. On the other hand, in TD learning, we bootstrap the next state's value estimate to get the **current state's value estimate**, so the estimate is biased by the estimated value of the next state. Further discussions on this can be found in section 6.2 of Sutton and Barto [1].

The variance of Monte Carlo evaluation is relatively higher than TD learning because in Monte Carlo evaluation, we consider many transitions in each episode with each transition contributing variance to our estimate. On the other hand, TD learning only considers one transition per update, so we do not accumulate variance as quickly.

Finally, Monte Carlo is generally more data efficient than TD(0). In Monte Carlo, we update the value of a state based on the returns of the entire episode, so if there are highly positive or negative rewards many trajectories in the future, these rewards will still be immediately incorporated into our update of the value of the state. On the other hand, in TD(0), we update the value of a state using only the reward in the current step and some previous estimate of the value at the next state. This means that if there are highly positive or negative rewards many trajectories in the future, we will only incorporate these into the current state's value update when that reward has been used to update the bootstrap estimate of the next state's value. This means that if a highly rewarding episode has length  $L$ , then we may need to experience that episode  $L$  times for the information of the highly rewarding episode to travel all the way back to the starting state.

In Table 1, we summarize the strengths and limitations of each method discussed here.

	Dynamic Programming	Monte Carlo	Temporal Difference
Model Free?	No	Yes	Yes
Non-episodic domains?	Yes	No	Yes
Non-Markovian domains?	No	Yes	No
Converges to true value	Yes	Yes	Yes
Unbiased Estimate	N/A	Yes	No
Variance	N/A	High	Low
Data efficient	No	Yes	No

Table 1: Summary of methods in this lecture.

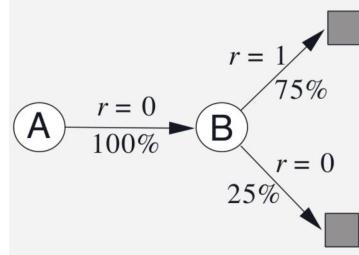


Figure 5: Example 6.4 from Sutton and Barto [1].

## 4.7 Batch Monte Carlo and Temporal Difference

We now look at the batch versions of the algorithms in today's lecture, where we have a set of histories that we use to make updates many times. Before looking at the batch cases in generality, let's first look at Example 6.4 from Sutton and Barto [1] to more closely examine the difference between Monte Carlo and TD(0). Suppose  $\gamma = 1$  and we have eight histories generated by policy  $\pi$ , take action  $act1$  in all states:

$$\begin{aligned} h_1 &= (A, act1, +0, B, act1, +0, \text{terminal}) \\ h_j &= (B, act1, +1, \text{terminal}) \text{ for } j = 2, \dots, 7 \\ h_8 &= (B, act1, +0, \text{terminal}). \end{aligned}$$

Then, using either batch Monte Carlo or TD(0) with  $\alpha = 1$ , we see that  $V(B) = 0.75$ . However, using Monte Carlo, we get that  $V(A) = 0$  since only the first episode visits state  $A$  and has return 0. On the other hand, TD(0) gives us  $V(A) = 0.75$  because we perform the update  $V(A) \leftarrow r_{1,1} + \gamma V(B)$ . Under a Markovian domain like the one shown in Figure 5, the estimate given by TD(0) makes more sense.

In this section, we consider the batch cases of Monte Carlo and TD(0). In the batch case, we are given a batch, or set of histories  $h_1, \dots, h_n$ , which we then feed through Monte Carlo or TD(0) many times.

The only difference from our formulations before is that we only update the value function after each time we process the entire batch. Thus in TD(0), the bootstrap estimate is updated only after each pass through the batch. *avg over entire batch, not just 1 history / episode.*

- 1) In the Monte Carlo batch setting, the value at each state converges to the value that minimizes the mean squared error with the observed returns. This follows directly from the fact that in Monte Carlo, we take an average over returns at each state, and in general, the MSE minimizer of samples is precisely the average of the samples. That is, given samples  $y_1, \dots, y_n$ , the value  $\sum_{i=1}^n (y_i - \hat{y})^2$  is minimized for  $\hat{y} = \frac{1}{n} \sum_{i=1}^n y_i$ . We can also see this in the example at the beginning of the section. We get that  $V(A) = 0$  for Monte Carlo because this is the only history visiting state  $A$ .

