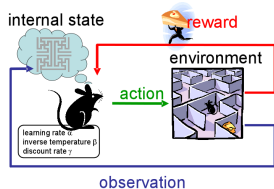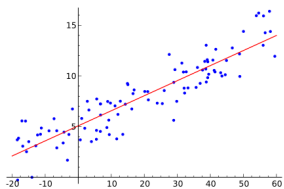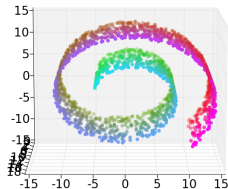# Tutorial on Reinforcement Learning

**Marc Deisenroth**
Department of Computing
Imperial College London

Department of Computer Science
TU Darmstadt

`m.deisenroth@imperial.ac.uk`

# Machine Learning



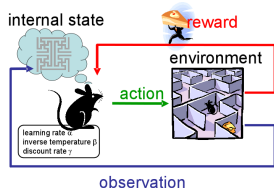- ‣ Unsupervised learning
- ‣ Supervised learning
- ‣ Reinforcement learning

# Machine Learning



- ‣ Unsupervised learning
- ‣ Supervised learning
- ‣ Reinforcement learning

RL makes decisions!

# RL Success Stories



- ‣ Games (e.g., G. Tesauro, D. Silver)
- ‣ Operations research and scheduling (e.g., W. Powell, P. Tadepalli)
- ‣ Recently: robotics (e.g., P. Abbeel, J. Peters, P. Stone, M. Riedmiller)

# Motivation



internal state — reward — environment — action — observation

learning rate $\alpha$
inverse temperature $\beta$
discount rate $\gamma$

- Learning system in an unknown environment
- Knowledge only through interacting with environment
- Explores the environment and receives rewards
- Find strategy/policy, which maximizes overall reward
  ▶▶ Optimal behavior

# Bayesian Decision Theory

‣ Make optimal decisions $a^*$ by maximizing an expected utility

$$a^* \in \arg\max_a \mathbb{E}[r(a)] = \arg\max_a \sum_{j=1}^{m} r(s_j, a)p(s_j)$$

$a$ : decision

$s$ : information about environment/state

# Bayesian Decision Theory

‣ Make optimal decisions $a^*$ by maximizing an expected utility

$$a^* \in \arg\max_a \mathbb{E}[r(a)] = \arg\max_a \sum_{j=1}^{m} r(s_j, a) p(s_j)$$

$a$ : decision

$s$ : information about environment/state

‣ Bayesian sequential decision theory (statistics)

‣ Optimal control theory (engineering)

‣ Reinforcement learning (computer science, psychology)

# Example: Winning the Lottery

| Actions | Outcomes |
|---------|----------|
| $a_1$: play | $s_1$: Win the lottery |
| $a_2$: don't play | $s_2$: Don't win the lottery |

# Example: Winning the Lottery

| Actions | Outcomes |
|---|---|
| $a_1$: play | $s_1$: Win the lottery |
| $a_2$: don't play | $s_2$: Don't win the lottery |

### Optimal action

$$a^* = \arg\max_{a_i} \sum_{j=1}^{2} r_{ij} p(s_j|a_i)$$

# Example: Winning the Lottery

| Actions | Outcomes |
|---------|----------|
| $a_1$: play | $s_1$: Win the lottery |
| $a_2$: don't play | $s_2$: Don't win the lottery |

**Optimal action**

$$a^* = \arg \max_{a_i} \sum_{j=1}^{2} r_{ij} p(s_j | a_i)$$

$$
\begin{array}{ll}
p(s_1|a_1) = 10^{-7} & r_{11} = 500,000 \text{ USD} \\
p(s_2|a_1) = 1 - 10^{-7} & r_{12} = -1 \text{ USD} \\
p(s_1|a_2) = 0 & r_{21} = 0 \text{ USD} \\
p(s_2|a_2) = 1 & r_{22} = 0 \text{ USD}
\end{array}
$$

▶▶ What is the optimal action for this decision problem?

# From Bayesian Decision Theory to RL

‣ So far: single decisions. How do we make a sequence of decisions in order to achieve some long-term rewards?

‣ What about state-dependent actions $a_i|s_j$ ?

# RL Set-up



- Agent interacts with environment to gain knowledge
- Explores and receives rewards
- Actions change the state of the environment
- Choose actions to maximize long-term reward

# RL Set-up



- ‣ Agent interacts with environment to gain knowledge
- ‣ Explores and receives rewards
- ‣ Actions change the state of the environment
- ‣ Choose actions to maximize long-term reward

# RL Set-up



internal state

reward

environment

action

learning rate $\alpha$
inverse temperature $\beta$
discount rate $\gamma$

observation

- ‣ Agent interacts with environment to gain knowledge
- ‣ Explores and receives rewards
- ‣ Actions change the state of the environment
- ‣ Choose actions to maximize long-term reward

▸▸ Markov Decision Process

# Markov Decision Process: Definition

- $\mathcal{S}$: State space (finite)

- $\mathcal{A}$: Action space (finite)
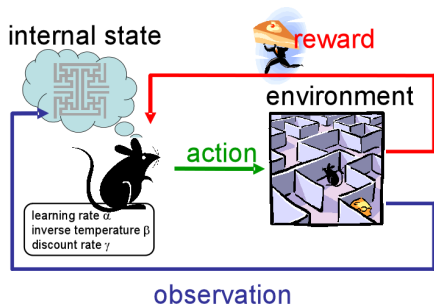
- $\mathcal{P}$: Transition probability $p(\boldsymbol{s}_{k+1}|\boldsymbol{s}_k, \boldsymbol{a}_k)$

- $r$: Reward function

- $\gamma \in [0, 1)$: Discount factor

- $\pi$: Policy
  - Deterministic: $\boldsymbol{a} = \pi(\boldsymbol{s})$
  - Stochastic: $\boldsymbol{a} \sim p_\pi(\boldsymbol{a}|\boldsymbol{s})$     alternative notation: $p_\pi(\boldsymbol{a}|\boldsymbol{s}) = \pi(\boldsymbol{a}|\boldsymbol{s})$

# Markov Decision Process: Definition

- $\mathcal{S}$: State space (finite)

- $\mathcal{A}$: Action space (finite)

- $\mathcal{P}$: Transition probability $p(\boldsymbol{s}_{k+1}|\boldsymbol{s}_k, \boldsymbol{a}_k)$

- $r$: Reward function

- $\gamma \in [0, 1)$: Discount factor

- $\pi$: Policy
  - Deterministic: $\boldsymbol{a} = \pi(\boldsymbol{s})$
  - Stochastic: $\boldsymbol{a} \sim p_\pi(\boldsymbol{a}|\boldsymbol{s})$      alternative notation: $p_\pi(\boldsymbol{a}|\boldsymbol{s}) = \pi(\boldsymbol{a}|\boldsymbol{s})$

## Objective

Find a policy $\pi^*$ that maximizes the expected long-term reward

$$V^\pi(\boldsymbol{s}) = \mathbb{E}\Big[\sum_{k=0}^{\infty} \gamma^k r_{k+1} \big| \boldsymbol{s}_0 = \boldsymbol{s}, \pi\Big], \qquad r_{k+1} = r_{k+1}(\boldsymbol{s}_k, \boldsymbol{a}_k, \boldsymbol{s}_{k+1})$$

# Categorization of RL Algorithms

‣ Value function methods
  ▶▶ Use structure of a value function to discover optimal policies
‣ Value function-free methods (e.g., policy search)
  ▶▶ Search in policy space directly

Motivation

Reinforcement Learning Set-Up

## Value Function Methods

Policy Search

# Value Functions

‣ State Value Function: How good is it to be in a particular state $s$? Well, this depends on the current policy:

$$V^\pi(s) = \mathbb{E}\big[R|s_0 = s\big] = \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|s_0 = s, \pi\big]$$
$$= \mathbb{E}[r_1 + \gamma V^\pi(s_1)|s_0 = s, \pi] \qquad \text{Self-consistency}$$

# Value Functions

‣ State Value Function: How good is it to be in a particular state $s$?
  Well, this depends on the current policy:

$$V^\pi(s) = \mathbb{E}[R|s_0 = s] = \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|s_0 = s, \pi\right]$$
$$= \mathbb{E}[r_1 + \gamma V^\pi(s_1)|s_0 = s, \pi] \qquad \text{Self-consistency}$$

‣ State-Action Value Function: How good is it to be in a particular
  state $s$ and apply action $a$, and afterwards follow policy $\pi$?

$$Q^\pi(s, a) = \mathbb{E}[R|s_0 = s, a_0 = a, \pi]$$
$$= \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|s_0 = s, a_0 = a, \pi\right]$$

# Value Functions

‣ State Value Function: How good is it to be in a particular state $s$? Well, this depends on the current policy:

$$V^\pi(s) = \mathbb{E}\big[R|s_0 = s\big] = \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|s_0 = s, \pi\big]$$
$$= \mathbb{E}[r_1 + \gamma V^\pi(s_1)|s_0 = s, \pi] \qquad \text{Self-consistency}$$

‣ State-Action Value Function: How good is it to be in a particular state $s$ and apply action $a$, and afterwards follow policy $\pi$?

$$Q^\pi(s, a) = \mathbb{E}\big[R|s_0 = s, a_0 = a, \pi\big]$$
$$= \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|s_0 = s, a_0 = a, \pi\big]$$
$$= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2}|s_0 = s, a_0 = a, \pi\big]$$
$$=$$

# Value Functions

‣ State Value Function: How good is it to be in a particular state $s$?
Well, this depends on the current policy:

$$V^\pi(s) = \mathbb{E}\big[R|s_0 = s\big] = \mathbb{E}\big[\sum_{k=0}^\infty \gamma^k r_{k+1}|s_0 = s, \pi\big]$$
$$= \mathbb{E}[r_1 + \gamma V^\pi(s_1)|s_0 = s, \pi] \qquad \text{Self-consistency}$$

‣ State-Action Value Function: How good is it to be in a particular
state $s$ and apply action $a$, and afterwards follow policy $\pi$?

$$Q^\pi(s, a) = \mathbb{E}\big[R|s_0 = s, a_0 = a, \pi\big]$$
$$= \mathbb{E}\big[\sum_{k=0}^\infty \gamma^k r_{k+1}|s_0 = s, a_0 = a, \pi\big]$$
$$= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma \sum_{k=0}^\infty \gamma^k r_{k+2}|s_0 = s, a_0 = a, \pi\big]$$
$$=$$

# Value Functions

‣ State Value Function: How good is it to be in a particular state $s$?
Well, this depends on the current policy:

$$V^\pi(s) = \mathbb{E}\big[R|s_0 = s\big] = \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|s_0 = s, \pi\big]$$
$$= \mathbb{E}[r_1 + \gamma V^\pi(s_1)|s_0 = s, \pi] \qquad \text{Self-consistency}$$

‣ State-Action Value Function: How good is it to be in a particular
state $s$ and apply action $a$, and afterwards follow policy $\pi$?

$$Q^\pi(s, a) = \mathbb{E}\big[R|s_0 = s, a_0 = a, \pi\big]$$
$$= \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|s_0 = s, a_0 = a, \pi\big]$$
$$= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2}|s_0 = s, a_0 = a, \pi\big]$$
$$= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1)|s_0 = s, a_0 = a, \pi\big]$$

# Bellman Operator

$$V^\pi(s) = \mathbb{E}[r_1 + \gamma V^\pi(s_1) | s_0 = s, \pi]$$
$$= T^\pi[V^\pi](s)$$
$$T^\pi[V^\pi](s) := \mathbb{E}[r_1 + \gamma V^\pi(s_1) | s_0 = s, \pi]$$

# Bellman Operator

$$V^\pi(s) = \mathbb{E}[r_1 + \gamma V^\pi(s_1)|s_0 = s, \pi]$$
$$= T^\pi[V^\pi](s)$$
$$T^\pi[V^\pi](s) := \mathbb{E}[r_1 + \gamma V^\pi(s_1)|s_0 = s, \pi]$$

▶▶ $V^\pi = T^\pi[V^\pi]$,     $T^\pi$ : Bellman operator (linear affine)

# Bellman Operator

$$V^\pi(\boldsymbol{s}) = \mathbb{E}[r_1 + \gamma V^\pi(\boldsymbol{s}_1)|\boldsymbol{s}_0 = \boldsymbol{s}, \pi]$$
$$= T^\pi[V^\pi](\boldsymbol{s})$$
$$T^\pi[V^\pi](\boldsymbol{s}) := \mathbb{E}[r_1 + \gamma V^\pi(\boldsymbol{s}_1)|\boldsymbol{s}_0 = \boldsymbol{s}, \pi]$$

▶▶ $V^\pi = T^\pi[V^\pi]$,     $T^\pi$ : Bellman operator (linear affine)

- Fixed point equation with a unique solution for $V^\pi$
  (Banach's FP theorem)

# Bellman Operator

$$V^{\pi}(\boldsymbol{s}) = \mathbb{E}[r_1 + \gamma V^{\pi}(\boldsymbol{s}_1)|\boldsymbol{s}_0 = \boldsymbol{s}, \pi]$$
$$= T^{\pi}[V^{\pi}](\boldsymbol{s})$$
$$T^{\pi}[V^{\pi}](\boldsymbol{s}) := \mathbb{E}[r_1 + \gamma V^{\pi}(\boldsymbol{s}_1)|\boldsymbol{s}_0 = \boldsymbol{s}, \pi]$$

▶ $V^{\pi} = T^{\pi}[V^{\pi}]$,     $T^{\pi}$ : Bellman operator (linear affine)

- Fixed point equation with a unique solution for $V^{\pi}$
  (Banach's FP theorem)
  $$T^{\pi}[V^{\pi}] = r^{\pi} + \gamma \mathcal{P}^{\pi} V^{\pi}$$
  $$\implies V^{\pi} = r^{\pi} + \gamma \mathcal{P}^{\pi} V^{\pi}$$

for suitable representations $r^{\pi}, \mathcal{P}^{\pi}, V^{\pi}$

# Optimal Policies and Value Functions

- Optimal policy $\pi^*$ ensures that $V^{\pi^*}(s) \geq V^{\pi}(s) \quad \forall s \in \mathcal{S}, \pi$
- Existence of $\pi^*$? Uniqueness?

