# Policy Iteration & Monte Carlo Prediction

CMPUT 366: Intelligent Systems

S&B §4.3-4.4, 5.0-5.2

# Lecture Outline

1. Recap & Logistics

2. Policy Iteration

3. Monte Carlo Prediction

# Assignment #3

- Assignment #3 is due **Mar 29 (next Monday)** at 11:59pm

- Reminder that TAs are available during office hours 5 days/week to help

- **TensorFlow tutorial** in today's office hour:

  - Wednesday Mar 24 at **2:00pm**

  - see eClass for Google Meet link

# Recap: Value Functions

**State-value function**

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

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

**Action-value function**

$$q_\pi(s, a) \doteq \mathbb{E}_\pi[G_t \,|\, S_t = s, A_t = a]$$

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

# Recap: Bellman Equations

Value functions satisfy a **recursive consistency condition** called the **Bellman equation**:

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

$$= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \,|\, S_t = s]$$

$$= \sum_a \pi(a\,|\,s) \sum_{s'} \sum_r p(s', r\,|\,s, a)\big[r + \gamma\mathbb{E}_\pi[G_{t+1}\,|\,S_{t+1} = s']\big]$$

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

- $v_\pi$ is the **unique solution** to $\pi$'s (state-value) Bellman equation

- There is also a Bellman equation for $\pi$'s **action-value function**

# In-Place Iterative Policy Evaluation

> **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 $V(terminal) = 0$
>
> Loop:
>     $\Delta \leftarrow 0$
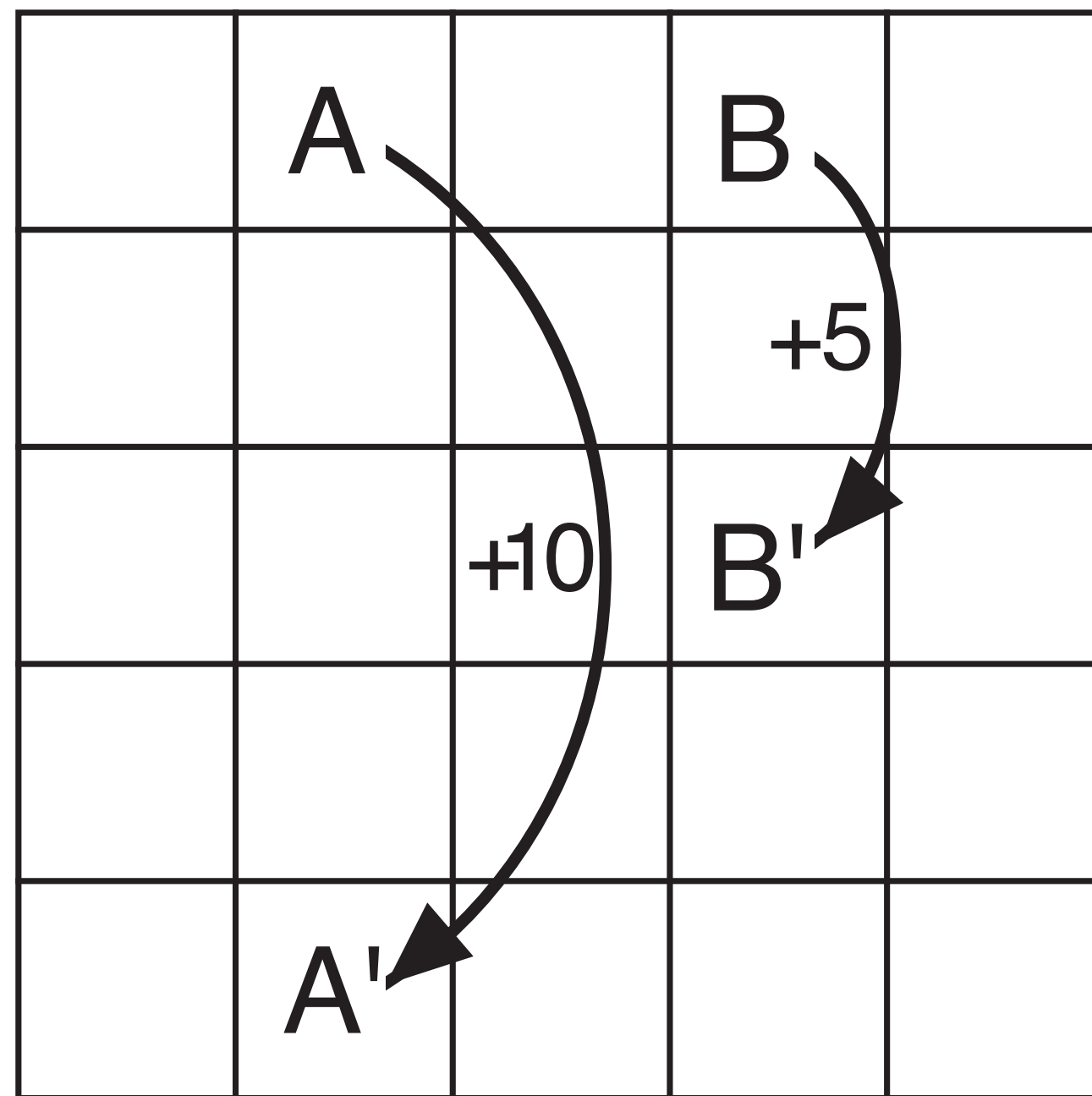>     Loop for each $s \in \mathcal{S}$:
>         $v \leftarrow V(s)$
>         $V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r\,|\,s,a)\big[r + \gamma V(s')\big]$
>         $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
> until $\Delta < \theta$

