

# **ST455: Reinforcement Learning**

## **Lecture 3: Elementary Solution Methods Dynamic Programming and Monte Carlo**

Chengchun Shi

# Lecture Outline

---

## 1. Preliminaries

## 2. Dynamic Programming

2.1 Policy Iteration

2.2 Value Iteration

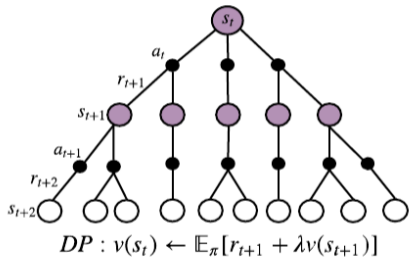
## 3. Monte Carlo Methods

3.1 MC Policy Evaluation (Prediction)

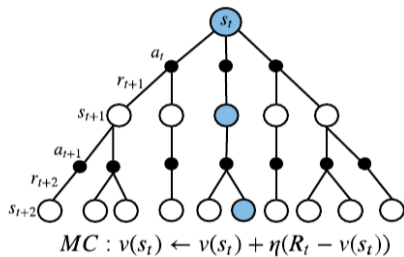
3.2 MC Policy Optimization (Control)

# Lecture Outline (Cont'd)

---



Dynamic Programming (DP)



Monte Carlo (MC)

# 1. Preliminaries

## 2. Dynamic Programming

2.1 Policy Iteration

2.2 Value Iteration

## 3. Monte Carlo Methods

3.1 MC Policy Evaluation (Prediction)

3.2 MC Policy Optimization (Control)

# Learning v.s. Planning

---

Two fundamental problems in sequential decision making

- **Planning**

- A model of the environment (e.g., state transition, reward function) is **known**
- The agent performs computations with its model, **without** any external interaction
- a.k.a. deliberation, reasoning, introspection, pondering, thought, search
- Example: **Dynamic Programming**

- **Learning**

- The environment is initially **unknown**
- The agent **interacts** with the model
- The agent **learns** the optimal policy from experience
- Example: **Monte Carlo methods, temporal difference learning, policy-based learning, model-based learning**

# Example: Go Game

---

- **Planning:** Rules of Go are known
- Exhaustive search of the optimal move
- No need to play Go with others



- **Learning:** No need to know the rules
- Learn the optimal move from experience
- Practice makes perfect



# Models: Finite MDPs

---

- Environment modelled by a finite MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
- MDP model assumption: **Markovianity** & **time-homogeneity**
- $\mathcal{S}$ : state space (a **finite** set of states)
- $\mathcal{A}$ : action space (a **finite** set of actions)
- $\mathcal{P}$ : state transition probability matrix,  $\mathcal{P}_{ss'}^a = \Pr(\mathbf{S}_{t+1} = s' | \mathbf{A}_t = a, \mathbf{S}_t = s)$
- $\mathcal{R}$ : reward function,  $\mathcal{R}_s^a = \mathbb{E}(\mathbf{R}_t | \mathbf{A}_t = a, \mathbf{S}_t = s)$
- $\gamma$ : discounted factor  $\in [0, 1]$ , allowed to be **1** if all sequences terminate (e.g., finite horizons)
- Dynamic Programming (DP) and Monte Carlo methods (MC) are **equally applicable** to settings with continuous state or action space

# Bellman Equations

---

- Bellman equation for the (state) value function:

$$V^\pi(s) = \mathbb{E}^\pi[R_t + \gamma V^\pi(S_{t+1}) | S_t = s],$$

- or equivalently,

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V^\pi(s') \right].$$

- Bellman optimality equation for the **optimal** value function:

$$V^{\pi^{\text{opt}}}(s) = \max_a \mathbb{E}[R_t + \gamma V^{\pi^{\text{opt}}}(S_{t+1}) | A_t = a, S_t = s],$$

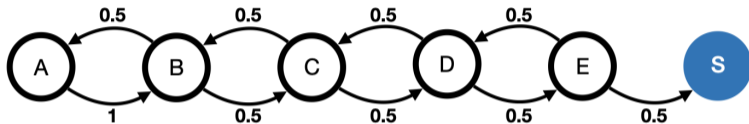
- or equivalently,

$$V^{\pi^{\text{opt}}}(s) = \max_{a \in \mathcal{A}} \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V^{\pi^{\text{opt}}}(s') \right].$$



# Bellman Equation: The Random Walk Example

- Consider a simple **random walk** on a path:



- Reward for transition to State **S** of value **1**, zero reward for other transitions
- Bellman equations:

$$V^\pi(\mathbf{A}) = \mathbb{E}^\pi[R_t + \gamma V^\pi(\mathbf{S}_{t+1}) | \mathbf{S}_t = \mathbf{A}] = \gamma V^\pi(\mathbf{B})$$

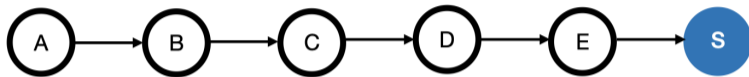
$$V^\pi(\mathbf{B}) = \mathbb{E}^\pi[R_t + \gamma V^\pi(\mathbf{S}_{t+1}) | \mathbf{S}_t = \mathbf{B}] = \frac{\gamma}{2} V^\pi(\mathbf{C}) + \frac{\gamma}{2} V^\pi(\mathbf{A})$$

⋮

$$V^\pi(\mathbf{S}) = \mathbb{E}^\pi[R_t + \gamma V^\pi(\mathbf{S}_{t+1}) | \mathbf{S}_t = \mathbf{S}] = 1$$