# Optimal Policies and Value Functions

- Optimal policy $\pi^*$ ensures that $V^{\pi^*}(s) \geqslant V^{\pi}(s) \quad \forall s \in \mathcal{S}, \pi$

- Existence of $\pi^*$? Uniqueness?

- Optimal state value function:
  $\forall s \in \mathcal{S} : V^*(s) = V^{\pi^*}(s) = \max_{\pi} V^{\pi}(s)$

# Optimal Policies and Value Functions

- Optimal policy $\pi^*$ ensures that $V^{\pi^*}(s) \geqslant V^{\pi}(s) \quad \forall s \in \mathcal{S}, \pi$

- Existence of $\pi^*$? Uniqueness?

- Optimal state value function:
  $\forall s \in \mathcal{S} : V^*(s) = V^{\pi^*}(s) = \max_{\pi} V^{\pi}(s)$

- Optimal state-action value function:
  $\forall s \in \mathcal{S} : Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a)$

  ▶ Expected return of choosing action $a$ in state $s$ and
    afterwards following the optimal policy $\pi^*$. Note that

  $$Q^*(s, a) = \mathbb{E}[r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a]$$

# Bellman Optimality Equations

$$V^*(s) = \max_{a} Q^*(s, a)$$

# Bellman Optimality Equations

$$V^*(s) = \max_{a} Q^*(s, a)$$

$$= \max_{a} \mathbb{E}\Big[ \sum_{k=0}^{\infty} \gamma^k r_{k+1} | s_0 = s, a_0 = a, \pi^* \Big]$$

# Bellman Optimality Equations

$$V^*(\boldsymbol{s}) = \max_{\boldsymbol{a}} Q^*(\boldsymbol{s}, \boldsymbol{a})$$

$$= \max_{\boldsymbol{a}} \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1} | \boldsymbol{s}_0 = \boldsymbol{s}, \boldsymbol{a}_0 = \boldsymbol{a}, \pi^*\big]$$

$$= \max_{\boldsymbol{a}} \mathbb{E}\big[r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2} | \boldsymbol{s}_0 = \boldsymbol{s}, \boldsymbol{a}_0 = \boldsymbol{a}, \pi^*\big]$$

# Bellman Optimality Equations

$$V^*(s) = \max_a Q^*(s, a)$$

$$= \max_a \mathbb{E}\Big[ \sum_{k=0}^{\infty} \gamma^k r_{k+1} | s_0 = s, a_0 = a, \pi^* \Big]$$

$$= \max_a \mathbb{E}\Big[ r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2} | s_0 = s, a_0 = a, \pi^* \Big]$$

$$= \max_a \mathbb{E}\Big[ r_1 + \gamma V^*(s_1) | s_0 = s, a_0 = a, \pi^* \Big]$$

# Bellman Optimality Equations

$$V^*(\boldsymbol{s}) = \max_{\boldsymbol{a}} Q^*(\boldsymbol{s}, \boldsymbol{a})$$

$$= \max_{\boldsymbol{a}} \mathbb{E}\Big[ \sum_{k=0}^{\infty} \gamma^k r_{k+1} | \boldsymbol{s}_0 = \boldsymbol{s}, \boldsymbol{a}_0 = \boldsymbol{a}, \pi^* \Big]$$

$$= \max_{\boldsymbol{a}} \mathbb{E}\Big[ r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2} | \boldsymbol{s}_0 = \boldsymbol{s}, \boldsymbol{a}_0 = \boldsymbol{a}, \pi^* \Big]$$

$$= \max_{\boldsymbol{a}} \mathbb{E}\Big[ r_1 + \gamma V^*(\boldsymbol{s}_1) | \boldsymbol{s}_0 = \boldsymbol{s}, \boldsymbol{a}_0 = \boldsymbol{a}, \pi^* \Big]$$

▶▶ $V^* = T^*[V^*]$ $\qquad T^*$ : Bellman optimality operator (nonlinear)

# Bellman Optimality Equations

$$V^*(s) = \max_a Q^*(s, a)$$

$$= \max_a \mathbb{E}\Big[\sum_{k=0}^{\infty} \gamma^k r_{k+1} | s_0 = s, a_0 = a, \pi^*\Big]$$

$$= \max_a \mathbb{E}\Big[r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2} | s_0 = s, a_0 = a, \pi^*\Big]$$

$$= \max_a \mathbb{E}\big[r_1 + \gamma V^*(s_1) | s_0 = s, a_0 = a, \pi^*\big]$$

▸ $V^* = T^*[V^*]$     $T^*$ : Bellman optimality operator (nonlinear)

$$Q^*(s, a) = \mathbb{E}[r_1 + \gamma \max_{a'} Q^*(s_1, a') | s_0 = s, a_0 = a]$$

$$=$$

# Bellman Optimality Equations

$$V^*(s) = \max_a Q^*(s, a)$$

$$= \max_a \mathbb{E}\Big[ \sum_{k=0}^{\infty} \gamma^k r_{k+1} | s_0 = s, a_0 = a, \pi^* \Big]$$

$$= \max_a \mathbb{E}\Big[ r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2} | s_0 = s, a_0 = a, \pi^* \Big]$$

$$= \max_a \mathbb{E}\big[ r_1 + \gamma V^*(s_1) | s_0 = s, a_0 = a, \pi^* \big]$$

▶▶ $V^* = T^*[V^*]$     $T^*$ : Bellman optimality operator (nonlinear)

$$Q^*(s, a) = \mathbb{E}\big[ r_1 + \gamma \max_{a'} Q^*(s_1, a') | s_0 = s, a_0 = a \big]$$

$$= \mathbb{E}\big[ r_1 + \gamma V^*(s_1) | s_0 = s, a_0 = a \big]$$

# Solving MDPs

- Assume you know $V^*$
- **One-step search**: be "greedy" with respect to $V^*$

$$\pi^*(s) = \arg\max_a \mathbb{E}[r(s, a, s') + \gamma V^*(s')]$$

# Solving MDPs

‣ Assume you know $V^*$

‣ **One-step search**: be "greedy" with respect to $V^*$

$$\pi^*(s) = \arg\max_a \mathbb{E}[r(s, a, s') + \gamma V^*(s')]$$

‣ Assume you know $Q^*$

‣ **Zero-step search**:

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

# Solving MDPs

- Assume you know $V^*$
- **One-step search**: be "greedy" with respect to $V^*$

$$\pi^*(s) = \arg\max_a \mathbb{E}[r(s, a, s') + \gamma V^*(s')]$$

- Assume you know $Q^*$
- **Zero-step search**:

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

Assumptions for solving the Bellman equations exactly:

- Know the transition probabilities $p(s'|s, a)$
- Sufficient computational resources available
- Markov property holds

# Solving MDPs

‣ Assume you know $V^*$

‣ **One-step search**: be "greedy" with respect to $V^*$

$$\pi^*(\boldsymbol{s}) = \arg\max_{\boldsymbol{a}} \mathbb{E}[r(\boldsymbol{s}, \boldsymbol{a}, \boldsymbol{s}') + \gamma V^*(\boldsymbol{s}')]$$

‣ Assume you know $Q^*$

‣ **Zero-step search**:

$$\pi^*(\boldsymbol{s}) = \max_{\boldsymbol{a}} Q^*(\boldsymbol{s}, \boldsymbol{a})$$

Assumptions for solving the Bellman equations exactly:

‣ Know the transition probabilities $p(\boldsymbol{s}'|\boldsymbol{s}, \boldsymbol{a})$

‣ Sufficient computational resources available

‣ Markov property holds

▶▶ Approximate solutions in practice

# Solving MDPs (2)

- Exact: Dynamic programming
- Approximate: Monte Carlo, Temporal Difference Learning

# Dynamic Programming

Assumptions:

- Perfect model $p(s'|s, a)$ is known
- Typically finite state spaces $\mathcal{S}$ and action spaces $\mathcal{A}$
- Expected immediate rewards $\mathbb{E}[r(s, a, s')]$ are known

▶▶ Use value functions to structure the search for good policies

# Policy Evaluation

## Objective

For a given policy $\pi$, find the corresponding value function $V^\pi$

- Exploit the fixed-point property of the value function
  $V^\pi = T^\pi[V^\pi]$ :
    - Initialize $V_0^\pi$ arbitrarily
    - Find $V^\pi$ as the limit of the sequence $V_0^\pi, V_1^\pi, \ldots$

# Policy Evaluation

## Objective

For a given policy $\pi$, find the corresponding value function $V^\pi$

- Exploit the fixed-point property of the value function
  $V^\pi = T^\pi[V^\pi]$ :
    - Initialize $V_0^\pi$ arbitrarily
    - Find $V^\pi$ as the limit of the sequence $V_0^\pi, V_1^\pi, \ldots$

## Update Rule

$$\forall s \in \mathcal{S} : V_{k+1}^\pi(s) \leftarrow \mathbb{E}[r(s, a, s') + \gamma V_k^\pi(s')|s, \pi]$$

$$a \sim p_\pi(a|s), \quad s' \sim p(s'|s, a)$$

▶▶ Bootstrapping

# Policy Improvement

- So far: We know $V^\pi$, but we want $V^*$
- Find a better policy $\pi'$

### Objective

Find a policy $\pi' \geqslant \pi$, i.e., $V^{\pi'} \geqslant V^\pi$

# Policy Improvement Theorem

## Policy Improvement Theorem

If $\pi, \pi'$ are two (deterministic) policies with

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geqslant Q^\pi(s, \pi(s)) = V^\pi(s)$$

then $\pi' \geqslant \pi$ and $V^{\pi'} \geqslant V^\pi$, i.e., $\pi'$ improves $\pi$.

# Policy Improvement Theorem

### Policy Improvement Theorem

If $\pi, \pi'$ are two (deterministic) policies with

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geqslant Q^\pi(s, \pi(s)) = V^\pi(s)$$

then $\pi' \geqslant \pi$ and $V^{\pi'} \geqslant V^\pi$, i.e., $\pi'$ improves $\pi$.

- For stochastic policies:

$$
\begin{aligned}
Q^\pi(s, \pi'(s)) &= \mathbb{E}_a[Q^\pi(s, a)] & a \sim p_{\pi'}(a|s) \\
&= \sum_a p_{\pi'}(a|s) Q^\pi(s, a)
\end{aligned}
$$

# Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geqslant V^\pi(s) = Q^\pi(s, \pi(s)) \qquad (*)$$

$$V^\pi(s) \overset{(*)}{\leqslant} Q^\pi(s, \pi'(s))$$

# Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geq V^\pi(s) = Q^\pi(s, \pi(s)) \tag{$*$}$$

$$V^\pi(s) \overset{(*)}{\leq} Q^\pi(s, \pi'(s))$$
$$= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1) | s_0 = s, a_0 = \pi'(s_0)\big]$$

# Proof

$$\forall \boldsymbol{s} \in \mathcal{S}: \quad Q^\pi(\boldsymbol{s}, \pi'(\boldsymbol{s})) \geqslant V^\pi(\boldsymbol{s}) = Q^\pi(\boldsymbol{s}, \pi(\boldsymbol{s})) \tag{$*$}$$

$$V^\pi(\boldsymbol{s}) \overset{(*)}{\leqslant} Q^\pi(\boldsymbol{s}, \pi'(\boldsymbol{s}))$$
$$= \mathbb{E}\big[r_1(\boldsymbol{s}_0, \boldsymbol{a}_0, \boldsymbol{s}_1) + \gamma V^\pi(\boldsymbol{s}_1) | \boldsymbol{s}_0 = \boldsymbol{s}, \boldsymbol{a}_0 = \pi'(\boldsymbol{s}_0)\big]$$

# Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geq V^\pi(s) = Q^\pi(s, \pi(s)) \tag{$*$}$$

$$
\begin{aligned}
V^\pi(s) &\overset{(*)}{\leq} Q^\pi(s, \pi'(s)) \\
&= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1) | s_0 = s, a_0 = \pi'(s_0)\big] \\
&\overset{(*)}{\leq} \mathbb{E}\big[r_1 + \gamma Q^\pi(s_1, \pi'(s_1)) | s_0 = s, a_0 = \pi'(s_0), a_1 = \pi'(s_1)\big]
\end{aligned}
$$

# Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geqslant V^\pi(s) = Q^\pi(s, \pi(s)) \qquad (*)$$

$$
\begin{aligned}
V^\pi(s) &\overset{(*)}{\leqslant} Q^\pi(s, \pi'(s)) \\
&= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1) | s_0 = s, a_0 = \pi'(s_0)\big] \\
&\overset{(*)}{\leqslant} \mathbb{E}\big[r_1 + \gamma Q^\pi(s_1, \pi'(s_1)) | s_0 = s, a_0 = \pi'(s_0), a_1 = \pi'(s_1)\big]
\end{aligned}
$$

# Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geqslant V^\pi(s) = Q^\pi(s, \pi(s)) \tag{*}$$

$$
\begin{aligned}
V^\pi(s) &\overset{(*)}{\leqslant} Q^\pi(s, \pi'(s)) \\
&= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1) | s_0 = s, a_0 = \pi'(s_0)\big] \\
&\overset{(*)}{\leqslant} \mathbb{E}\big[r_1 + \gamma Q^\pi(s_1, \pi'(s_1)) | s_0 = s, a_0 = \pi'(s_0), a_1 = \pi'(s_1)\big] \\
&= \mathbb{E}\big[r_1 + \gamma \mathbb{E}[r_2 + \gamma V^\pi(s_2)] | s_0 = s\big]
\end{aligned}
$$

# Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geqslant V^\pi(s) = Q^\pi(s, \pi(s)) \tag{$*$}$$

$$
\begin{aligned}
V^\pi(s) &\overset{(*)}{\leqslant} Q^\pi(s, \pi'(s)) \\
&= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1) | s_0 = s, a_0 = \pi'(s_0)\big] \\
&\overset{(*)}{\leqslant} \mathbb{E}\big[r_1 + \gamma Q^\pi(s_1, \pi'(s_1)) | s_0 = s, a_0 = \pi'(s_0), a_1 = \pi'(s_1)\big] \\
&= \mathbb{E}\big[r_1 + \gamma \mathbb{E}[r_2 + \gamma V^\pi(s_2)] | s_0 = s\big]
\end{aligned}
$$

# Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geqslant V^\pi(s) = Q^\pi(s, \pi(s)) \tag{$*$}$$

$$
\begin{aligned}
V^\pi(s) &\stackrel{(*)}{\leqslant} Q^\pi(s, \pi'(s)) \\
&= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1) | s_0 = s, a_0 = \pi'(s_0)\big] \\
&\stackrel{(*)}{\leqslant} \mathbb{E}\big[r_1 + \gamma Q^\pi(s_1, \pi'(s_1)) | s_0 = s, a_0 = \pi'(s_0), a_1 = \pi'(s_1)\big] \\
&= \mathbb{E}\big[r_1 + \gamma \mathbb{E}[r_2 + \gamma V^\pi(s_2)] | s_0 = s\big] \\
&= \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 V^\pi(s_2) | s_0 = s\big]
\end{aligned}
$$

# Proof

$$\forall \boldsymbol{s} \in \mathcal{S}: \quad Q^\pi(\boldsymbol{s}, \pi'(\boldsymbol{s})) \geqslant V^\pi(\boldsymbol{s}) = Q^\pi(\boldsymbol{s}, \pi(\boldsymbol{s})) \tag{$*$}$$

$$
\begin{aligned}
V^\pi(\boldsymbol{s}) &\overset{(*)}{\leqslant} Q^\pi(\boldsymbol{s}, \pi'(\boldsymbol{s})) \\
&= \mathbb{E}\big[r_1(\boldsymbol{s}_0, \boldsymbol{a}_0, \boldsymbol{s}_1) + \gamma V^\pi(\boldsymbol{s}_1)|\boldsymbol{s}_0 = \boldsymbol{s}, \boldsymbol{a}_0 = \pi'(\boldsymbol{s}_0)\big] \\
&\overset{(*)}{\leqslant} \mathbb{E}\big[r_1 + \gamma Q^\pi(\boldsymbol{s}_1, \pi'(\boldsymbol{s}_1))|\boldsymbol{s}_0 = \boldsymbol{s}, \boldsymbol{a}_0 = \pi'(\boldsymbol{s}_0), \boldsymbol{a}_1 = \pi'(\boldsymbol{s}_1)\big] \\
&= \mathbb{E}\big[r_1 + \gamma \mathbb{E}[r_2 + \gamma V^\pi(\boldsymbol{s}_2)]|\boldsymbol{s}_0 = \boldsymbol{s}\big] \\
&= \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 V^\pi(\boldsymbol{s}_2)|\boldsymbol{s}_0 = \boldsymbol{s}\big]
\end{aligned}
$$

# Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geqslant V^\pi(s) = Q^\pi(s, \pi(s)) \tag{$*$}$$

$$
\begin{aligned}
V^\pi(s) &\overset{(*)}{\leqslant} Q^\pi(s, \pi'(s)) \\
&= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1)|s_0 = s, a_0 = \pi'(s_0)\big] \\
&\overset{(*)}{\leqslant} \mathbb{E}\big[r_1 + \gamma Q^\pi(s_1, \pi'(s_1))|s_0 = s, a_0 = \pi'(s_0), a_1 = \pi'(s_1)\big] \\
&= \mathbb{E}\big[r_1 + \gamma \mathbb{E}[r_2 + \gamma V^\pi(s_2)]|s_0 = s\big] \\
&= \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 V^\pi(s_2)|s_0 = s\big] \\
&\overset{(*)}{\leqslant} \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 Q^\pi(s_2, \pi'(s_2))|s_0 = s\big]
\end{aligned}
$$

# Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geq V^\pi(s) = Q^\pi(s, \pi(s)) \tag{$*$}$$

$$
\begin{aligned}
V^\pi(s) &\overset{(*)}{\leq} Q^\pi(s, \pi'(s)) \\
&= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1) | s_0 = s, a_0 = \pi'(s_0)\big] \\
&\overset{(*)}{\leq} \mathbb{E}\big[r_1 + \gamma Q^\pi(s_1, \pi'(s_1)) | s_0 = s, a_0 = \pi'(s_0), a_1 = \pi'(s_1)\big] \\
&= \mathbb{E}\big[r_1 + \gamma \mathbb{E}[r_2 + \gamma V^\pi(s_2)] | s_0 = s\big] \\
&= \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 V^\pi(s_2) | s_0 = s\big] \\
&\overset{(*)}{\leq} \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 Q^\pi(s_2, \pi'(s_2)) | s_0 = s\big]
\end{aligned}
$$

# Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geqslant V^\pi(s) = Q^\pi(s, \pi(s)) \tag{$*$}$$

$$
\begin{aligned}
V^\pi(s) &\overset{(*)}{\leqslant} Q^\pi(s, \pi'(s)) \\
&= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1)|s_0 = s, a_0 = \pi'(s_0)\big] \\
&\overset{(*)}{\leqslant} \mathbb{E}\big[r_1 + \gamma Q^\pi(s_1, \pi'(s_1))|s_0 = s, a_0 = \pi'(s_0), a_1 = \pi'(s_1)\big] \\
&= \mathbb{E}\big[r_1 + \gamma \mathbb{E}[r_2 + \gamma V^\pi(s_2)]|s_0 = s\big] \\
&= \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 V^\pi(s_2)|s_0 = s\big] \\
&\overset{(*)}{\leqslant} \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 Q^\pi(s_2, \pi'(s_2))|s_0 = s\big] \\
&= \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 V^\pi(s_3)|s_0 = s\big]
\end{aligned}
$$

## Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geq V^\pi(s) = Q^\pi(s, \pi(s)) \tag{$*$}$$

$$
\begin{aligned}
V^\pi(s) &\overset{(*)}{\leq} Q^\pi(s, \pi'(s)) \\
&= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1)|s_0 = s, a_0 = \pi'(s_0)\big] \\
&\overset{(*)}{\leq} \mathbb{E}\big[r_1 + \gamma Q^\pi(s_1, \pi'(s_1))|s_0 = s, a_0 = \pi'(s_0), a_1 = \pi'(s_1)\big] \\
&= \mathbb{E}\big[r_1 + \gamma \mathbb{E}[r_2 + \gamma V^\pi(s_2)]|s_0 = s\big] \\
&= \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 V^\pi(s_2)|s_0 = s\big] \\
&\overset{(*)}{\leq} \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 Q^\pi(s_2, \pi'(s_2))|s_0 = s\big] \\
&= \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 V^\pi(s_3)|s_0 = s\big]
\end{aligned}
$$

## Proof

$$\forall s \in \mathcal{S}: \quad Q^\pi(s, \pi'(s)) \geqslant V^\pi(s) = Q^\pi(s, \pi(s)) \tag{$*$}$$

$$
\begin{aligned}
V^\pi(s) &\overset{(*)}{\leqslant} Q^\pi(s, \pi'(s)) \\
&= \mathbb{E}\big[r_1(s_0, a_0, s_1) + \gamma V^\pi(s_1) | s_0 = s, a_0 = \pi'(s_0)\big] \\
&\overset{(*)}{\leqslant} \mathbb{E}\big[r_1 + \gamma Q^\pi(s_1, \pi'(s_1)) | s_0 = s, a_0 = \pi'(s_0), a_1 = \pi'(s_1)\big] \\
&= \mathbb{E}\big[r_1 + \gamma \mathbb{E}[r_2 + \gamma V^\pi(s_2)] | s_0 = s\big] \\
&= \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 V^\pi(s_2) | s_0 = s\big] \\
&\overset{(*)}{\leqslant} \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 Q^\pi(s_2, \pi'(s_2)) | s_0 = s\big] \\
&= \mathbb{E}\big[r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 V^\pi(s_3) | s_0 = s\big] \\
&\cdots \\
&\leqslant \mathbb{E}\big[\textstyle\sum_{k=0}^\infty \gamma^k r_{k+1} | s_0 = s\big] = V^{\pi'}(s)
\end{aligned}
$$

# Policy Improvement

‣ Easy to evaluate a change in the policy at a single state $s$

# Policy Improvement

‣ Easy to evaluate a change in the policy at a single state $s$

‣ Extend this idea to all states:

$$\forall s \in \mathcal{S}: \quad \pi'(s) = \arg\max_a \mathbb{E}\big[r(s, a, s') + \gamma V^\pi(s')\big]$$
$$= \arg\max_a Q^\pi(s, a)$$

# Policy Improvement

‣ Easy to evaluate a change in the policy at a single state $s$

‣ Extend this idea to all states:

$$\forall s \in \mathcal{S} : \quad \pi'(s) = \arg \max_a \mathbb{E}\big[r(s, a, s') + \gamma V^\pi(s')\big]$$
$$= \arg \max_a Q^\pi(s, a)$$

‣ When $\pi' = \pi$ then $V^{\pi'} = V^\pi$ ▶▶ Convergence (Why?)

# Policy Improvement

‣ Easy to evaluate a change in the policy at a single state $s$

‣ Extend this idea to all states:

$$\forall s \in \mathcal{S}: \quad \pi'(s) = \arg\max_{a} \mathbb{E}\big[r(s, a, s') + \gamma V^{\pi}(s')\big]$$
$$= \arg\max_{a} Q^{\pi}(s, a)$$

‣ When $\pi' = \pi$ then $V^{\pi'} = V^{\pi}$ ▶▶ Convergence (Why?)

‣ If $V^{\pi}$ is known, we need a model $p(s'|s, a)$ for the one-step look-ahead

# Policy Improvement

‣ Easy to evaluate a change in the policy at a single state $s$

‣ Extend this idea to all states:

$$\forall s \in \mathcal{S} : \quad \pi'(s) = \arg \max_a \mathbb{E}\big[r(s, a, s') + \gamma V^\pi(s')\big]$$
$$= \arg \max_a Q^\pi(s, a)$$

‣ When $\pi' = \pi$ then $V^{\pi'} = V^\pi$ ▶▶ Convergence (Why?)

‣ If $V^\pi$ is known, we need a model $p(s'|s, a)$ for the one-step look-ahead

‣ If $Q^\pi$ is known, we don't need a model (no prediction required)

# Policy Improvement

‣ Easy to evaluate a change in the policy at a single state $s$

‣ Extend this idea to all states:

$$\forall s \in \mathcal{S} : \quad \pi'(s) = \arg \max_a \mathbb{E}\big[r(s, a, s') + \gamma V^\pi(s')\big]$$
$$= \arg \max_a Q^\pi(s, a)$$

‣ When $\pi' = \pi$ then $V^{\pi'} = V^\pi$ ▶▶ Convergence (Why?)

‣ If $V^\pi$ is known, we need a model $p(s'|s, a)$ for the one-step look-ahead

‣ If $Q^\pi$ is known, we don't need a model (no prediction required)

‣ Greedy policy update with respect to the value function (but look implicitly at long-term rewards)

# Policy Iteration

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \cdots \xrightarrow{E} V^* \xrightarrow{I} \pi^*$$

‣ E: policy evaluation

‣ I: policy improvement

# Policy Iteration

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \cdots \xrightarrow{E} V^* \xrightarrow{I} \pi^*$$

- E: policy evaluation
- I: policy improvement
- Strict policy improvement at each step
  (unless policy is already optimal)
- Converges often after a few iterations

# Policy Iteration

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \cdots \xrightarrow{E} V^* \xrightarrow{I} \pi^*$$

- E: policy evaluation
- I: policy improvement
- Strict policy improvement at each step
  (unless policy is already optimal)
- Converges often after a few iterations
- Each policy evaluation is itself an iterative process
  ▶▶ Can be really slow!