- The updates are **in-place**: we use new values for $V(s)$ **immediately** instead of waiting for the current sweep to complete (**why?**)

- These are **expected updates**: Based on a weighted average (expectation) of **all possible next states** (**instead of what?**)
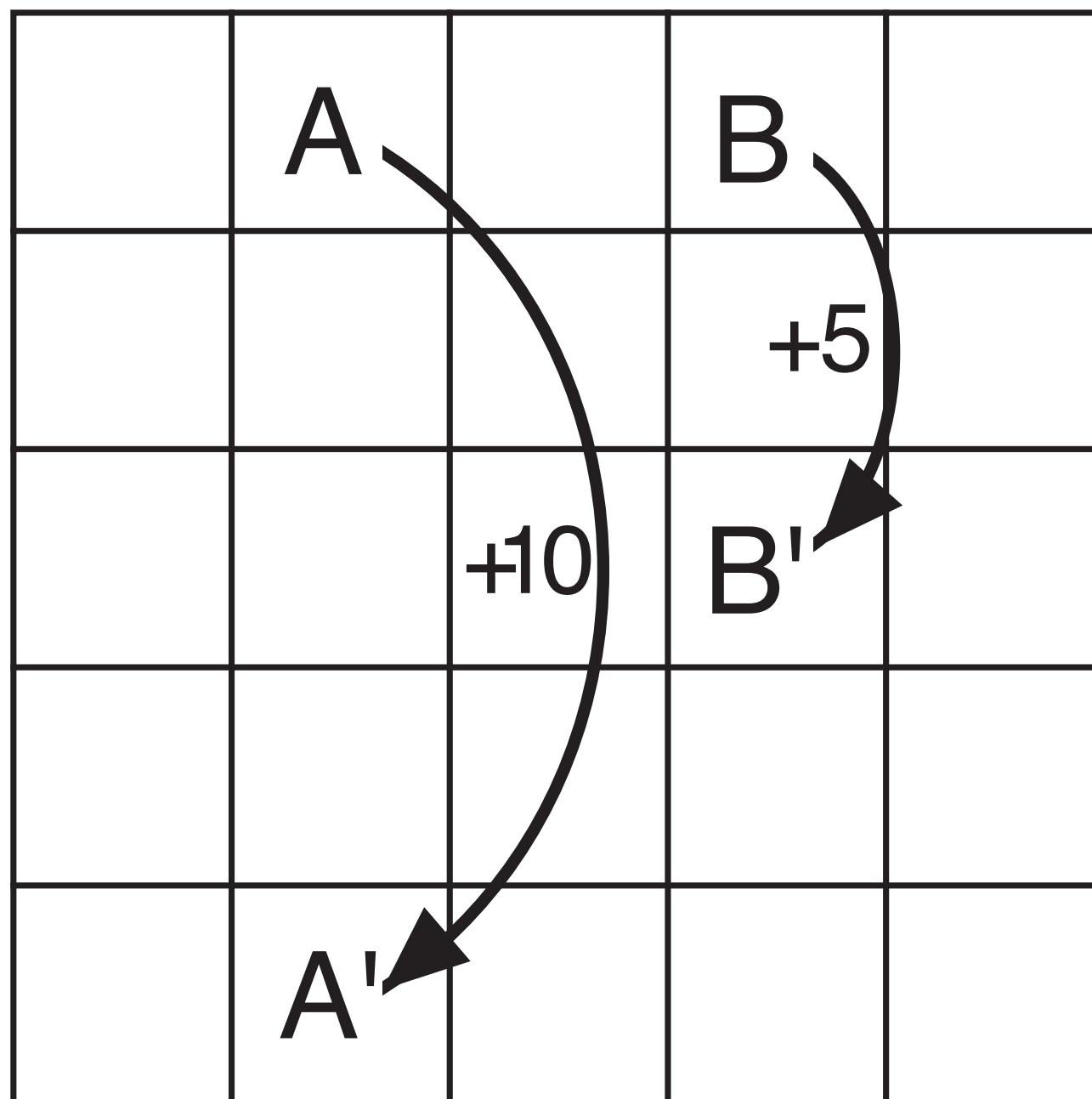
# Iterative Policy Evaluation