# Bellman Optimality Equation: Random Walk

- The **random walk** example:



- Reward for transition to State **S** of value **1**, zero reward for other transitions
- Bellman optimality equations:

$$V^{\pi^{\text{opt}}}(\mathbf{A}) = \max_a \mathbb{E}[R_t + \gamma V^{\pi^{\text{opt}}}(\mathbf{S}_{t+1}) | \mathbf{A}_t = \mathbf{a}, \mathbf{S}_t = \mathbf{A}] = \gamma V^{\pi^{\text{opt}}}(\mathbf{B})$$

$$V^{\pi^{\text{opt}}}(\mathbf{B}) = \max_a \mathbb{E}[R_t + \gamma V^{\pi^{\text{opt}}}(\mathbf{S}_{t+1}) | \mathbf{A}_t = \mathbf{a}, \mathbf{S}_t = \mathbf{B}] = \gamma V^{\pi^{\text{opt}}}(\mathbf{C})$$

⋮

$$V^{\pi^{\text{opt}}}(\mathbf{S}) = \max_a \mathbb{E}[R_t + \gamma V^{\pi^{\text{opt}}}(\mathbf{S}_{t+1}) | \mathbf{A}_t = \mathbf{a}, \mathbf{S}_t = \mathbf{S}] = 1$$

# State-Action Value Function

## Definition

The state-action value function (better known as the **Q-function**) is expected return starting from  $s$  and  $a$  under  $\pi$ ,

$$Q^\pi(s, a) = \mathbb{E}^\pi(G_t | A_t = a, S_t = s) = \mathbb{E}^\pi \left( \sum_{i=0}^{+\infty} \gamma^i R_{i+t} | A_t = a, S_t = s \right).$$

- $Q^\pi$  is **independent** of the time  $t$  in its definition, under **time-homogeneity**
- $Q^\pi$  is the state value  $V^\pi$  under a Markov policy that implements  $a$  at the first time and follows  $\pi$  afterwards
- Reduces to action value function  $\mathbb{E}^\pi(R_t | A_t = a)$  in Lecture 1 when  $\gamma = 0$ ,  $S = \emptyset$

# State-Action Value Function (Cont'd)

---

Relationships between  $V^\pi$  and  $Q^\pi$

- $Q^\pi \rightarrow V^\pi$ :

$$V^\pi(s) = \mathbb{E}^\pi(G_t | S_t = s) = \sum_{a \in \mathcal{A}} \pi(a|s) \mathbb{E}^\pi(G_t | A_t = a, S_t = s) = \sum_{a \in \mathcal{A}} \pi(a|s) Q^\pi(s, a)$$

- $V^\pi \rightarrow Q^\pi$ :

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E}(R_t | A_t = a, S_t = s) + \gamma \mathbb{E}(G_{t+1} | A_t = a, S_t = s) \\ &= \mathbb{E}(R_t | A_t = a, S_t = s) + \gamma \mathbb{E}[\mathbb{E}^\pi(G_{t+1} | S_{t+1}) | A_t = a, S_t = s] \\ &= \mathbb{E}[R_t + \gamma V^\pi(S_{t+1}) | A_t = a, S_t = s] \end{aligned}$$

# 1. Preliminaries

## 2. Dynamic Programming

- 2.1 Policy Iteration
- 2.2 Value Iteration

## 3. Monte Carlo Methods

- 3.1 MC Policy Evaluation (Prediction)
- 3.2 MC Policy Optimization (Control)

# Dynamic Programming

---

## Definition (Dynamic Programming)

A collection of algorithms used to compute optimal policies given **perfect** knowledge of the environment

- **Dynamic**: sequential or temporal component to the problem
- **Programming**: optimise a “program”, i.e., a policy
- Dynamic programming (DP) is **rarely** used in practice (the environment is usually unknown)
- However, they provide a foundation for other solution methods

# Dynamic Programming (Cont'd)

---

“Dynamic programming” is used to solve many other statistical learning problems

- Learning optimal **dynamic treatment regimes** (DTRs)
- Multi-scale **change point detection**
- De Boor algorithm for evaluating **B-spline** basis functions

Also used in bioinformatics, optimisation, control theory (see [wiki page](#))

# Dynamic Programming Methods

---

- **Policy Iteration**: an iterative method that alternates between
  - **Policy Evaluation**
  - **Policy Improvement**

$$\pi_0 \longrightarrow V^{\pi_0} \longrightarrow \pi_1 \longrightarrow V^{\pi_1} \longrightarrow \dots \longrightarrow \pi^{opt} \longrightarrow V^{\pi^{opt}}$$

- **Value Iteration**: simultaneously combine **policy evaluation** and **policy improvement**

$$V^{\pi_0} \longrightarrow V^{\pi_1} \longrightarrow V^{\pi_2} \longrightarrow \dots \longrightarrow V^{\pi^{opt}} \longrightarrow \pi^{opt}$$



# 1. Preliminaries

## 2. Dynamic Programming

2.1 Policy Iteration

2.2 Value Iteration

## 3. Monte Carlo Methods

3.1 MC Policy Evaluation (Prediction)

3.2 MC Policy Optimization (Control)

# Policy Iteration: Policy Evaluation

- Computation of the (state) value function  $\mathbf{V}^\pi$  for a given  $\pi$
- According to the Bellman equation, for any  $s$ ,

$$\mathbf{V}^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a \mathbf{V}^\pi(s') \right],$$