# Value Iteration

$$\pi_0 \xrightarrow{\;E\;} V^{\pi_0} \xrightarrow{\;I\;} \pi_1 \xrightarrow{\;E\;} V^{\pi_1} \xrightarrow{\;I\;} \cdots \xrightarrow{\;E\;} V^* \xrightarrow{\;I\;} \pi^*$$

‣ Stop policy evaluation after a single update
  ▶▶ No longer an iterative process

# Value Iteration

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \cdots \xrightarrow{E} V^* \xrightarrow{I} \pi^*$$

‣ Stop policy evaluation after a single update
  ▶ No longer an iterative process

## Update Rule

$$V_{k+1}(\boldsymbol{s}) = \max_{\boldsymbol{a}} \mathbb{E}\big[r(\boldsymbol{s}, \boldsymbol{a}, \boldsymbol{s}') + \gamma V_k(\boldsymbol{s}')|\boldsymbol{s}, \boldsymbol{a}\big]$$

$$= \max_{\boldsymbol{a}} \sum_{\boldsymbol{s}'} p(\boldsymbol{s}'|\boldsymbol{s}, \boldsymbol{a})\big(\mathbb{E}[r(\boldsymbol{s}, \boldsymbol{a}, \boldsymbol{s}')] + \gamma V_k(\boldsymbol{s}')\big)$$

# Value Iteration

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \cdots \xrightarrow{E} V^* \xrightarrow{I} \pi^*$$
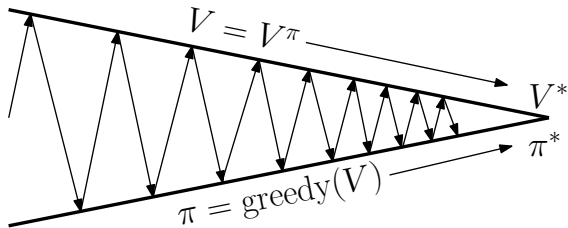
‣ Stop policy evaluation after a single update
  ▶ No longer an iterative process

### Update Rule

$$V_{k+1}(\boldsymbol{s}) = \max_{\boldsymbol{a}} \mathbb{E}\big[r(\boldsymbol{s}, \boldsymbol{a}, \boldsymbol{s}') + \gamma V_k(\boldsymbol{s}')|\boldsymbol{s}, \boldsymbol{a}\big]$$
$$= \max_{\boldsymbol{a}} \sum_{\boldsymbol{s}'} p(\boldsymbol{s}'|\boldsymbol{s}, \boldsymbol{a})\big(\mathbb{E}[r(\boldsymbol{s}, \boldsymbol{a}, \boldsymbol{s}')] + \gamma V_k(\boldsymbol{s}')\big)$$
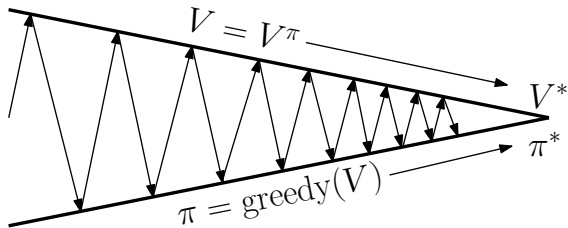
‣ Bootstrapping
‣ Bellman optimality equation as an update rule: $V_{k+1} \leftarrow T^*[V_k]$
‣ Identical to policy evaluation backup if you add the max operator
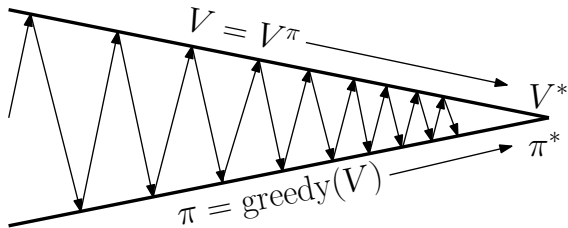
# Generalized Policy Iteration (GPI)



- Abstraction/generalization of policy iteration
- Two interacting processes: policy evaluation/improvement

# Generalized Policy Iteration (GPI)



- Abstraction/generalization of policy iteration
- Two interacting processes: policy evaluation/improvement
- Update details abstracted away (policy needs to be greedy)
- Don't need to go full way to $\pi$ or $V^\pi$, just "in the direction"

# Generalized Policy Iteration (GPI)



- Abstraction/generalization of policy iteration
- Two interacting processes: policy evaluation/improvement
- Update details abstracted away (policy needs to be greedy)
- Don't need to go full way to $\pi$ or $V^\pi$, just "in the direction"
- Both processes converge to a single joint solution $(V^*, \pi^*)$
- Value iteration (incomplete value function updates) is one example of GPI

# Summary: Dynamic Programming

- Find optimal policies via value functions and bootstrapping
- Exact method (standard method in optimal control)

# Summary: Dynamic Programming

- Find optimal policies via value functions and bootstrapping
- Exact method (standard method in optimal control)
- Computationally expensive (sweeps through state-action spaces)
  ▶▶ Curse of dimensionality

# Summary: Dynamic Programming

‣ Find optimal policies via value functions and bootstrapping

‣ Exact method (standard method in optimal control)

‣ Computationally expensive (sweeps through state-action spaces)
   ▶▶ Curse of dimensionality

‣ Exponentially faster than any direct policy search method (if the policy space is not restricted)

‣ Requires a model $p(s'|s, a)$ and the knowledge of $\mathbb{E}[r(s, a, s')]$ before DP can be applied.

‣ 2 algorithms: Policy iteration, value iteration

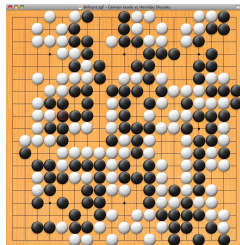# Approximate Value Function Methods

Look trajectory samples

- Monte Carlo methods
- Temporal difference learning

# Monte Carlo Methods

## Key Idea

$$V^\pi(\boldsymbol{s}) = \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1} | \pi, \boldsymbol{s}_0 = \boldsymbol{s}\big]$$

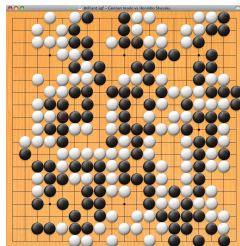Estimate value function by averaging returns of sampled trajectories

# Monte Carlo Methods

## Key Idea

$$V^\pi(s) = \mathbb{E}\big[ \sum_{k=0}^{\infty} \gamma^k r_{k+1} | \pi, s_0 = s \big]$$

Estimate value function by averaging returns of sampled trajectories
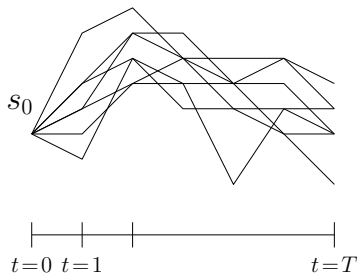


Properties:

- Model free (no knowledge of $p(s'|s, a)$ required) ▸▸ very general
- Learn from online experience (sampled trajectories of states, actions, rewards)
- Compute the same value function as DP (in the limit)

# Episodic Set-up

- Consider a finite time horizon of length $T$
- Usually a fixed set of initial states $s_0$
- Observe rewards $r_1, r_2, \ldots, r_T | s_0, \pi$
- Value estimates $V^\pi(s_0)$ updated at the end of an episode (not after each time step)

# First-Visit Monte Carlo Policy Evaluation

- Generate trajectories with policy $\pi$
- Record reward after visiting state $s$, average at the end

1: **for** $i = 1$ to $\infty$ **do**
2:      Generate trajectory $\tau_i^\pi$ with policy $\pi$
3:      **for** all states $s \in \tau_i$ **do**
4:          $r \leftarrow$ sum of rewards following the <u>first</u> occurrence of $s$
5:          $R(s) \leftarrow [R(s), r]$             $\rhd$ Append $r$ to array
6:      **end for**
7:      $V^\pi(s) \approx \mathbb{E}[R(s)|\pi] \approx \frac{1}{|R(s)|} \sum_{j=1}^{|R(s)|} R(s)[j]$
8: **end for**

- Convergence to the correct value $V^\pi$ in the limit

# Properties

‣ Computational complexity independent of the size of the state space

‣ Very valuable if we are only interested in values starting from a small set of states $s_0$ ▶ episodic set-up

‣ No bootstrapping (unlike DP)

‣ Learn from actual experience (DP only uses the model)

# Properties

- Computational complexity independent of the size of the state space
- Very valuable if we are only interested in values starting from a small set of states $s_0$ ▶ episodic set-up
- No bootstrapping (unlike DP)
- Learn from actual experience (DP only uses the model)

Next steps?

- Low-cost (computation) solution to computer $V^\pi$
- Is this useful for policy improvement? Why? (not?)

# Properties

- Computational complexity independent of the size of the state space
- Very valuable if we are only interested in values starting from a small set of states $s_0$ ▶ episodic set-up
- No bootstrapping (unlike DP)
- Learn from actual experience (DP only uses the model)

Next steps?

- Low-cost (computation) solution to computer $V^\pi$
- Is this useful for policy improvement? Why? (not?)
- If we only know $V^\pi$, we need to perform a one-step search for policy improvement...
  ▶ Need an estimate of $Q^\pi$ if we don't have a model $p(s'|s, a)$

## Monte Carlo for Q Function

1: **for** $i = 1$ to $\infty$ **do**
2:     Generate trajectory $\tau_i^\pi$ with policy $\pi$
3:     **for** all states $s \in \tau_i$ **do**
4:         $r \leftarrow$ sum of rewards following the <u>first</u> occurrence of $s$
5:         $R(s) \leftarrow [R(s), r]$                        $\rhd$ Append $r$ to array
6:     **end for**
7:     $V^\pi(s) \approx \mathbb{E}[R(s)|\pi] \approx \frac{1}{|R(s)|} \sum\limits_{j=1}^{|R(s)|} R(s)[j]$
8: **end for**

# Monte Carlo for Q Function

1: **for** $i = 1$ to $\infty$ **do**
2:      Generate trajectory $\tau_i^\pi$ with policy $\pi$
3:      **for** all state-action pairs $(s, a) \in \tau_i$ **do**
4:          $r \leftarrow$ sum of rewards following the <u>first</u> occurrence of $(s, a)$
5:          $R(s, a) \leftarrow [R(s, a), r]$                 $\triangleright$ Append $r$ to array
6:      **end for**
7:      $Q^\pi(s, a) \approx \mathbb{E}[R(s, a)|\pi] \approx \frac{1}{|R(s,a)|} \sum_{j=1}^{|R(s,a)|} R(s, a)[j]$
8: **end for**

# Monte Carlo for Q Function

1: **for** $i = 1$ to $\infty$ **do**
2:      Generate trajectory $\tau_i^\pi$ with policy $\pi$
3:      **for** all state-action pairs $(s, a) \in \tau_i$ **do**
4:          $r \leftarrow$ sum of rewards following the <u>first</u> occurrence of $(s, a)$
5:          $R(s, a) \leftarrow [R(s, a), \, r]$          $\triangleright$ Append $r$ to array
6:      **end for**
7:      $Q^\pi(s, a) \approx \mathbb{E}[R(s, a)|\pi] \approx \frac{1}{|R(s,a)|} \sum_{j=1}^{|R(s,a)|} R(s, a)[j]$
8: **end for**

‣ Any potential problems?

# Monte Carlo for Q Function

1: **for** $i = 1$ to $\infty$ **do**
2:      Generate trajectory $\tau_i^\pi$ with policy $\pi$
3:      **for** all state-action pairs $(s, a) \in \tau_i$ **do**
4:          $r \leftarrow$ sum of rewards following the <u>first</u> occurrence of $(s, a)$
5:          $R(s, a) \leftarrow [R(s, a), \, r]$               $\triangleright$ Append $r$ to array
6:      **end for**
7:      $Q^\pi(s, a) \approx \mathbb{E}[R(s, a)|\pi] \approx \frac{1}{|R(s,a)|} \sum_{j=1}^{|R(s,a)|} R(s, a)[j]$
8: **end for**

- Any potential problems?

▶▶ $\infty$ many trajectories (use GPI idea)

# Monte Carlo for Q Function

1: **for** $i = 1$ to $\infty$ **do**
2:      Generate trajectory $\tau_i^\pi$ with policy $\pi$
3:      **for** all state-action pairs $(s, a) \in \tau_i$ **do**
4:          $r \leftarrow$ sum of rewards following the <u>first</u> occurrence of $(s, a)$
5:          $R(s, a) \leftarrow [R(s, a), \, r]$          $\triangleright$ Append $r$ to array
6:      **end for**
7:      $Q^\pi(s, a) \approx \mathbb{E}[R(s, a)|\pi] \approx \frac{1}{|R(s,a)|} \sum_{j=1}^{|R(s,a)|} R(s, a)[j]$
8: **end for**

- Any potential problems?
- Deterministic policy?

▶▶ $\infty$ many trajectories (use GPI idea)

# Monte Carlo for Q Function

1: **for** $i = 1$ to $\infty$ **do**
2:      Generate trajectory $\tau_i^\pi$ with policy $\pi$
3:      **for** all state-action pairs $(\boldsymbol{s}, \boldsymbol{a}) \in \tau_i$ **do**
4:          $r \leftarrow$ sum of rewards following the <u>first</u> occurrence of $(\boldsymbol{s}, \boldsymbol{a})$
5:          $R(\boldsymbol{s}, \boldsymbol{a}) \leftarrow [R(\boldsymbol{s}, \boldsymbol{a}), \, r]$          $\triangleright$ Append $r$ to array
6:      **end for**
7:      $Q^\pi(\boldsymbol{s}, \boldsymbol{a}) \approx \mathbb{E}[R(\boldsymbol{s}, \boldsymbol{a})|\pi] \approx \frac{1}{|R(\boldsymbol{s},\boldsymbol{a})|} \sum\limits_{j=1}^{|R(\boldsymbol{s},\boldsymbol{a})|} R(\boldsymbol{s}, \boldsymbol{a})[j]$
8: **end for**

- Any potential problems?
- Deterministic policy? Not many state-action pairs are visited!