Reward dynamics

$V$ at $k = 0$

# Iterative Policy Evaluation
# in GridWorld

| | | | | |
|---|---|---|---|---|
| | A | | B | |
| | | | +5 | |
| | +10 | B' | | |
| | | | | |
| | A' | | | |

Reward dynamics

| -0.5 | 10 | 2 | 5 | 0.6 |
|---|---|---|---|---|
| -0.3 | 2.1 | 0.9 | 1.3 | 0.2 |
| -0.3 | 0.4 | 0.3 | 0.4 | -0.1 |
| -0.3 | 0.0 | 0.0 | 0.1 | -0.2 |
| -0.5 | -0.3 | -0.3 | -0.3 | -0.6 |

$V$ at $k = 1$

# Iterative Policy Evaluation in GridWorld



Reward dynamics

| 1.4 | 9.7 | 3.7 | 5.3 | 1.0 |
|------|------|------|------|------|
| 0.4 | 2.5 | 1.8 | 1.7 | 0.4 |
| -0.2 | 0.6 | 0.6 | 0.5 | -0.1 |
| -0.5 | 0.0 | 0.0 | 0.0 | -0.5 |
| -1.0 | -0.6 | -0.5 | -0.5 | -1.0 |

$V$ at $k = 2$

# Iterative Policy Evaluation in GridWorld



Reward dynamics

| | | | | |
|---|---|---|---|---|
| 3.4 | 8.9 | 4.5 | 5.3 | 1.5 |
| 1.6 | 3.0 | 2.3 | 1.9 | 0.6 |
| 0.1 | 0.8 | 0.7 | 0.4 | -0.4 |
| -1.0 | -0.4 | -0.3 | -0.6 | -1.2 |
| -1.9 | -1.3 | -1.2 | -1.4 | -2.0 |

$V$ at $k = 10\,000$

# Policy Improvement Theorem

**Theorem:**

Let $\pi$ and $\pi'$ be any pair of deterministic policies.

If $q_\pi(s, \pi'(s)) \geq v_\pi(s) \quad \forall s \in \mathcal{S}$,

then $v_{\pi'}(s) \geq v_\pi(s) \quad \forall s \in \mathcal{S}$.

If you are never worse off **at any state** by following $\pi'$ for **one step** and then following $\pi$ forever after, then following $\pi'$ **forever** has a higher expected value **at every state**.

# Policy Improvement Theorem Proof

$$v_\pi(s) \leq q_\pi(s, \pi'(s))$$

# Policy Improvement Theorem Proof

$$v_\pi(s) \le q_\pi(s, \pi'(s))$$
$$= \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = \pi'(s)]$$
$$= \mathbb{E}_{\pi'}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s]$$
$$\le \mathbb{E}_{\pi'}[R_{t+1} + \gamma q_\pi(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s]$$
$$= \mathbb{E}_{\pi'}[R_{t+1} + \gamma \mathbb{E}_{\pi'}[R_{t+2} + \gamma v_\pi(S_{t+2}) | S_{t+1}, A_{t+1} = \pi'(S_{t+1})] \mid S_t = s]$$
$$= \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma R_{t+2} + \gamma^2 v_\pi(S_{t+2}) \mid S_t = s\right]$$
$$\le \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v_\pi(S_{t+3}) \mid S_t = s\right]$$
$$\vdots$$
$$\le \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \cdots \mid S_t = s\right]$$
$$= v_{\pi'}(s).$$

# Greedy Policy Improvement

Given any policy $\pi$, we can construct a new greedy policy $\pi'$ that is guaranteed to be **at least as good**:

$$\pi'(s) \doteq \arg\max_a q_\pi(s, a)$$

$$= \arg\max_a \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \,|\, S_t = s, A_t = a]$$

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

- If this new policy is **not better** than the old policy, then $v_\pi(s) = v_{\pi'}(s)$ for all $s \in \mathcal{S}$ (**why?**)

- Also means that the new (and old) policies are **optimal** (**why?**)

# Policy Iteration

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

**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:
   $\quad \Delta \leftarrow 0$
   $\quad$ Loop for each $s \in \mathcal{S}$:
   $\quad\quad v \leftarrow V(s)$
   $\quad\quad V(s) \leftarrow \sum_{s',r} p(s',r \,|\, s, \pi(s))\big[r + \gamma V(s')\big]$
   $\quad\quad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
   until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement
   $policy\text{-}stable \leftarrow true$
   For each $s \in \mathcal{S}$:
   $\quad old\text{-}action \leftarrow \pi(s)$
   $\quad \pi(s) \leftarrow \arg\max_a \sum_{s',r} p(s',r \,|\, s, a)\big[r + \gamma V(s')\big]$
   $\quad$ If $old\text{-}action \neq \pi(s)$, then $policy\text{-}stable \leftarrow false$
   If $policy\text{-}stable$, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

This is a lot of iterations! Is it necessary to run to completion?