- written in matrix form,  $\mathbf{V}^\pi = \mathcal{R} + \gamma \mathcal{P} \mathbf{V}^\pi$
- $\mathbf{V}^\pi$  is a column vector with one entry per state

$$\begin{bmatrix} \mathbf{V}^\pi(1) \\ \vdots \\ \mathbf{V}^\pi(n) \end{bmatrix} = \sum_{a \in \mathcal{A}} \begin{bmatrix} \pi(a|1) \mathcal{R}_1^a \\ \vdots \\ \pi(a|n) \mathcal{R}_n^a \end{bmatrix} + \gamma \begin{bmatrix} \mathcal{P}_{11} & \cdots & \mathcal{P}_{1n} \\ \vdots & & \vdots \\ \mathcal{P}_{n1} & \cdots & \mathcal{P}_{nn} \end{bmatrix} \begin{bmatrix} \mathbf{V}^\pi(1) \\ \vdots \\ \mathbf{V}^\pi(n) \end{bmatrix},$$

where  $\mathcal{P}_{ij} = \sum_{a \in \mathcal{A}} \pi(a|i) \mathcal{P}_{ij}^a$

# Policy Evaluation (Cont'd)

---

- $\mathbf{V}^\pi$  is a solution of a system of  $n$  linear equations with  $n$  unknowns
- It can be computed directly

$$\begin{aligned}\mathbf{V}^\pi &= \mathcal{R} + \gamma \mathcal{P} \mathbf{V}^\pi \\ (I - \gamma \mathcal{P}) \mathbf{V}^\pi &= \mathcal{R} \\ \mathbf{V}^\pi &= (I - \gamma \mathcal{P})^{-1} \mathcal{R}\end{aligned}$$

- $I - \gamma \mathcal{P}$  is **invertible** when  $\gamma$  is strictly smaller than  $\mathbf{1}$ , since

$$\mathbf{x}^\top (I - \gamma \mathcal{P}) \mathbf{x} = (\mathbf{1} - \gamma) \|\mathbf{x}\|_2^2 + \gamma \sum_{ij} \mathcal{P}_{ij} (\mathbf{x}_i - \mathbf{x}_j)^2 > \mathbf{0},$$

when  $\mathbf{x} \neq \mathbf{0}$ . The equality holds due to that each row of  $\mathcal{P}$  sums up to  $\mathbf{1}$ .

# Policy Evaluation: Algorithm

---

- **Iterative Policy Evaluation:** an iterative method that outputs a sequence of value functions  $V_0, V_1, V_2, \dots, V_k \rightarrow V^\pi$
- **Initial** value function  $V_0$  is chosen arbitrarily subject to the **constraint** that at terminal state it has value  $0$
- **Iterative** update rule (according to the Bellman equation):

$$V_{k+1} = \mathcal{R} + \gamma \mathcal{P} V_k$$

- **Convergence** is guaranteed when  $\gamma$  is strictly smaller than  $1$  (more in appendix), or eventual termination is guaranteed from all states under  $\pi$

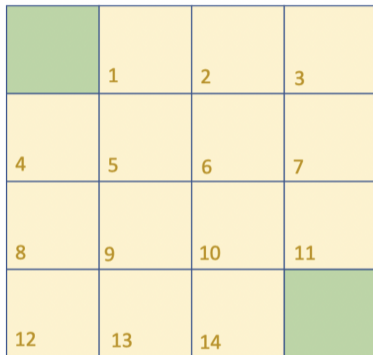
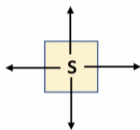
# Policy Evaluation: Pseudocode

---

- **Input:** a policy  $\pi$ , a threshold parameter  $\epsilon > 0$
- **Initialization:**  $V(s) = 0$  for any  $s \in \mathcal{S}$
- **Repeat:**
  - $\Delta \leftarrow 0$
  - For each**  $s \in \mathcal{S}$ 
    - $v \leftarrow V(s)$
    - $V(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V(s') \right]$
    - $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
  - until**  $\Delta < \epsilon$
- **Output**  $V$

# GridWorld Example

---



terminal state

- Undiscounted, episodic, finite MDP task
- $\mathcal{A} = \{\text{up, down, right, left}\}$ . Actions leading out of the grid leave state unchanged
- Rewards: for each transition, the reward of value  $-1$

# GridWorld Example (Cont'd)

---

0	1 -14	2 -20	3 -22
4 -14	5 -18	6 -20	7 -20
8 -20	9 -20	10 -18	11 -14
12 -22	13 -20	14 -14	0

Figure: Values of uniform random policy

$$\pi(\mathbf{n}|\cdot) = \pi(\mathbf{s}|\cdot) = \pi(\mathbf{w}|\cdot) = \pi(\mathbf{e}|\cdot) = 0.25$$

# GridWorld Example (Cont'd)

By **symmetry** and Bellman equation,

0	$v_1$	$v_2$	$v_3$
$v_1$	$v_5$	$v_6$	$v_2$
$v_2$	$v_6$	$v_5$	$v_1$
$v_3$	$v_2$	$v_1$	0

$$v_1 = \frac{1}{4}(-1 + v_2) + \frac{1}{4}(-1 + v_5) + \frac{1}{4}(-1 + 0) + \frac{1}{4}(-1 + v_1)$$

$$v_2 = \frac{1}{4}(-1 + v_3) + \frac{1}{4}(-1 + v_6) + \frac{1}{4}(-1 + v_1) + \frac{1}{4}(-1 + v_2)$$

$$v_3 = \frac{1}{4}(-1 + v_3) + \frac{1}{4}(-1 + v_2) + \frac{1}{4}(-1 + v_2) + \frac{1}{4}(-1 + v_3)$$