▶▶ $\infty$ many trajectories (use GPI idea)

## Monte Carlo for Q Function

1: **for** $i = 1$ to $\infty$ **do**
2:     Generate trajectory $\tau_i^\pi$ with policy $\pi$
3:     **for** all state-action pairs $(s, a) \in \tau_i$ **do**
4:         $r \leftarrow$ sum of rewards following the <u>first</u> occurrence of $(s, a)$
5:         $R(s, a) \leftarrow [R(s, a), r]$           $\triangleright$ Append $r$ to array
6:     **end for**
7:     $Q^\pi(s, a) \approx \mathbb{E}[R(s, a)|\pi] \approx \frac{1}{|R(s,a)|} \sum_{j=1}^{|R(s,a)|} R(s, a)[j]$
8: **end for**

- Any potential problems?
- Deterministic policy? Not many state-action pairs are visited!
- Ideally, estimate $Q^\pi(s, a)$ for all actions $a \in \mathcal{A}$

▶▶ $\infty$ many trajectories (use GPI idea)

# Monte Carlo for Q Function

1: **for** $i = 1$ to $\infty$ **do**
2:      Generate trajectory $\tau_i^\pi$ with policy $\pi$
3:      **for** all state-action pairs $(\boldsymbol{s}, \boldsymbol{a}) \in \tau_i$ **do**
4:          $r \leftarrow$ sum of rewards following the <u>first</u> occurrence of $(\boldsymbol{s}, \boldsymbol{a})$
5:          $R(\boldsymbol{s}, \boldsymbol{a}) \leftarrow [R(\boldsymbol{s}, \boldsymbol{a}), r]$        $\triangleright$ Append $r$ to array
6:      **end for**
7:      $Q^\pi(\boldsymbol{s}, \boldsymbol{a}) \approx \mathbb{E}[R(\boldsymbol{s}, \boldsymbol{a})|\pi] \approx \frac{1}{|R(\boldsymbol{s},\boldsymbol{a})|} \sum\limits_{j=1}^{|R(\boldsymbol{s},\boldsymbol{a})|} R(\boldsymbol{s}, \boldsymbol{a})[j]$
8: **end for**

- Any potential problems?
- Deterministic policy? Not many state-action pairs are visited!
- Ideally, estimate $Q^\pi(\boldsymbol{s}, \boldsymbol{a})$ for all actions $\boldsymbol{a} \in \mathcal{A}$
▶▶ $\infty$ many trajectories (use GPI idea)
▶▶ Maintain exploration

# Exploration

‣ Exploring starts: non-zero probability that each state-action pair is chosen as the start (difficult in practice, not commonly used)

# Exploration

- Exploring starts: non-zero probability that each state-action pair is chosen as the start (difficult in practice, not commonly used)
- Stochastic policies with non-zero probability on each action
    - $\forall \boldsymbol{s} \in \mathcal{S}, \boldsymbol{a} \in \mathcal{A}: \quad p_\pi(\boldsymbol{a}|\boldsymbol{s}) > 0$
    - Example: $\epsilon$-greedy policies, softmax policies
    - $\epsilon$-greedy:

$$p_\pi(\boldsymbol{a}|\boldsymbol{s}) = \begin{cases} \pi(\boldsymbol{s}) & \text{with probability } 1-\epsilon \\ \mathcal{U}(\mathcal{A}) & \text{with probability } \epsilon \end{cases}$$

# Exploration

- Exploring starts: non-zero probability that each state-action pair is chosen as the start (difficult in practice, not commonly used)
- Stochastic policies with non-zero probability on each action
  - $\forall s \in \mathcal{S}, a \in \mathcal{A}: \quad p_\pi(a|s) > 0$
  - Example: $\epsilon$-greedy policies, softmax policies
  - $\epsilon$-greedy:

  $$p_\pi(a|s) = \begin{cases} \pi(s) & \text{with probability } 1 - \epsilon \\ \mathcal{U}(\mathcal{A}) & \text{with probability } \epsilon \end{cases}$$

▶▶ Trade off exploring the world and exploiting current knowledge

# On-Policy Monte Carlo Control

$$\pi_0 \xrightarrow{E} Q^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} Q^{\pi_1} \xrightarrow{I} \cdots \xrightarrow{E} Q^* \xrightarrow{I} \pi^*$$

‣ Approximate optimal policies
‣ Follow GPI idea: $\epsilon$-greedy policy

# On-Policy Monte Carlo Control

$$\pi_0 \xrightarrow{E} Q^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} Q^{\pi_1} \xrightarrow{I} \cdots \xrightarrow{E} Q^* \xrightarrow{I} \pi^*$$

- Approximate optimal policies
- Follow GPI idea: $\epsilon$-greedy policy
- MC version of policy iteration:
    - Policy evaluation via Monte Carlo estimates (stochastic policy)
      ▶ get $Q^\pi$
    - Policy improvement: select greedy policy with respect to $Q^\pi$

# Off-Policy Monte Carlo Control

‣ So far on-policy: Evaluate $Q^\pi$ and, subsequently, apply $\pi$ when interacting with the environment

# Off-Policy Monte Carlo Control

- So far on-policy: Evaluate $Q^\pi$ and, subsequently, apply $\pi$ when interacting with the environment

- Off-policy: Evaluate $Q^\pi$ but when interacting with the environment, follow $\pi' \neq \pi$

  ▶ Experience only generated from $\pi'$. Can we still estimate $Q^\pi$?

# Off-Policy Monte Carlo Control

- So far on-policy: Evaluate $Q^\pi$ and, subsequently, apply $\pi$ when interacting with the environment

- Off-policy: Evaluate $Q^\pi$ but when interacting with the environment, follow $\pi' \neq \pi$

  ▶ Experience only generated from $\pi'$. Can we still estimate $Q^\pi$?

- Yes, if we ensure that every action under $\pi$ is also (occasionally) taken under $\pi'$, i.e.,

$$\text{if } p_\pi(\boldsymbol{a}|\boldsymbol{s}) > 0 \text{ then } p_{\pi'}(\boldsymbol{a}|\boldsymbol{s}) > 0$$

# Off-Policy Monte Carlo Control

- So far on-policy: Evaluate $Q^\pi$ and, subsequently, apply $\pi$ when interacting with the environment

- Off-policy: Evaluate $Q^\pi$ but when interacting with the environment, follow $\pi' \neq \pi$
  ▶ Experience only generated from $\pi'$. Can we still estimate $Q^\pi$?

- Yes, if we ensure that every action under $\pi$ is also (occasionally) taken under $\pi'$, i.e.,

$$\text{if } p_\pi(\boldsymbol{a}|\boldsymbol{s}) > 0 \text{ then } p_{\pi'}(\boldsymbol{a}|\boldsymbol{s}) > 0$$

- Learning can be slow if the policies are very explorative, i.e., $\epsilon \gg 0$

# Summary: Monte Carlo Methods

‣ Learn optimal behavior from interaction

‣ Easy to focus them on a small set of start states

‣ Incremental implementation (updates after each episode) possible

‣ Exploration required!

‣ No bootstrapping (unlike DP)

‣ Model free (unlike DP)

# Temporal-Difference Learning



- Between MC and DP
- Bootstrapping
- Model-free
- MC waits until the end of the episode to update $V^\pi$, $Q^\pi$
- TD only waits until the next time step. Update value functions based on observed reward and the current estimate of $V^\pi$, $Q^\pi$

# TD Policy Evaluation

Generic Update Rule

$$V(\boldsymbol{s}) \leftarrow V(\boldsymbol{s}) + \alpha\big(\kappa - V(\boldsymbol{s})\big) = (1 - \alpha)V(\boldsymbol{s}) + \alpha\kappa$$

$\kappa$: "target"

# TD Policy Evaluation

## Generic Update Rule

$$V(\boldsymbol{s}) \leftarrow V(\boldsymbol{s}) + \alpha\big(\kappa - V(\boldsymbol{s})\big) = (1-\alpha)V(\boldsymbol{s}) + \alpha\kappa$$

$\kappa$: "target"

$$\begin{aligned}
V^{\pi}(\boldsymbol{s}) &= \mathbb{E}[R(\boldsymbol{s})|\boldsymbol{s}_0 = \boldsymbol{s}] \\
&= \mathbb{E}\big[\sum_k \gamma^k r_{k+1}|\boldsymbol{s}_0 = \boldsymbol{s}\big] \\
&= \mathbb{E}\big[r_1 + \gamma V^{\pi}(\boldsymbol{s}_1)|\boldsymbol{s}_0 = \boldsymbol{s}\big]
\end{aligned}$$

# TD Policy Evaluation

## Generic Update Rule

$$V(\boldsymbol{s}) \leftarrow V(\boldsymbol{s}) + \alpha\big(\kappa - V(\boldsymbol{s})\big) = (1 - \alpha)V(\boldsymbol{s}) + \alpha\kappa$$

$\kappa$: "target"

$$
\begin{aligned}
V^{\pi}(\boldsymbol{s}) &= \mathbb{E}[R(\boldsymbol{s})|\boldsymbol{s}_0 = \boldsymbol{s}] \\
&= \mathbb{E}\big[\sum_k \gamma^k r_{k+1}|\boldsymbol{s}_0 = \boldsymbol{s}\big] \\
&= \mathbb{E}\big[r_1 + \gamma V^{\pi}(\boldsymbol{s}_1)|\boldsymbol{s}_0 = \boldsymbol{s}\big]
\end{aligned}
$$

‣ MC uses an estimate of $R(\boldsymbol{s})$ as target $\kappa$ (sample average)

# TD Policy Evaluation

## Generic Update Rule

$$V(s) \leftarrow V(s) + \alpha\big(\kappa - V(s)\big) = (1 - \alpha)V(s) + \alpha\kappa$$

$\kappa$: "target"

$$\begin{aligned}
V^{\pi}(s) &= \mathbb{E}[R(s)|s_0 = s] \\
&= \mathbb{E}\big[\sum_k \gamma^k r_{k+1}|s_0 = s\big] \\
&= \mathbb{E}\big[r_1 + \gamma V^{\pi}(s_1)|s_0 = s\big]
\end{aligned}$$

- MC uses an estimate of $R(s)$ as target $\kappa$ (sample average)
- DP uses an estimate of $r_1 + \gamma V^{\pi}(s_1)$ as target $\kappa$ ($V_k$ instead of $V^{\pi}$)

# TD Policy Evaluation

## Generic Update Rule

$$V(\boldsymbol{s}) \leftarrow V(\boldsymbol{s}) + \alpha\big(\kappa - V(\boldsymbol{s})\big) = (1 - \alpha)V(\boldsymbol{s}) + \alpha\kappa$$

$\kappa$: "target"

$$
\begin{aligned}
V^\pi(\boldsymbol{s}) &= \mathbb{E}[R(\boldsymbol{s})|\boldsymbol{s}_0 = \boldsymbol{s}] \\
&= \mathbb{E}\big[\sum_k \gamma^k r_{k+1}|\boldsymbol{s}_0 = \boldsymbol{s}\big] \\
&= \mathbb{E}\big[r_1 + \gamma V^\pi(\boldsymbol{s}_1)|\boldsymbol{s}_0 = \boldsymbol{s}\big]
\end{aligned}
$$

- MC uses an estimate of $R(\boldsymbol{s})$ as target $\kappa$ (sample average)
- DP uses an estimate of $r_1 + \gamma V^\pi(\boldsymbol{s}_1)$ as target $\kappa$ ($V_k$ instead of $V^\pi$)
- TD target: $\kappa = r_1 + \gamma V^\pi(\boldsymbol{s}_1)$

# TD Policy Evaluation

## Generic Update Rule

$$V(s) \leftarrow V(s) + \alpha\big(\kappa - V(s)\big) = (1 - \alpha)V(s) + \alpha\kappa$$

$\kappa$: "target"

$$V^{\pi}(s) = \mathbb{E}[R(s)|s_0 = s]$$
$$= \mathbb{E}\big[\sum_k \gamma^k r_{k+1}|s_0 = s\big]$$
$$= \mathbb{E}\big[r_1 + \gamma V^{\pi}(s_1)|s_0 = s\big]$$

- MC uses an estimate of $R(s)$ as target $\kappa$ (sample average)
- DP uses an estimate of $r_1 + \gamma V^{\pi}(s_1)$ as target $\kappa$ ($V_k$ instead of $V^{\pi}$)
- TD target: $\kappa = r_1 + \gamma V^{\pi}(s_1)$
  - Approximate $V^{\pi}$ by $V_k$ (same as DP)

# TD Policy Evaluation

## Generic Update Rule

$$V(s) \leftarrow V(s) + \alpha\big(\kappa - V(s)\big) = (1 - \alpha)V(s) + \alpha\kappa$$

$\kappa$: "target"

$$V^{\pi}(s) = \mathbb{E}[R(s)|s_0 = s]$$
$$= \mathbb{E}\big[\sum_k \gamma^k r_{k+1}|s_0 = s\big]$$
$$= \mathbb{E}\big[r_1 + \gamma V^{\pi}(s_1)|s_0 = s\big]$$

- MC uses an estimate of $R(s)$ as target $\kappa$ (sample average)
- DP uses an estimate of $r_1 + \gamma V^{\pi}(s_1)$ as target $\kappa$ ($V_k$ instead of $V^{\pi}$)
- TD target: $\kappa = r_1 + \gamma V^{\pi}(s_1)$
  - Approximate $V^{\pi}$ by $V_k$ (same as DP)
  - Sample trajectories (same as MC)
  ▶▶ Combine MC sampling with DP bootstrapping

## TD(0)

1: **Init.:** Set $V(s)$ arbitrarily
2: **repeat** for each episode
3:     $a \leftarrow p_\pi(a|s)$                  ▷ Sample action in current state
4:     Apply $a$, observe $r, s'$                  ▷ Transition to next state
5:     $V(s) \leftarrow V(s) + \alpha\big(r + \gamma V(s') - V(s)\big)$                  ▷ Update $V$
6:     $s \leftarrow s'$                  ▷ Re-set current state
7: **until** $s$ is terminal

## TD(0)

1: **Init.:** Set $V(s)$ arbitrarily
2: **repeat** for each episode
3:     $a \leftarrow p_\pi(a|s)$                   ▷ Sample action in current state
4:     Apply $a$, observe $r, s'$                   ▷ Transition to next state
5:     $V(s) \leftarrow V(s) + \alpha\big(r + \gamma V(s') - V(s)\big)$                   ▷ Update $V$
6:     $s \leftarrow s'$                   ▷ Re-set current state
7: **until** $s$ is terminal

- Temporal difference error: $r + \gamma V(s') - V(s)$     $(= V^{\text{new}} - V^{\text{old}})$
- Difference to MC and DP:
    - MC waits until the end of the episode to update $V$
    - DP needs complete distribution $p(s'|s, a)$ of successor states to update $V$

# TD($\lambda$)

- TD(0) uses 1-step returns to update $V^\pi$
- MC uses full trajectories to update $V^\pi$
- TD($\lambda$), $\lambda \in [0, 1]$, blends between them
  - $\lambda = 0$: TD(0), $\quad \lambda = 1$: MC
  - TD($\lambda$) update rule given as a mixture of multi-step returns
  - Mixing coefficients $(1 - \lambda)\lambda^k$, $k \geqslant 0$