# Value Iteration

**Value iteration** <span style="color:red">interleaves</span> the estimation and improvement steps:

$$v_{k+1}(s) \doteq \max_a \mathbb{E}\left[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s, A_t = a\right]$$

$$= \max_a \sum_{s',r} p(s', r \mid s, a)\left[r + \gamma v_k(s')\right]$$

---

### 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 $V(terminal) = 0$

Loop:
|    $\Delta \leftarrow 0$
|    Loop for each $s \in \mathcal{S}$:
|       $v \leftarrow V(s)$
|       $V(s) \leftarrow \max_a \sum_{s',r} p(s', r \mid s, a)\left[r + \gamma V(s')\right]$
|       $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that
    $\pi(s) = \arg\max_a \sum_{s',r} p(s', r \mid s, a)\left[r + \gamma V(s')\right]$

# Policy Iteration Summary

- An **optimal policy** has higher state value than any other policy **at every state**

- A policy's state-value function can be computed by **iterating** an **expected update** based on the Bellman equation

- Given any policy $\pi$, we can compute a **greedy improvement** $\pi'$ by choosing highest expected value action based on $v_\pi$

- **Policy iteration:** Repeat:
  Greedy improvement using $v_\pi$, then recompute $v_\pi$

- **Value iteration:** Repeat:
  Recompute $v_\pi$ by assuming greedy improvement at every update

# Example: Blackjack

- Player gets two cards, dealer gets 1

- Player can `hit` (get a new card) as many times as they like, or `stick` (stop hitting)

- After the player is done, the dealer hits / sticks according to a fixed rule

- Whoever has the most points (sum of card values) wins

- But, if you have more than 21 points, you **lose immediately** ("bust")

# Simulating Blackjack

- Given a policy for the player, it is **very easy** to simulate a game of Blackjack

- **Question:** Is it easy to **compute** the full **dynamics**?

- **Question:** Is it easy to run **iterative policy evaluation**?

# Experience vs. Expectation

- In order to compute **expected updates**, we need to know the exact **probability** of **every** possible transition

- Often we don't have access to the full probability distribution, but we do have access to **samples of experience**

  1. **Actual experience:** We want to learn based on interactions with a **real environment**, without knowing its dynamics

  2. **Simulated experience:** We can **simulate** the dynamics, but we don't have an **explicit representation** of transition probabilities, or there are **too many states**

# Monte Carlo Estimation

- Instead of estimating expectations by a **weighted sum** over **all possibilities**, estimate expectation by **averaging** over a **sample** drawn from the distribution:

$$\mathbb{E}[X] = \sum_{x} f(x)x \approx \frac{1}{n} \sum_{i=1}^{n} x_i \quad \text{where } x_i \sim f$$

# Monte Carlo Prediction

- Use a **large sample** of **episodes** generated by a policy $\pi$ to estimate the state-values $v_\pi(s)$ for each state $s$

  - We will consider only **episodic** tasks for now

- **Question:** What is the **return** $G_t$ for state $S_t = s$ in a given episode?

- We can estimate the expected return $v_\pi(s) = \mathbb{E}[G_t \mid S_t = s]$ by averaging the returns for that state in every episode containing a visit to $s$

# First-visit Monte Carlo Prediction

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

Input: a policy $\pi$ to be evaluated

Initialize:
$\quad V(s) \in \mathbb{R}$, arbitrarily, for all $s \in \mathcal{S}$
$\quad Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

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

# Monte Carlo vs. Dynamic Programming

- **Iterative policy evaluation** uses the estimates of the **next state's** value to update the value of this state

  - Only needs to compute a **single transition** to update a state's estimate

- **Monte Carlo** estimate of each state's value is **independent** from estimates of **other states'** values

  - Needs the **entire episode** to compute an update

  - Can focus on evaluating a **subset of states** if desired

# Summary

**Monte Carlo estimation** estimates values by averaging returns over **sample episodes**

- Does not require access to full model of **dynamics**

- Does require access to an entire **episode** for each sample