$$v_5 = \frac{1}{4}(-1 + v_6) + \frac{1}{4}(-1 + v_6) + \frac{1}{4}(-1 + v_1) + \frac{1}{4}(-1 + v_1)$$

$$v_6 = \frac{1}{4}(-1 + v_2) + \frac{1}{4}(-1 + v_5) + \frac{1}{4}(-1 + v_5) + \frac{1}{4}(-1 + v_2)$$

$$\Rightarrow (v_1, v_2, v_3, v_5, v_6) = (-14, -20, -22, -18, -20)$$



# GridWorld Example (Cont'd)

---

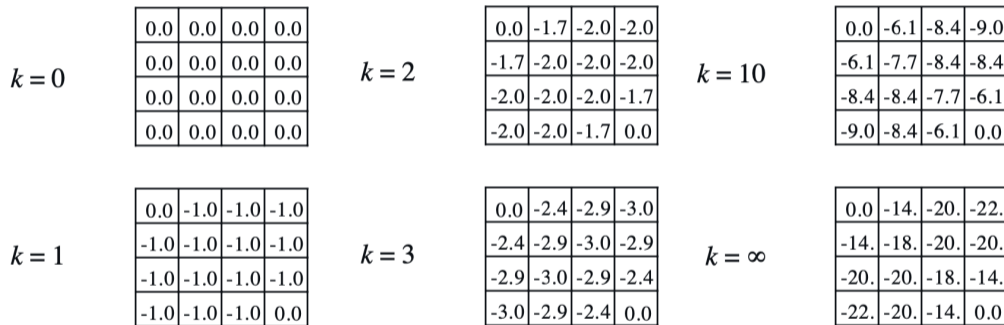


Figure: Value functions at each iteration

# Policy Iteration: Policy Improvement

---

- Identify some  $\pi'$  that is no worse than  $\pi$  based on  $V^\pi$
- For any  $s$ , consider a hybrid policy
  - implements  $a$  at the first time
  - follows  $\pi$  afterwards
- Its value is given by  $Q^\pi(s, a)$  (can be computed based on  $V^\pi$ )
- Select  $\pi'$  among the class of hybrid policies that **maximizes** the value

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

- Its value is given by  $Q^\pi(s, \pi'(s)) \geq V^\pi(s)$ , since the hybrid policy class contains  $\pi$
- Surprisingly, according to **policy improvement theorem**,  $V^{\pi'}(s) \geq V^\pi(s)$  for any  $s$ !

## Policy Improvement (Cont'd)

Given a policy  $\pi$ , improve  $\pi$  by acting **greedily** with respect to  $\mathbf{V}^\pi$ ,

$$\begin{aligned}\pi'(s) &= \arg \max_a Q^\pi(s, a) = \arg \max_a \mathbb{E}[R_t + \gamma V^\pi(S_{t+1}) | A_t = a, S_t = s] \\ &= \arg \max_a [R_s^a + \gamma \sum_{s'} P_{ss'}^a V^\pi(s')]\end{aligned}$$

### Theorem

The greedy policy  $\pi'$  with respect to  $\mathbf{V}^\pi$  is as good as or better than  $\pi$ ,

$$\mathbf{V}^{\pi'}(s) \geq \mathbf{V}^\pi(s),$$

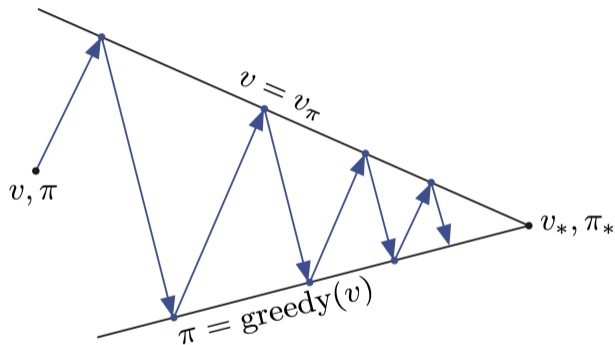
for any  $s \in \mathcal{S}$ .

Proof can be found in the Appendix.



# Policy Iteration: Revisit

---



- **Policy Evaluation:** Compute  $\mathbf{V}^\pi$  via iterative policy evaluation
- **Policy Improvement:** Generate  $\pi'$  via greedy policy improvement

# Policy Iteration: Pseudocode

---

- **Initialization:**  $V(s) = 0$ ,  $\pi(s) \in \mathcal{A}$  arbitrarily for any  $s \in \mathcal{S}$
- **Repeat:**
  - $\Delta \leftarrow 0$
  - For each**  $s \in \mathcal{S}$ 
    - $v \leftarrow V(s)$
    - $V(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V(s') \right]$
    - $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
  - until**  $\Delta < \epsilon$
- **polycystable**  $\leftarrow$  True
- **For each**  $s \in \mathcal{S}$ :
  - $b \leftarrow \pi(s)$
  - $\pi(s) \leftarrow \arg \max_a \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V(s') \right]$
  - If**  $b \neq \pi(s)$  **then** **polycystable**  $\leftarrow$  False
- **If** **polycystable**, **then Return**  $\pi$ , **else** go to bullet point #2

# 1. Preliminaries

## 2. Dynamic Programming

2.1 Policy Iteration

2.2 Value Iteration

## 3. Monte Carlo Methods

3.1 MC Policy Evaluation (Prediction)

3.2 MC Policy Optimization (Control)

# Value Iteration

---

- Policy iteration is **computationally inefficient**, as each iteration requires executing policy evaluation which requires multiple iterations
- According to the Bellman optimality equation,

$$V^{\pi^{\text{opt}}}(s) = \max_{a \in \mathcal{A}} \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V^{\pi^{\text{opt}}}(s') \right].$$

- **Value iteration** idea: iteratively apply the above updates

$$V_{k+1}(s) = \max_{a \in \mathcal{A}} \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V_k(s') \right].$$

- Drive the optimal deterministic policy

$$\pi^{\text{opt}}(s) = \arg \max_{a \in \mathcal{A}} \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V^{\pi^{\text{opt}}}(s') \right].$$

- **Convergence** is guaranteed when  $\gamma$  is strictly smaller than **1** (more in Appendix), or eventual termination is guaranteed from all states.