# SARSA: On-Policy TD Control

- Thus far, we only learned $V^\pi$ using TD
- Not very useful without model when we want to do control
  ▶ Learn $Q^\pi$

# SARSA: On-Policy TD Control

- Thus far, we only learned $V^\pi$ using TD
- Not very useful without model when we want to do control
  ▶ Learn $Q^\pi$

TD Update for $Q^\pi$

$$Q(s, a) \leftarrow \underbrace{Q(s, a) + \alpha \left( r(s, a, s') + \gamma Q(s', a') - Q(s, a) \right)}_{=(1-\alpha)Q^{\text{old}} + \alpha Q^{\text{new}}}$$

# SARSA: On-Policy TD Control

- Thus far, we only learned $V^\pi$ using TD
- Not very useful without model when we want to do control
  ▶▶ Learn $Q^\pi$

TD Update for $Q^\pi$

$$Q(s, a) \leftarrow \underbrace{Q(s, a) + \alpha\big(r(s, a, s') + \gamma Q(s', a') - Q(s, a)\big)}_{=(1-\alpha)Q^{\text{old}} + \alpha Q^{\text{new}}}$$

- Update rule needs $(s, a, r, s', a')$ ▶▶ SARSA
- On-policy algorithm
- $a'$ chosen from $s'$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
- Convergence proofs for $\epsilon$-greedy policies

# SARSA: On-Policy TD Control

- Thus far, we only learned $V^\pi$ using TD
- Not very useful without model when we want to do control
  ▶ Learn $Q^\pi$

### TD Update for $Q^\pi$

$$Q(s, a) \leftarrow \underbrace{Q(s, a) + \alpha \left( r(s, a, s') + \gamma Q(s', a') - Q(s, a) \right)}_{= (1-\alpha) Q^{\text{old}} + \alpha Q^{\text{new}}}$$

- Update rule needs $(s, a, r, s', a')$ ▶ SARSA
- On-policy algorithm
- $a'$ chosen from $s'$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
- Convergence proofs for $\epsilon$-greedy policies

**DEMO gridworld**

# Q-Learning: Off-Policy TD Control

## TD Update for $Q^*$

$$Q(s, a) \leftarrow Q(s, a) + \alpha\big(r(s, a, s') + \gamma\max_{a'}Q(s', a') - Q(s, a)\big)$$

# Q-Learning: Off-Policy TD Control

## TD Update for $Q^*$

$$Q(s, a) \leftarrow Q(s, a) + \alpha\big(r(s, a, s') + \gamma\max_{a'}Q(s', a') - Q(s, a)\big)$$

- Off-policy TD control
- SARSA: learn $Q^\pi$, Q-learning: learn $Q^*$
- Update the value function $Q$ independent of the policy the agent actually follows to generate the samples (max ...)
- For convergence: keep updating all state-action pairs

# Q-Learning: Off-Policy TD Control

### TD Update for $Q^*$

$$Q(s, a) \leftarrow Q(s, a) + \alpha\big(r(s, a, s') + \gamma \max_{a'} Q(s', a') - Q(s, a)\big)$$

- Off-policy TD control
- SARSA: learn $Q^\pi$, Q-learning: learn $Q^*$
- Update the value function $Q$ independent of the policy the agent actually follows to generate the samples (max ...)
- For convergence: keep updating all state-action pairs

**DEMO gridworld 2**

# RL in Continuous Spaces: Function Approximation

- So far: discrete states and actions
  ▶ Table representation sufficient
- In continuous spaces: Function approximation for better generalization:

# Function Approximation

- Typically: (linear) basis function representation

$$V(\boldsymbol{s}) = \sum_i \theta_i \phi_i(\boldsymbol{s})$$

- Basis functions $\phi_i$ are fixed, only parameters $\theta_i$ need to be learned

# Function Approximation

‣ Typically: (linear) basis function representation

$$V(\boldsymbol{s}) = \sum_i \theta_i \phi_i(\boldsymbol{s})$$

‣ Basis functions $\phi_i$ are fixed, only parameters $\theta_i$ need to be learned

‣ Gradient descent to update the parameters. Example TD(0):

$$\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \alpha \big(v_k - V_k(\boldsymbol{s}_k)\big) \frac{\partial V_k(\boldsymbol{s}_k)}{\partial \boldsymbol{\theta}_k}$$

$v_k$: approximation/estimate of $V^\pi(\boldsymbol{s}_k)$, e.g., MC estimate

$\alpha$: learning rate

‣ Convergence if $v_k$ is unbiased, i.e., $\mathbb{E}[v_k] = V^\pi(\boldsymbol{s}_k)$

# Summary: RL with Value Functions

- Learn policies exploiting properties of the value functions $V$, $Q$
- Bellman equations/optimality principle
- Policy evaluation/improvement
- Exact solution: Dynamic programming
- Approximation solution: MC, TD
- Exploration
- Function approximation

# Applications



- ‣ Board games (e.g., Tesauro, Silver, Riedmiller)
- ‣ Power systems (e.g., D. Ernst)
- ‣ Robocup (e.g., Stone, Riedmiller)
- ‣ Scheduling tasks (e.g., W. Powell)

Motivation

Reinforcement Learning Set-Up

Value Function Methods

Policy Search

# RL without Value Functions: Policy Search

‣ Value function approximation is hard, especially in continuous
  domains
‣ Search directly in policy space?

# RL without Value Functions: Policy Search

‣ Value function approximation is hard, especially in continuous domains
‣ Search directly in policy space?
‣ Only feasible in a restricted class $\Pi$ of policies

$$\pi = \pi(\boldsymbol{\theta})$$

$$V^\pi(\boldsymbol{\theta}) = \mathbb{E}\big[\sum_{k=0}^{T} r_{k+1}|\boldsymbol{\theta}\big]$$

‣ Consider an episodic set-up, i.e., $p(\boldsymbol{s}_0)$ is given

# RL without Value Functions: Policy Search

- Value function approximation is hard, especially in continuous domains
- Search directly in policy space?
- Only feasible in a restricted class $\Pi$ of policies

$$\pi = \pi(\boldsymbol{\theta})$$

$$V^\pi(\boldsymbol{\theta}) = \mathbb{E}\Big[\sum_{k=0}^{T} r_{k+1} | \boldsymbol{\theta}\Big]$$

- Consider an episodic set-up, i.e., $p(\boldsymbol{s}_0)$ is given

## Objective (Policy Search)

For a given class $\Pi(\boldsymbol{\theta})$ of policies, find an optimal policy

$$\pi^* \in \arg \max_{\pi \in \Pi(\boldsymbol{\theta})} V^\pi(\boldsymbol{\theta})$$

# Categorization of Policy Search Methods



Model-free policy search

Model-based policy search

- Data: $(s_i, a_i, r_i), i = 1, \ldots, n$
- Model-based vs. model-free policy search
- Policy evaluation (red), policy improvement (green)

# Model-based Set-up: Interaction and Simulation



*interaction*   *simulation*

Two alternating phases:

‣ **Interaction:** internal model is refined using experience from interacting with the real system

# Model-based Set-up: Interaction and Simulation



*interaction*                    *simulation*

Two alternating phases:

- **Interaction:** internal model is refined using experience from interacting with the real system
- **Simulation:** internal model simulates consequences of actions in the real system, policy is refined

# Model Building



Model-free policy search

Model-based policy search

# Model Learning



*interaction*

*simulation*

- ‣ Pro:

# Model Learning



*interaction*  *simulation*

- Pro: No "real" experiments with robot for policy evaluation and improvement (just simulate!) ▶▶ Protect hardware
- Con:

# Model Learning



*interaction*  *simulation*

- Pro: No "real" experiments with robot for policy evaluation and improvement (just simulate!) ▶▶ Protect hardware
- Con: Model errors ▶▶ Effects of "wrong" models?

# Model Learning



*interaction*

*simulation*

- Pro: No "real" experiments with robot for policy evaluation and improvement (just simulate!) ▶ Protect hardware
- Con: Model errors ▶ Effects of "wrong" models?
- What are good models?

# Model Learning

Model learning problem: Find a function $f : x \mapsto f(x) = y$



Observed function values

# Model Learning

Model learning problem: Find a function $f : x \mapsto f(x) = y$



Plausible function approximators

# Model Learning

Model learning problem: Find a function $f : x \mapsto f(x) = y$



Plausible function approximators

**Predictions? Decision Making?**

# Model Learning

Model learning problem: Find a function $f : x \mapsto f(x) = y$



Plausible function approximators

**Predictions? Decision Making? Model Errors!**

# Model Learning

Model learning problem: Find a function $f : x \mapsto f(x) = y$



Distribution over plausible functions

# Model Learning

Model learning problem: Find a function $f : x \mapsto f(x) = y$



Distribution over plausible functions

▶▶ Express uncertainty about the underlying function
▶▶ Bayesian models (Bayesian linear regression, Gaussian process, ...)

# Useful Models

- Probabilistic models!
- Reduce model errors and simulation/optimization bias
- Examples of probabilistic models
  - Bayesian linear regression
  - Gaussian process

# Bayesian Linear Regression: Model

- Model: $y = \theta^\top \phi(x) = \sum_i \theta_i \phi_i(x)$
- $\phi(x)$: "Features", e.g., $\phi(x) = [x, x^2]^\top$

# Bayesian Linear Regression: Model

- Model: $y = \boldsymbol{\theta}^\top \phi(\boldsymbol{x}) = \sum_i \theta_i \phi_i(\boldsymbol{x})$
- $\phi(\boldsymbol{x})$: "Features", e.g., $\phi(x) = [x, x^2]^\top$
- Distribution over model parameters $\boldsymbol{\theta}$:

$$p(\boldsymbol{\theta}) = \mathcal{N}\left(\boldsymbol{\theta} \,|\, \boldsymbol{m}, \boldsymbol{S}\right)$$

# Bayesian Linear Regression: Model

- Model: $y = \boldsymbol{\theta}^\top \phi(\boldsymbol{x}) = \sum_i \theta_i \phi_i(\boldsymbol{x})$

- $\phi(\boldsymbol{x})$: "Features", e.g., $\phi(x) = [x, x^2]^\top$

- Distribution over model parameters $\boldsymbol{\theta}$:

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} \mid \boldsymbol{m}, \boldsymbol{S})$$

▶ For any $\boldsymbol{\theta}$, one particular function is defined
▶ Distribution over $\boldsymbol{\theta}$ induces distribution over functions

# Bayesian Linear Regression: Model

- Model: $y = \boldsymbol{\theta}^\top \phi(\boldsymbol{x}) = \sum_i \theta_i \phi_i(\boldsymbol{x})$
- $\phi(\boldsymbol{x})$: "Features", e.g., $\phi(x) = [x, \, x^2]^\top$
- Distribution over model parameters $\boldsymbol{\theta}$:

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} \,|\, \boldsymbol{m}, \, \boldsymbol{S})$$

▶ For any $\boldsymbol{\theta}$, one particular function is defined
▶ Distribution over $\boldsymbol{\theta}$ induces distribution over functions

**DEMO Bayesian regression**

# Bayesian Linear Regression: Predictions

‣ Model: $y = \phi^\top(x)\theta, \quad \theta \sim \mathcal{N}(m, S)$

# Bayesian Linear Regression: Predictions

- Model: $y = \phi^\top(x)\theta, \quad \theta \sim \mathcal{N}(m, S)$
- Predict function values $y = [y_1, \ldots, y_n]^\top$ at inputs
  $X = [x_1, \ldots, x_n]$

Define $\boldsymbol{\Phi} = \phi(X) \Rightarrow y = \boldsymbol{\Phi}^\top \theta$

$$\mathbb{E}[y] =$$
$$\mathbb{V}[y] =$$

# Bayesian Linear Regression: Predictions

- Model: $y = \phi^\top(x)\theta, \quad \theta \sim \mathcal{N}(m, S)$
- Predict function values $y = [y_1, \ldots, y_n]^\top$ at inputs $X = [x_1, \ldots, x_n]$

Define $\boldsymbol{\Phi} = \phi(X) \Rightarrow y = \boldsymbol{\Phi}^\top \theta$
$$\mathbb{E}[y] = \mathbb{E}[\boldsymbol{\Phi}^\top \theta] = \boldsymbol{\Phi}^\top \mathbb{E}[\theta] = \boldsymbol{\Phi}^\top m$$
$$\mathbb{V}[y] =$$

# Bayesian Linear Regression: Predictions

- Model: $y = \phi^\top(x)\theta, \quad \theta \sim \mathcal{N}(m, S)$
- Predict function values $y = [y_1, \ldots, y_n]^\top$ at inputs $X = [x_1, \ldots, x_n]$

Define $\boldsymbol{\Phi} = \phi(X) \Rightarrow y = \boldsymbol{\Phi}^\top \theta$

$$\mathbb{E}[y] = \mathbb{E}[\boldsymbol{\Phi}^\top \theta] = \boldsymbol{\Phi}^\top \mathbb{E}[\theta] = \boldsymbol{\Phi}^\top m$$

$$\mathbb{V}[y] = \mathbb{V}[\boldsymbol{\Phi}^\top \theta] = \boldsymbol{\Phi}^\top \mathbb{V}[\theta]\boldsymbol{\Phi} = \boldsymbol{\Phi}^\top S \boldsymbol{\Phi}$$

# Introduction to Gaussian Processes

‣ Generalization of Bayesian linear regression
‣ Nonparametric Bayesian regression method
‣ Probability distribution over functions
‣ Fully specified by
  ‣ Mean function *m* (average function)
  ‣ Covariance function/kernel *k* (assumptions on structure)

$$\text{Cov}[f(\boldsymbol{x}_p), f(\boldsymbol{x}_q)] = k(\boldsymbol{x}_p, \boldsymbol{x}_q)$$

# Introduction to Gaussian Processes

- Generalization of Bayesian linear regression
- Nonparametric Bayesian regression method
- Probability distribution over functions
- Fully specified by
    - Mean function $m$ (average function)
    - Covariance function/kernel $k$ (assumptions on structure)

$$\mathrm{Cov}[f(\boldsymbol{x}_p), f(\boldsymbol{x}_q)] = k(\boldsymbol{x}_p, \boldsymbol{x}_q)$$

- Posterior predictive distribution at $\boldsymbol{x}_*$ is Gaussian:

$$p(f(\boldsymbol{x}_*) | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) = \mathcal{N}\big(f(\boldsymbol{x}_*) | m(\boldsymbol{x}_*), \sigma^2(\boldsymbol{x}_*)\big)$$

Test input   Training data

# Gaussian Process: Definition

### Definition

A Gaussian process is a collection of random variables, any finite number of which has a joint Gaussian distribution.