- 2) In the TD(0) batch setting, we do not converge to the same result as in Monte Carlo. In this case, we

converge to the value  $V^\pi$  that is the value of policy  $\pi$  on the maximum likelihood MDP model where

$$\hat{P}(s'|s, a) = \frac{1}{N(s, a)} \sum_{j=1}^n \sum_{t=1}^{L_j-1} \mathbb{1}(s_{j,t} = s, a_{j,t}, s_{j,t+1} = s') \quad (22)$$

$$\hat{r}(s, a) = \frac{1}{N(s, a)} \sum_{j=1}^n \sum_{t=1}^{L_j-1} \mathbb{1}(s_{j,t} = s, a_{j,t}) r_{j,t}. \quad (23)$$

Kinda like MC, simply taking avg

In other words, the maximum likelihood MDP model is the most naive model we can create based on the batch - the transition probability  $\hat{P}(s'|s, a)$  is the fraction of times that we see the transition  $(s, a, s')$  after we take action  $a$  at state  $s$  in the batch, and the reward  $\hat{r}(s, a)$  is the average reward experienced after taking action  $a$  at state  $s$  in the batch.

We also see this result in the example from the beginning of the section. In this case, our maximum likelihood model is

$$\hat{P}(B|A, act1) = 1 \quad (24)$$

$$\hat{P}(terminal|B, act1) = 1 \quad (25)$$

$$\hat{r}(A, act1) = 0 \quad (26)$$

$$\hat{r}(B, act1) = 0.75. \quad (27)$$

This gives us  $V^\pi(A) = 0.75$ , like we stated before.

The value function derived from the maximum likelihood MDP model is known as the **certainty equivalence estimate**. Using this relationship, we have another method for evaluating the policy. We can first compute the maximum likelihood MDP model using the batch. Then we can compute  $V^\pi$  using this model and the model-based policy evaluation methods discussed in last lecture. This method is highly data efficient but is computationally expensive because it involves solving the MDP which takes time  $O(|S|^3)$  analytically and  $(|S|^2|A|)$  via dynamic programming.

**Exercise 4.3.** Consider again the Mars Rover example in Figure 3. Suppose that our estimate for the value of each state is currently 0. If our batch consists of two histories

$$h_1 = (S3, TL, +0, S2, TL, +0, S1, TL, +1, terminal)$$

$$h_2 = (S3, TL, +0, S2, TL, +0, S2, TL, +0, S1, TL, +1, terminal)$$

and our policy is  $TL$ , then what is the certainty equivalence estimate?

## References

- [1] Sutton, Richard S. and Andrew G. Barto. *Introduction to Reinforcement Learning*. 2nd ed., MIT Press, 2017. Draft. <http://incompleteideas.net/book/the-book-2nd.html>.

$$\begin{aligned} \hat{p}(S_2|S_3, TL) &= 1 & \hat{r}(S_3, TL) &= 0 & V^\pi(S) &= V^\pi(S) + \alpha(R(S) + \gamma V^\pi(S') - V^\pi(S)) \\ \hat{p}(S_2|S_2, TL) &= \frac{1}{3} & \hat{r}(S_2, TL) &= 0 & V^\pi(S_1) &= 0 + \alpha(1 + 0 - 0) = \alpha \\ \hat{p}(S_1|S_2, TL) &= \frac{2}{3} & \hat{r}(S_1, TL) &= 1 & V^\pi(S_2) &= 0 + \alpha(0 + \frac{2}{3} \gamma \alpha - 0) = \frac{2}{3} \gamma \alpha^2 \\ \hat{p}(\text{terminal}|S_1, TL) &= 1 & & & V^\pi(S_3) &= 0 + \alpha(0 + \gamma \cdot \frac{2}{3} \gamma \alpha^2 - 0) \\ & & & & & = \frac{2}{3} \gamma^2 \alpha^2 \end{aligned}$$

$$\text{solution is } V^\pi(S_1) = 1 \quad V^\pi(S_2) = \frac{2\gamma}{3-\gamma} \quad V^\pi(S_3) = \frac{2\gamma^2}{3-\gamma}$$