# Value Iteration: Pseudocode

---

- **Initialization:**  $V(s) = 0$ ,  $\pi(s) \in \mathcal{A}$  arbitrarily for any  $s \in \mathcal{S}$
- **Repeat:**
  - $\Delta \leftarrow 0$
  - For each**  $s \in \mathcal{S}$ 
    - $v \leftarrow V(s)$
    - $V(s) \leftarrow \max_{a \in \mathcal{A}} \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V(s') \right]$
    - $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
  - until**  $\Delta < \epsilon$
- **Output:** optimal deterministic policy given by

$$\pi^{\text{opt}}(s) = \arg \max_{a \in \mathcal{A}} \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V^{\pi^{\text{opt}}}(s') \right].$$

# Example: Gambler's Problem

---



- A gambler makes bets on the outcomes of a sequence of coin flips
- The gambler must decide for each coin flip what proportion of capital to stake
- If **the outcome of the coin flip = heads**, then:
  - The gambler **wins** as much money as they have staked on this flip
- Else:
  - The gambler **loses** their stake
- The game ends when the gambler reaches the goal of **£100** or **runs out of money**

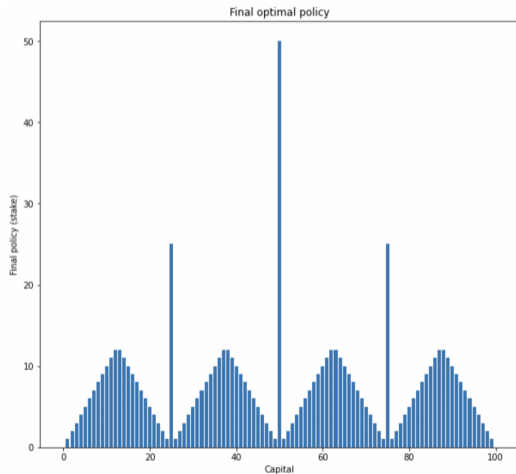
## Example: Gambler's Problem (Cont'd)

---

- Undiscounted, episodic, finite MDP task
- $\mathcal{S}$ :  $\{0, 1, \dots, 99, 100\}$ , termination states  $0$  and  $100$
- $\mathcal{A}(s)$ :  $\{1, 2, \dots, \min(s, 100 - s)\}$ , depends on the state
- $\Pr(\text{outcome of coin flip is heads}) = p$  (known parameter)
- Seminars:
  - Show the value function for different iterations
  - Show the optimal policy

# Example: Gambler's Problem, the Optimal Policy

---



# Some Technical Questions

---

- How do we know that value iteration converges to  $V^{\pi^{\text{opt}}}$ ?
- Or that iterative policy evaluation converges to  $V^{\pi}$ ?
- And therefore that policy iteration converges to  $V^{\pi^{\text{opt}}}$ ?
- Is the solution unique?
- These questions are resolved by **Banach fixed-point theorem** (or **contraction mapping theorem**), mentioned in Seminar 2 (more in the appendix)

# 1. Preliminaries

## 2. Dynamic Programming

2.1 Policy Iteration

2.2 Value Iteration

## 3. Monte Carlo Methods

3.1 MC Policy Evaluation (Prediction)

3.2 MC Policy Optimization (Control)

# Monte Carlo (MC) Methods

---

- Learning methods for solving the RL problem based on **averaging sample returns**
  - Estimating value functions and discovering optimal policies
  - Not assuming a model of the environment, based only on **experiences (model free)**
- Defined for **episodic** tasks
  - Value functions and policies are updated upon completion of an episode
  - Different from **step-by-step** methods (e.g., temporal difference learning)

# 1. Preliminaries

## 2. Dynamic Programming

2.1 Policy Iteration

2.2 Value Iteration

## 3. Monte Carlo Methods

3.1 MC Policy Evaluation (Prediction)

3.2 MC Policy Optimization (Control)



# MC Policy Evaluation

---

- **Objective:** estimate the value function  $V^\pi$  for a given policy  $\pi$ , from a set of episodes obtained by following  $\pi$

$$S_0, A_0, R_0, \dots, S_T \sim \pi$$

- $V^\pi$  is the expected return  $\mathbb{E}^\pi(\sum_{0 \leq t \leq T} \gamma^t R_t | S_0 = s)$
- Monte Carlo idea: use **empirical mean** return to approximate **expected return**
- Convergence is guaranteed by **law of large numbers**
- Types of MC methods:
  - **First-visit MC method:**  $V^\pi(s)$  estimated by the average of returns following **each first visit** to  $s$  in a set of episodes
  - **Every-visit MC method:**  $V^\pi(s)$  estimated by the average of returns following **each visit** to  $s$  in a set of episodes