# Gaussian Process: Definition

### Definition

A Gaussian process is a collection of random variables, any finite number of which has a joint Gaussian distribution.

- Look at Gaussian distributions of function values $f_1, f_2, \ldots$
- All of them are jointly Gaussian distributed

$$\mathbb{E}[f(\boldsymbol{x})] = m(\boldsymbol{x})$$
$$\mathrm{Cov}[f(\boldsymbol{x}_i), f(\boldsymbol{x}_j)] = k(\boldsymbol{x}_i, \boldsymbol{x}_j)$$

# Gaussian Process: Predictions

‣ Given a training set $\left(x_i, f(x_i)\right)_{i=1}^{n}$, we can predict function values $f_{*_j}$ at test inputs $x_{*_j}$

‣ First, compute the joint distribution:

$$p(f, f_* | X, X_*) = \mathcal{N}\left(\begin{bmatrix} m(X) \\ m(X_*) \end{bmatrix}, \begin{bmatrix} K & k(X, X_*) \\ k(X_*, X) & K_* \end{bmatrix}\right)$$

$$K_{ij} = k(x_i, x_j) = \text{Cov}[f(x_i), f(x_j)]$$

# Gaussian Process: Predictions

- Given a training set $(x_i, f(x_i))_{i=1}^{n}$, we can predict function values $f_{*_j}$ at test inputs $x_{*_j}$

- First, compute the joint distribution:

$$p(f, f_* | X, X_*) = \mathcal{N} \left( \begin{bmatrix} m(X) \\ m(X_*) \end{bmatrix}, \begin{bmatrix} K & k(X, X_*) \\ k(X_*, X) & K_* \end{bmatrix} \right)$$
$$K_{ij} = k(x_i, x_j) = \mathrm{Cov}[f(x_i), f(x_j)]$$

- Second, compute the conditional $p(f_* | X, X_*, f)$ by plain Gaussian conditioning:

# Gaussian Process: Predictions

- Given a training set $(x_i, f(x_i))_{i=1}^{n}$, we can predict function values $f_{*_j}$ at test inputs $x_{*_j}$
- First, compute the joint distribution:

$$p(f, f_* | X, X_*) = \mathcal{N}\left(\begin{bmatrix} m(X) \\ m(X_*) \end{bmatrix}, \begin{bmatrix} K & k(X, X_*) \\ k(X_*, X) & K_* \end{bmatrix}\right)$$

$$K_{ij} = k(x_i, x_j) = \text{Cov}[f(x_i), f(x_j)]$$

- Second, compute the conditional $p(f_* | X, X_*, f)$ by plain Gaussian conditioning:

$$p(f_* | X, X_*, f) = \mathcal{N}(\mu_*, \Sigma_*)$$

$$\mu_* = m(X_*) + k(X_*, X)K^{-1}(f - m(X))$$

$$\Sigma_* = K_* - k(X_*, X)K^{-1}k(X, X_*)$$

# Intuitive Introduction to Gaussian Processes



Prior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\boldsymbol{x}_*)|\varnothing] = m(\boldsymbol{x}_*) = 0$$
$$\mathbb{V}[f(\boldsymbol{x}_*)|\varnothing] = \sigma^2(\boldsymbol{x}_*) = \mathrm{Cov}[f(\boldsymbol{x}_*), f(\boldsymbol{x}_*)] = k(\boldsymbol{x}_*, \boldsymbol{x}_*)$$

# Intuitive Introduction to Gaussian Processes



Prior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\boldsymbol{x}_*)|\varnothing] = m(\boldsymbol{x}_*) = 0$$
$$\mathbb{V}[f(\boldsymbol{x}_*)|\varnothing] = \sigma^2(\boldsymbol{x}_*) = \text{Cov}[f(\boldsymbol{x}_*), f(\boldsymbol{x}_*)] = k(\boldsymbol{x}_*, \boldsymbol{x}_*)$$

# Intuitive Introduction to Gaussian Processes



Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = m(\boldsymbol{x}_*) = k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1}\boldsymbol{y}$$
$$\mathbb{V}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = \sigma^2(\boldsymbol{x}_*) = k(\boldsymbol{x}_*, \boldsymbol{x}_*) - k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1}k(\boldsymbol{X}, \boldsymbol{x}_*)$$

# Intuitive Introduction to Gaussian Processes



Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = m(\boldsymbol{x}_*) = k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1}\boldsymbol{y}$$
$$\mathbb{V}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = \sigma^2(\boldsymbol{x}_*) = k(\boldsymbol{x}_*, \boldsymbol{x}_*) - k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1}k(\boldsymbol{X}, \boldsymbol{x}_*)$$

# Intuitive Introduction to Gaussian Processes



Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = m(\boldsymbol{x}_*) = k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1} \boldsymbol{y}$$
$$\mathbb{V}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = \sigma^2(\boldsymbol{x}_*) = k(\boldsymbol{x}_*, \boldsymbol{x}_*) - k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1} k(\boldsymbol{X}, \boldsymbol{x}_*)$$

# Intuitive Introduction to Gaussian Processes



Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\boldsymbol{x}_*)|\boldsymbol{X},\boldsymbol{y}] = m(\boldsymbol{x}_*) = k(\boldsymbol{X},\boldsymbol{x}_*)^\top k(\boldsymbol{X},\boldsymbol{X})^{-1}\boldsymbol{y}$$
$$\mathbb{V}[f(\boldsymbol{x}_*)|\boldsymbol{X},\boldsymbol{y}] = \sigma^2(\boldsymbol{x}_*) = k(\boldsymbol{x}_*,\boldsymbol{x}_*) - k(\boldsymbol{X},\boldsymbol{x}_*)^\top k(\boldsymbol{X},\boldsymbol{X})^{-1}k(\boldsymbol{X},\boldsymbol{x}_*)$$

# Intuitive Introduction to Gaussian Processes



Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\boldsymbol{x}_*)|\boldsymbol{X},\boldsymbol{y}] = m(\boldsymbol{x}_*) = k(\boldsymbol{X},\boldsymbol{x}_*)^\top k(\boldsymbol{X},\boldsymbol{X})^{-1}\boldsymbol{y}$$
$$\mathbb{V}[f(\boldsymbol{x}_*)|\boldsymbol{X},\boldsymbol{y}] = \sigma^2(\boldsymbol{x}_*) = k(\boldsymbol{x}_*,\boldsymbol{x}_*) - k(\boldsymbol{X},\boldsymbol{x}_*)^\top k(\boldsymbol{X},\boldsymbol{X})^{-1}k(\boldsymbol{X},\boldsymbol{x}_*)$$

# Intuitive Introduction to Gaussian Processes



Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = m(\boldsymbol{x}_*) = k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1}\boldsymbol{y}$$
$$\mathbb{V}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = \sigma^2(\boldsymbol{x}_*) = k(\boldsymbol{x}_*, \boldsymbol{x}_*) - k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1}k(\boldsymbol{X}, \boldsymbol{x}_*)$$

# Intuitive Introduction to Gaussian Processes



Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = m(\boldsymbol{x}_*) = k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1} \boldsymbol{y}$$
$$\mathbb{V}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = \sigma^2(\boldsymbol{x}_*) = k(\boldsymbol{x}_*, \boldsymbol{x}_*) - k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1} k(\boldsymbol{X}, \boldsymbol{x}_*)$$

# Intuitive Introduction to Gaussian Processes



Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = m(\boldsymbol{x}_*) = k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1} \boldsymbol{y}$$
$$\mathbb{V}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = \sigma^2(\boldsymbol{x}_*) = k(\boldsymbol{x}_*, \boldsymbol{x}_*) - k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1} k(\boldsymbol{X}, \boldsymbol{x}_*)$$

# Intuitive Introduction to Gaussian Processes



Posterior belief about the function

Predictive (marginal) mean and variance:

$$\mathbb{E}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = m(\boldsymbol{x}_*) = k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1}\boldsymbol{y}$$
$$\mathbb{V}[f(\boldsymbol{x}_*)|\boldsymbol{X}, \boldsymbol{y}] = \sigma^2(\boldsymbol{x}_*) = k(\boldsymbol{x}_*, \boldsymbol{x}_*) - k(\boldsymbol{X}, \boldsymbol{x}_*)^\top k(\boldsymbol{X}, \boldsymbol{X})^{-1}k(\boldsymbol{X}, \boldsymbol{x}_*)$$

# Properties

- Universal function approximator ▶▶ extremely expressive
- Model gives "free" variance estimates

# Properties

- Universal function approximator ▶▶ extremely expressive
- Model gives "free" variance estimates
- Computationally involved:
  - Training: $\mathcal{O}(N^3)$ ▶▶ Repeated inversion of $N \times N$ matrix
  - Mean prediction: $\mathcal{O}(N)$ ▶▶ Scalar product
  - Variance prediction: $\mathcal{O}(N^2)$ ▶▶ Matrix-vector multiplication

# Properties

- Universal function approximator ▶▶ extremely expressive
- Model gives "free" variance estimates
- Computationally involved:
  - Training: $\mathcal{O}(N^3)$ ▶▶ Repeated inversion of $N \times N$ matrix
  - Mean prediction: $\mathcal{O}(N)$ ▶▶ Scalar product
  - Variance prediction: $\mathcal{O}(N^2)$ ▶▶ Matrix-vector multiplication
- Sparse approximations exist
- Code/book online: `http://www.gaussianprocess.org`

# Policy Evaluation



Model-free policy search

Model-based policy search

Policy evaluation: Compute expected long-term reward

- Stochastic trajectory evaluation
- Deterministic trajectory evaluation

# Policy Evaluation

- Stochastic inference (sampling) using either the learned model (simulator) or the real system

- Deterministic inference—only with a learned model

# Stochastic Inference

- Sample trajectories $(s_i, a_i, r_i)$ ▶▶ Monte Carlo
- Conceptually very simple
- Requires a lot of "interactions" (if you don't have a model or a good simulator)
  ▶▶ Potentially impractical (e.g., in robotics)

# Deterministic Inference



- ‣ Analytically propagate uncertainty through the model
- ‣ Computationally/mathematically more involved
- ‣ Can't do this for arbitrary systems, but for some.

# Deterministic Inference: Example

Linear system

$$p(\boldsymbol{s}_t) = \mathcal{N}\big(\boldsymbol{s}_t \,|\, \boldsymbol{m}_t, \, \boldsymbol{S}_t\big)$$
$$\boldsymbol{s}_{t+1} = \boldsymbol{A}\boldsymbol{s}_t$$

Successor state distribution $p(\boldsymbol{s}_{t+1})$?

# Deterministic Inference: Example

Linear system

$$p(\boldsymbol{s}_t) = \mathcal{N}\big(\boldsymbol{s}_t \,|\, \boldsymbol{m}_t, \boldsymbol{S}_t\big)$$
$$\boldsymbol{s}_{t+1} = \boldsymbol{A}\boldsymbol{s}_t$$

Successor state distribution $p(\boldsymbol{s}_{t+1})$?

$$p(\boldsymbol{s}_{t+1}) = \mathcal{N}\big(\boldsymbol{s}_{t+1} \,|\, \boldsymbol{m}_{t+1}, \boldsymbol{S}_{t+1}\big)$$
$$\boldsymbol{m}_{t+1} = \boldsymbol{A}\boldsymbol{m}_t, \quad \boldsymbol{S}_{t+1} = \boldsymbol{A}\boldsymbol{S}_t\boldsymbol{A}^\top$$

‣ In nonlinear/non-Gaussian systems, we need approximations (e.g., linearization, moment matching)

# Policy Improvement



Model-free policy search

Model-based policy search

Policy improvement (green)

- Policy gradients
- Expectation Maximization
- Information theory

# Policy Search: Policy Improvement

## Objective

Find policy parameters $\boldsymbol{\theta}^*$, which maximize the expected long-term reward

$$V^\pi(\boldsymbol{\theta}) = \mathbb{E}\Big[\sum_{k=0}^{T} \gamma^k r_{k+1} | \boldsymbol{\theta}\Big], \qquad \boldsymbol{s}_0 \sim p(\boldsymbol{s}_0)$$

- No global value function model $V^\pi$ or $Q^\pi$
- Search directly in (policy) parameter space

# Policy Search: Policy Improvement

## Objective

Find policy parameters $\boldsymbol{\theta}^*$, which maximize the expected long-term reward

$$V^\pi(\boldsymbol{\theta}) = \mathbb{E}\Big[\sum_{k=0}^{T} \gamma^k r_{k+1} | \boldsymbol{\theta}\Big], \qquad \boldsymbol{s}_0 \sim p(\boldsymbol{s}_0)$$

- No global value function model $V^\pi$ or $Q^\pi$
- Search directly in (policy) parameter space

▶▶ One way: gradient-based optimization

- Compute $V^\pi$ with corresponding gradients $\mathrm{d}V^\pi/\mathrm{d}\boldsymbol{\theta}$
- Gradient-based optimizer for maximization (e.g., CG, BFGS)

# Gradient Estimation for Stochastic Inference

# Gradient Estimation for Stochastic Inference

▸ Finite (central) differences

$$\frac{\mathrm{d}V^\pi(\boldsymbol{\theta})}{\mathrm{d}\boldsymbol{\theta}} \approx \frac{V^\pi(\boldsymbol{\theta} + \epsilon) - V^\pi(\boldsymbol{\theta} - \epsilon)}{2\epsilon}$$

