

# Monte Carlo Prediction

CMPUT 366: Intelligent Systems

S&B §4.3-4.4, 5.0-5.2

# Lecture Outline

1. Recap
2. Policy Iteration
3. Monte Carlo Prediction

# Recap: 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  $V(\text{terminal}) = 0$

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  $\Delta < \theta$

- 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?**)

# Recap: 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**.

# Recap: Policy Iteration

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

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:

$\Delta \leftarrow 0$

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

$policy\text{-}stable \leftarrow true$

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

$old\text{-}action \leftarrow \pi(s)$

$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

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

# 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

- **Question:** What was **Monte Carlo estimation** the last time we studied it (in Supervised Learning?)
- 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:

$V(s) \in \mathbb{R}$ , 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,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + R_{t+1}$

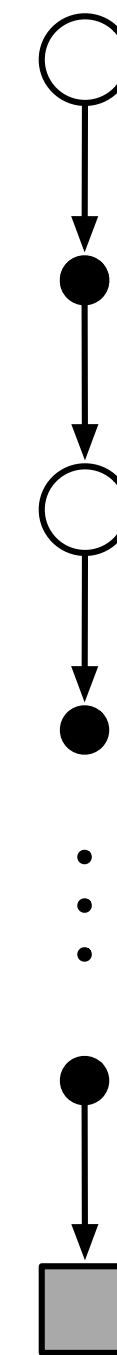
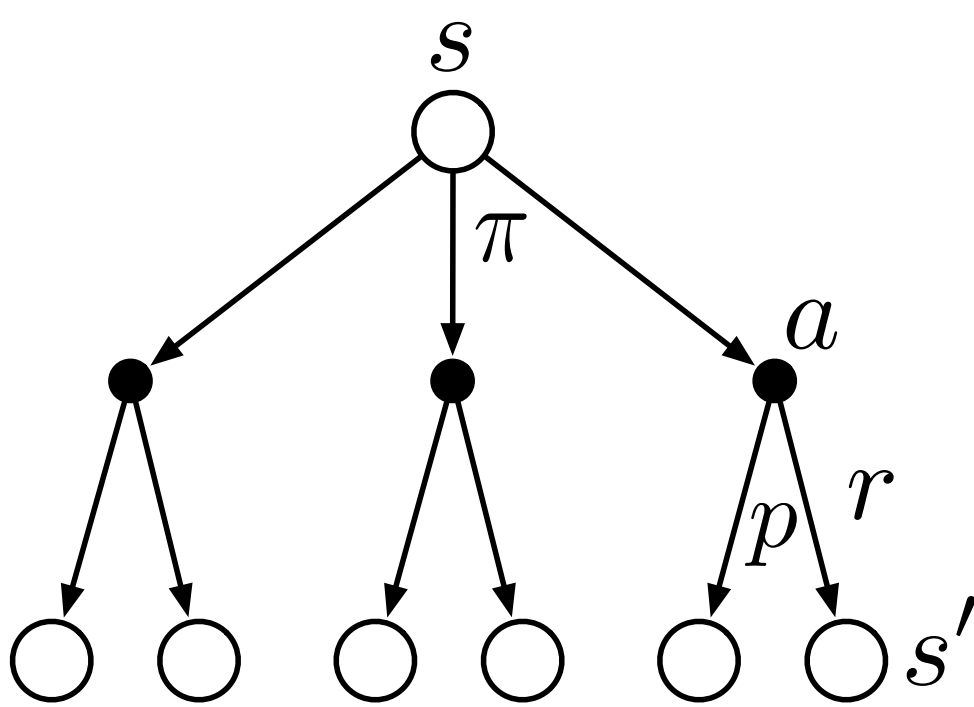
Unless  $S_t$  appears in  $S_0, S_1, \dots, S_{t-1}$ :

Append  $G$  to  $Returns(S_t)$

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

# 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