# First-Visit MC Policy Evaluation: Pseudocode

---

- **Initialization:**

$N$  (counter),  $N(s) \leftarrow 0$  for all  $s \in \mathcal{S}$

$\text{Returns}(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$

- **Repeat:**

**Generate** an episode following policy  $\pi$

**For each** distinct  $s$  appearing in the episode

$G \leftarrow$  return following the first occurrence of  $s$

$N(s) \leftarrow N(s) + 1$

$\text{Returns}(s) \leftarrow \text{Returns}(s) + G$

- **Output:**

**For each** distinct  $s$

$N^{-1}(s)\text{Returns}(s)$

# Every-Visit MC Policy Evaluation: Pseudocode

---

- **Initialization:**

$N \leftarrow$  counter,  $N(s) \leftarrow \mathbf{0}$  for all  $s \in \mathcal{S}$

$\text{Returns}(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$

- **Repeat:**

**Generate** an episode following policy  $\pi$

**For each**  $s$  appearing in the episode

$G \leftarrow$  return following the occurrence of  $s$

$N(s) \leftarrow N(s) + \mathbf{1}$

$\text{Returns}(s) \leftarrow \text{Returns}(s) + G$

- **Output:**

**For each** distinct  $s$

$N^{-1}(s)\text{Returns}(s)$

# 1. Preliminaries

## 2. Dynamic Programming

2.1 Policy Iteration

2.2 Value Iteration

## 3. Monte Carlo Methods

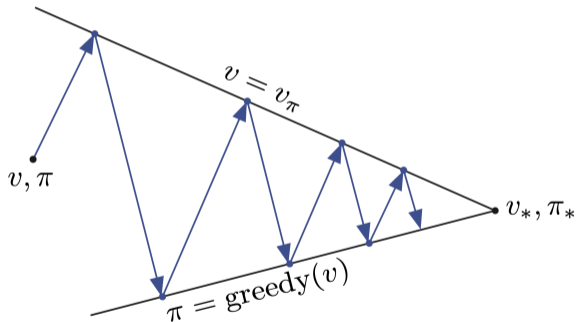
3.1 MC Policy Evaluation (Prediction)

3.2 MC Policy Optimization (Control)

# MC Control

---

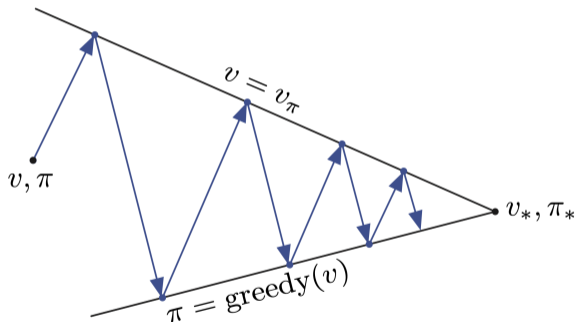
- **Objective:** use MC estimation to learn the optimal policy.
- Recall the policy iteration algorithm



- **Policy Evaluation:** Compute  $\mathbf{V}^\pi$  via iterative policy evaluation
- **Policy Improvement:** Generate  $\pi'$  via greedy policy improvement

# MC Control with Generalized Policy Iteration

- **Objective:** use MC estimation to learn the optimal policy.
- Integrate policy iteration with MC methods



- **Policy Evaluation:** Compute  $\mathbf{V}^\pi$  via MC policy evaluation
- **Policy Improvement:** Generate  $\pi'$  via greedy policy improvement?

# Policy Iteration Using State-Action Value Function

---

- Greedy policy improvement over  $V^\pi$  requires model of MDP

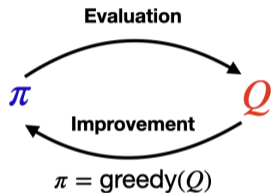
$$\pi'(s) = \arg \max_a [\mathcal{R}_s^a + \gamma \sum_{s'} P_{ss'}^a V^\pi(s')]$$

- Greedy policy improvement over  $Q^\pi(s, a)$  is model free

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

# MC Version of Policy Iteration

---



$$\pi_0 \rightarrow Q^{\pi_0} \rightarrow \pi_1 \rightarrow Q^{\pi_1} \rightarrow \dots \rightarrow \pi^{opt} \rightarrow Q^{\pi^{opt}}$$

- **Policy Evaluation:** MC estimation of state-action value function
- **Policy Improvement:** Improve the policy wrt the current state-action value function



# MC Estimation of State-Action Values

---

- Many state-action pairs may never be visited under a policy
  - Ex. if  $\pi$  is deterministic, only **one** state-action pair is observed for each distinct state
  - Need to ensure **exploration!**
- Two approaches for ensuring exploration:
  - **Exploring starts**: the first step of each episode starts at a state-action pair and every such pair has non-zero probability of being selected at the start
  - **Stochastic policies**: use policies that ensures a non-zero probability of selecting each action from the set of available actions in each given state

# MC Control with Exploring Starts

---

- Initialization:

$N$  (counter),  $N(s, a) \leftarrow 0$  for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}$

$\text{Returns}(s, a) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}$

$\pi \leftarrow$  arbitrary

$Q \leftarrow$  arbitrary

- Repeat:

**Generate** an episode using exploring starts and policy  $\pi$

**For each** distinct  $(s, a)$  appearing in the episode

$G \leftarrow$  return following the first occurrence of  $(s, a)$

$N(s, a) \leftarrow N(s, a) + 1$

$\text{Returns}(s, a) \leftarrow \text{Returns}(s, a) + G$

$Q(s, a) \leftarrow \text{Returns}(s, a) / N(s, a)$

$\pi(s) \leftarrow \arg \max_a Q(s, a)$  for all  $s$

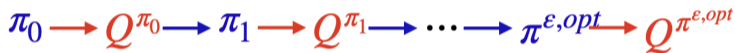
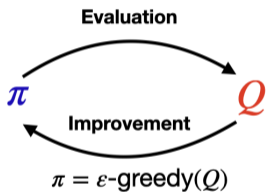
# MC Control with $\epsilon$ -Greedy Exploration

---

- Simplest idea for ensuring continual exploration
- All  $m$  actions are tried with non-zero probabilities
- With probability  $1 - \epsilon$  choose the **greedy** action
- With probability  $\epsilon$  choose an action at **random**

$$\pi(\mathbf{a}|\mathbf{s}) = \begin{cases} \epsilon/m + 1 - \epsilon, & \text{if } \mathbf{a} = \arg \max_{\mathbf{a}'} Q(\mathbf{s}, \mathbf{a}') \\ \epsilon/m, & \text{otherwise} \end{cases}$$

# MC Control with $\epsilon$ -Greedy Exploration (Cont'd)



# Pseudocode

---

- **Initialization:**

$N$  (counter),  $N(\mathbf{s}, \mathbf{a}) \leftarrow \mathbf{0}$  for all  $\mathbf{s} \in \mathcal{S}$ ,  $\mathbf{a} \in \mathcal{A}$

$\text{Returns}(\mathbf{s}, \mathbf{a}) \leftarrow$  empty lists, for all  $\mathbf{s} \in \mathcal{S}$ ,  $\mathbf{a} \in \mathcal{A}$

$\pi \leftarrow$  arbitrary  $\epsilon$ -greedy policy

$Q \leftarrow$  arbitrary

- **Repeat:**

**Generate** an episode using exploring starts and policy  $\pi$

**For each** distinct  $(\mathbf{s}, \mathbf{a})$  appearing in the episode

$G \leftarrow$  return following the first occurrence of  $(\mathbf{s}, \mathbf{a})$

$N(\mathbf{s}, \mathbf{a}) \leftarrow N(\mathbf{s}, \mathbf{a}) + 1$

$\text{Returns}(\mathbf{s}, \mathbf{a}) \leftarrow \text{Returns}(\mathbf{s}, \mathbf{a}) + G$

$Q(\mathbf{s}, \mathbf{a}) \leftarrow \text{Returns}(\mathbf{s}, \mathbf{a}) / N(\mathbf{s}, \mathbf{a})$

**For each** distinct  $\mathbf{s}$ :

$$\pi(\mathbf{a}|\mathbf{s}) \leftarrow \begin{cases} \epsilon/m + 1 - \epsilon, & \text{if } \mathbf{a} = \arg \max Q(\mathbf{s}, \mathbf{a}) \\ \epsilon/m, & \text{otherwise} \end{cases}$$

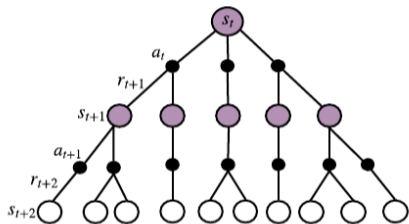
# Summary

---

- Planning v.s. Learning
- Dynamic programming v.s. Monte Carlo Methods
- Policy Iteration v.s. Value Iteration
- Policy Evaluation v.s. Policy Improvement
- MC Policy Evaluation v.s. MC Control
- $\gamma$ -Contraction, Banach Fixed Point Theorem

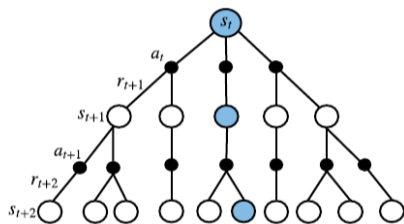
# Summary (Cont'd)

---



$$DP : v(s_t) \leftarrow \mathbb{E}_{\pi} [r_{t+1} + \lambda v(s_{t+1})]$$

Dynamic Programming (DP)

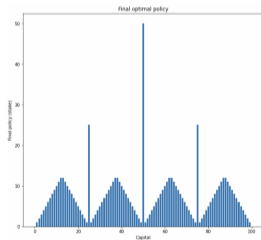
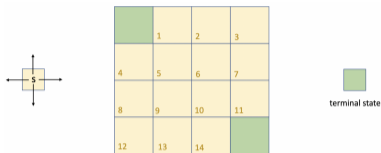


$$MC : v(s_t) \leftarrow v(s_t) + \eta (R_t - v(s_t))$$

Monte Carlo (MC)

# Seminar

- Solution to HW2 (due Wed 12pm)
- Iterative policy evaluation: Gridworld problem
- Value iteration: Gambler's problem



- Monte Carlo prediction & control: Black jack example



# Questions

# Appendix: Proof of Policy Improvement Theorem

---

Consider a sequence of policies:

- $\pi_0$ : a given stationary policy  $\pi$
- $\pi_k$ : a Markov policy that implements  $\pi'$  at the first  $k$  times and follows  $\pi$  afterwards
- $\pi_\infty$ : the greedy policy  $\pi'$

We show in the appendix

- **Step 1:**  $\pi_1$  is no worse than  $\pi_0$ , i.e.,  $Q^\pi(\mathbf{s}, \pi'(\mathbf{s})) \geq V^\pi(\mathbf{s})$
- **Step 2:**  $\pi_{k+1}$  is no worse than  $\pi_k$  for any  $k \geq 1$

This proves the policy improvement theorem

# Appendix: Policy Improvement Theorem, Step 1

---

- $\pi_0$ : a given stationary policy  $\pi$
- $\pi_1$ : a Markov policy that implements  $\pi'$  at the initial time and follows  $\pi$  afterwards
- By definition,

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

- This yields

$$Q^\pi(s, \pi'(s)) = \max_a Q^\pi(s, a) \geq \sum_a \pi(a|s) Q^\pi(s, a) = V^\pi(s)$$

- i.e.,  $\pi_1$  is no worse than  $\pi_0$

## Appendix: Policy Improvement Theorem, Step 2

---

- $\pi_k$ : a Markov policy that implements  $\pi'$  at the first  $k$  times and follows  $\pi$  afterwards
- The difference between two value functions is given by

$$\mathbf{V}^{\pi_{k+1}}(\mathbf{s}) - \mathbf{V}^{\pi_k}(\mathbf{s}) = \gamma^k \mathbb{E}^{\pi'}[\mathbf{Q}^\pi(\mathbf{S}_k, \pi'(\mathbf{S}_k)) | \mathbf{S}_0 = \mathbf{s}] - \gamma^k \mathbb{E}^{\pi'}[\mathbf{V}^\pi(\mathbf{S}_k) | \mathbf{S}_0 = \mathbf{s}]$$

- Results in Step 1 yield  $\mathbf{Q}^\pi(\mathbf{S}_k, \pi'(\mathbf{S}_k)) \geq \mathbf{V}^\pi(\mathbf{S}_k)$ , and hence  $\mathbf{V}^{\pi_{k+1}}(\mathbf{s}) \geq \mathbf{V}^{\pi_k}(\mathbf{s})$
- i.e.,  $\pi_{k+1}$  is no worse than  $\pi_k$

## Appendix: Value Function $\infty$ -Norm

---

- Measure distance between two value functions  $\mathbf{V}_1$  and  $\mathbf{V}_2$  by the  $\infty$ -norm
- i.e., the **largest** difference between state values,

$$\|\mathbf{V}_1 - \mathbf{V}_2\|_{\infty} = \max_{s \in \mathcal{S}} |\mathbf{V}_1(s) - \mathbf{V}_2(s)|$$

- Given a sequence of values  $\{\mathbf{V}_k\}_k$ , convergence requires  $\|\mathbf{V}_k - \mathbf{V}^*\|_{\infty} \rightarrow \mathbf{0}$  for some  $\mathbf{V}^*$  as  $k \rightarrow \infty$

# Appendix: Bellman Expectation Operator

## Definition

Define the Bellman Expectation Operator  $T^\pi$  as a function that maps a given value function  $V$  into another value function  $T^\pi V$  such that

$$T^\pi V(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V(s') \right], \quad \forall s \in \mathcal{S}.$$

- The Bellman equation can be rewritten as  $V^\pi = T^\pi V^\pi$
- This operator is a  $\gamma$ -**contraction**, i.e. it makes value function closer by at least  $\gamma$

$$\begin{aligned} \max_s |T^\pi V_1(s) - T^\pi V_2(s)| &= \gamma \max_s \left| \sum_{a,s'} \pi(a|s) \mathcal{P}_{ss'}^a [V_1(s') - V_2(s')] \right| \\ &\leq \gamma \max_s |V_1(s) - V_2(s)| \max_s \left| \sum_{a,s'} \pi(a|s) \mathcal{P}_{ss'}^a \right| = \gamma \max_s |V_1(s) - V_2(s)| \end{aligned}$$

- Iterative Policy Evaluation:  $V_0 \rightarrow T^\pi V_0 \rightarrow T^\pi T^\pi V_0 \rightarrow \dots$

# Appendix: Banach Fix Point Theorem

---

## Theorem

Suppose  $T$  is a  $\gamma$ -contraction. Then under certain conditions,

- $T$  admits a **unique** fix point  $V^*$ , i.e.  $TV^* = V^*$ ;
- $V^*$  can be found as follows: define a sequence  $\{V_k\}_k$  such that  $V_{k+1} = TV_k$ .  
Then  $V^* = \lim_k V_k$

- Proof can be found [here](#)
- $T^\pi$  is has a unique fix point
- $V^\pi$  is the fix point, according to the Bellman equation
- Iterative policy evaluation converges to  $V^\pi$

# Appendix: Bellman Optimality Operator

## Definition

Define the Bellman Expectation Operator  $T$  as a function that maps a given value function  $V$  into another value function  $TV$  such that

$$TV(s) = \max_{a \in \mathcal{A}} \left[ \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V(s') \right], \quad \forall s \in \mathcal{S}.$$

- The Bellman optimality equation can be rewritten as  $V^{\pi^{\text{opt}}} = TV^{\pi^{\text{opt}}}$
- This operator is a  $\gamma$ -**contraction** as well

$$\begin{aligned} \max_s |TV_1(s) - TV_2(s)| &= \gamma \max_{s,a} \left| \sum_{s'} \mathcal{P}_{ss'}^a [V_1(s') - V_2(s')] \right| \\ &\leq \gamma \max_{s'} |V_1(s') - V_2(s')| \end{aligned}$$



# Appendix: Convergence of Dynamic Programming

---

- $T$  has a unique fix point
- $V^{\pi^{opt}}$  is the fix point, according to the Bellman optimality equation
- According to the Banach fix point theorem, **value iteration** converges to  $V^{\pi^{opt}}$
- **Policy iteration** (that integrates iterative policy evaluation & policy improvement) converges to  $\pi^{opt}$