  ▸ Model $p(\boldsymbol{s}_{k+1}|\boldsymbol{s}_k, \boldsymbol{a}_k)$ not required but useful
  ▸ Large variance of estimator ▶▶ many samples needed

# Gradient Estimation for Stochastic Inference

‣ Finite (central) differences

$$\frac{\mathrm{d}V^{\pi}(\boldsymbol{\theta})}{\mathrm{d}\boldsymbol{\theta}} \approx \frac{V^{\pi}(\boldsymbol{\theta} + \epsilon) - V^{\pi}(\boldsymbol{\theta} - \epsilon)}{2\epsilon}$$

  ‣ Model $p(\boldsymbol{s}_{k+1}|\boldsymbol{s}_k, \boldsymbol{a}_k)$ not required but useful
  ‣ Large variance of estimator ▶▶ many samples needed

‣ PEGASUS trick:

  ‣ Fix the random seed and re-set
  ‣ Smaller variance of the estimate of $V^{\pi}$ and its gradient
  ‣ No model $p(\boldsymbol{s}_{k+1}|\boldsymbol{s}_k, \boldsymbol{a}_k)$ required
  ‣ Only with simulator (where we can run exactly the same experiment)

# Gradient Estimation for Deterministic Inference

- Finite differences and PEGASUS still work

- Analytic (=exact) gradients. Example (assume $r = r(s)$):

$$\frac{\mathrm{d}V(\boldsymbol{\theta})}{\mathrm{d}\boldsymbol{\theta}} = \sum_t \gamma^t \frac{\mathrm{d}r(s_t)}{\mathrm{d}\boldsymbol{\theta}} = \sum_t \gamma^t \frac{\partial r(s_t)}{\partial s_t} \frac{\mathrm{d}s_t}{\mathrm{d}\boldsymbol{\theta}}$$

$$= \sum_t \gamma^t \frac{\partial r(s_t)}{\partial s_t} \left( \frac{\partial s_t}{\partial s_{t-1}} \frac{\mathrm{d}s_{t-1}}{\mathrm{d}\boldsymbol{\theta}} + \frac{\partial s_t}{\partial a_{t-1}} \frac{\partial a_{t-1}}{\partial \boldsymbol{\theta}} \right)$$

# Gradient Estimation for Deterministic Inference

- Finite differences and PEGASUS still work
- Analytic (=exact) gradients. Example (assume $r = r(s)$):

$$\frac{\mathrm{d}V(\boldsymbol{\theta})}{\mathrm{d}\boldsymbol{\theta}} = \sum_t \gamma^t \frac{\mathrm{d}r(s_t)}{\mathrm{d}\boldsymbol{\theta}} = \sum_t \gamma^t \frac{\partial r(s_t)}{\partial s_t} \frac{\mathrm{d}s_t}{\mathrm{d}\boldsymbol{\theta}}$$

$$= \sum_t \gamma^t \frac{\partial r(s_t)}{\partial s_t} \left( \frac{\partial s_t}{\partial s_{t-1}} \frac{\mathrm{d}s_{t-1}}{\mathrm{d}\boldsymbol{\theta}} + \frac{\partial s_t}{\partial a_{t-1}} \frac{\partial a_{t-1}}{\partial \boldsymbol{\theta}} \right)$$

- Requires
    - Forward model $s_t = f(s_{t-1}, a_{t-1})$
    - Differentiable policy $a = \pi(s, \boldsymbol{\theta})$

# Gradient Estimation for Deterministic Inference

‣ Finite differences and PEGASUS still work

‣ Analytic (=exact) gradients. Example (assume $r = r(\boldsymbol{s})$):

$$\frac{\mathrm{d}V(\boldsymbol{\theta})}{\mathrm{d}\boldsymbol{\theta}} = \sum_t \gamma^t \frac{\mathrm{d}r(\boldsymbol{s}_t)}{\mathrm{d}\boldsymbol{\theta}} = \sum_t \gamma^t \frac{\partial r(\boldsymbol{s}_t)}{\partial \boldsymbol{s}_t} \frac{\mathrm{d}\boldsymbol{s}_t}{\mathrm{d}\boldsymbol{\theta}}$$

$$= \sum_t \gamma^t \frac{\partial r(\boldsymbol{s}_t)}{\partial \boldsymbol{s}_t} \left( \frac{\partial \boldsymbol{s}_t}{\partial \boldsymbol{s}_{t-1}} \frac{\mathrm{d}\boldsymbol{s}_{t-1}}{\mathrm{d}\boldsymbol{\theta}} + \frac{\partial \boldsymbol{s}_t}{\partial \boldsymbol{a}_{t-1}} \frac{\partial \boldsymbol{a}_{t-1}}{\partial \boldsymbol{\theta}} \right)$$

‣ Requires
  ‣ Forward model $\boldsymbol{s}_t = f(\boldsymbol{s}_{t-1}, \boldsymbol{a}_{t-1})$
  ‣ Differentiable policy $\boldsymbol{a} = \pi(\boldsymbol{s}, \boldsymbol{\theta})$

‣ Mathematically more involved

‣ Gradients are exact (no variance): single trajectory evaluation

# Policy Search



Model-free policy search

Model-based policy search

# Applications in Robotics and Control



- ‣ Cart-pole: e.g., Riedmiller (2005), Deisenroth & Rasmussen (2011)
- ‣ Throttle valve control: Bischoff et al. (2013)
- ‣ Autonomous helicopter: e.g., Abbeel, Ng et al. (2003–2010), Bagnell & Schneider (2001)
- ‣ Pancake flipping (Kormushev et al., 2010)
- ‣ Throwing and catching balls (Kober et al., 2012)

# Summary



- RL is a principled framework for sequential decision making under uncertainty
- Value functions $V, Q$
- Exact RL: Dynamic programming
- Approximate RL: Monte Carlo, TD
- Policy Search with applications in robotics

m.deisenroth@imperial.ac.uk

**Thank you for your attention**

# Key References



- Sutton, Barto: *Reinforcement Learning: An Introduction* (online)
- Bertsekas: *Dynamic Programming and Optimal Control, Vol. 1–2*
- Szepesvári: *Algorithms for Reinforcement Learning* (online)
- Deisenroth et al.: *A Survey on Policy Search for Robotics* (online)

# RL Software Packages

- RLGlue: `http://glue.rl-community.org/`
- RLPy: `http://acl.mit.edu/RLPy/`
- CLSquare: `http://www.ni.uos.de/index.php?id=70`
- PIQLE: `http://piqle.sourceforge.net/`
- RL Toolbox: `http://www.igi.tugraz.at/ril-toolbox/`
- LibPG: `http://code.google.com/p/libpgrl/`
- PILCO (policy search): `http://mlg.eng.cam.ac.uk/pilco`

# References

[1] P. Abbeel, A. Coates, M. Quigley, and A. Y. Ng. An Application of Reinforcement Learning to Aerobatic Helicopter Flight. In B. Schölkopf, J. C. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19*, volume 19, page 2007. The MIT Press, Cambridge, MA, USA, 2007.

[2] D. P. Bertsekas. *Dynamic Programming and Optimal Control*, volume 1 of *Optimization and Computation Series*. Athena Scientific, Belmont, MA, USA, 3rd edition, 2005.

[3] D. P. Bertsekas. *Dynamic Programming and Optimal Control*, volume 2 of *Optimization and Computation Series*. Athena Scientific, Belmont, MA, USA, 3rd edition, 2007.

[4] B. Bischoff, D. Nguyen-Tuong, T. Koller, H. Markert, and A. Knoll. Learning Throttle Valve Control Using Policy Search. In *Proceedings of the European Conference on Machine Learning and Knowledge Discovery in Databases*, 2013.

[5] M. P. Deisenroth, G. Neumann, and J. Peters. *A Survey on Policy Search for Robotics*, volume 2 of *Foundations and Trends in Robotics*. NOW Publishers, 2013.

[6] M. P. Deisenroth and C. E. Rasmussen. PILCO: A Model-Based and Data-Efficient Approach to Policy Search. In *Proceedings of the International Conference on Machine Learning*, pages 465–472, New York, NY, USA, June 2011. ACM.

[7] S. Gelly, M. Schoenauer, M. Sebag, O. Teytaud, L. Kocsis, D. Silver, and C. Szepesvári. The Grand Challenge of Computer Go: Monte Carlo Tree Search and Extensions. *Communications of the ACM*, 55(3), 2012.

[8] Q. Gemine, E. Karangelos, D. Ernst, and B. Cornélusse. Active Network Management: Planning under Uncertainty for Exploiting Load Modulation. In *Proceedings of the 2013 IREP Symposium-Bulk Power System Dynamics and Control*, 2013.

[9] J. Kober, K. Mülling, and J. Peters. Learning Throwing and Catching Skills. In *Proceedings of the IEEE/RSJ International Conference on Robot Systems*, 2012.

[10] P. Kormushev, S. Calinon, and D. G. Caldwell. Robot Motor Skill Coordination with EM-based Reinforcement Learning. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2010.

[11] W. B. Powell. *Approximate Dynamic Programming: Solving the Curses of Dimensionality*, volume 703. John Wiley & Sons, 2007.

[12] M. Riedmiller. Neural Fitted $Q$ Iteration—First Experiences with a Data Efficient Neural Reinforcement Learning Method. In *Proceedings of the 16th European Conference on Machine Learning*, Porto, Portugal, 2005.

[13] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. Adaptive Computation and Machine Learning. The MIT Press, Cambridge, MA, USA, 1998.

[14] C. Szepesvári. *Algorithms for Reinforcement Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2010.

[15] G. J. Tesauro. Td-gammon, a self-teaching backgammon program, achieves master-level play. *Neural Computation*, 6(2):215–219, 1994.

# Appendix

# Self-Consistency of Value Functions

$$V^\pi(\boldsymbol{s}) = \mathbb{E}\big[R|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

# Self-Consistency of Value Functions

$$V^\pi(s) = \mathbb{E}\big[R|s_0 = s, \pi\big]$$

# Self-Consistency of Value Functions

$$V^\pi(\boldsymbol{s}) = \mathbb{E}\big[R|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$
$$= \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

# Self-Consistency of Value Functions

$$V^\pi(s) = \mathbb{E}\big[R | s_0 = s, \pi\big]$$
$$= \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1} | s_0 = s, \pi\big]$$

# Self-Consistency of Value Functions

$$V^\pi(s) = \mathbb{E}\big[R|s_0 = s, \pi\big]$$

$$= \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|s_0 = s, \pi\big]$$

$$= \mathbb{E}\big[r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2}|s_0 = s, \pi\big]$$

# Self-Consistency of Value Functions

$$V^\pi(\boldsymbol{s}) = \mathbb{E}\big[R|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

$$= \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

$$= \mathbb{E}\big[r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2}|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

# Self-Consistency of Value Functions

$$V^\pi(\boldsymbol{s}) = \mathbb{E}\big[R|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

$$= \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

$$= \mathbb{E}\big[r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2}|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

$$= \sum_a p_\pi(\boldsymbol{a}|\boldsymbol{s}) \sum_{\boldsymbol{s}'} p(\boldsymbol{s}'|\boldsymbol{s}, \boldsymbol{a}) \left(\mathbb{E}[r_1] + \gamma\mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+2}|\boldsymbol{s}_1 = \boldsymbol{s}', \pi\big]\right)$$

# Self-Consistency of Value Functions

$$V^\pi(\boldsymbol{s}) = \mathbb{E}\big[R|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

$$= \mathbb{E}\big[\sum_{k=0}^\infty \gamma^k r_{k+1}|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

$$= \mathbb{E}\big[r_1 + \gamma \sum_{k=0}^\infty \gamma^k r_{k+2}|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

$$= \sum_a p_\pi(\boldsymbol{a}|\boldsymbol{s}) \sum_{\boldsymbol{s}'} p(\boldsymbol{s}'|\boldsymbol{s}, \boldsymbol{a}) \left( \mathbb{E}[r_1] + \gamma \mathbb{E}\big[\sum_{k=0}^\infty \gamma^k r_{k+2}|\boldsymbol{s}_1 = \boldsymbol{s}', \pi\big] \right)$$

# Self-Consistency of Value Functions

$$V^\pi(s) = \mathbb{E}\big[R | s_0 = s, \pi\big]$$

$$= \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1} | s_0 = s, \pi\big]$$

$$= \mathbb{E}\big[r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2} | s_0 = s, \pi\big]$$

$$= \sum_a p_\pi(a|s) \sum_{s'} p(s'|s, a) \left( \mathbb{E}[r_1] + \gamma \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+2} | s_1 = s', \pi\big] \right)$$

$$= \sum_a p_\pi(a|s) \sum_{s'} p(s'|s, a) \big( \mathbb{E}[r_1(s, a, s')] + \gamma V^\pi(s') \big)$$

# Self-Consistency of Value Functions

$$V^{\pi}(\boldsymbol{s}) = \mathbb{E}\big[R|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

$$= \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

$$= \mathbb{E}\big[r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2}|\boldsymbol{s}_0 = \boldsymbol{s}, \pi\big]$$

$$= \sum_a p_{\pi}(\boldsymbol{a}|\boldsymbol{s}) \sum_{\boldsymbol{s}'} p(\boldsymbol{s}'|\boldsymbol{s}, \boldsymbol{a}) \left( \mathbb{E}[r_1] + \gamma \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+2}|\boldsymbol{s}_1 = \boldsymbol{s}', \pi\big] \right)$$

$$= \sum_a p_{\pi}(\boldsymbol{a}|\boldsymbol{s}) \sum_{\boldsymbol{s}'} p(\boldsymbol{s}'|\boldsymbol{s}, \boldsymbol{a}) \left( \mathbb{E}[r_1(\boldsymbol{s}, \boldsymbol{a}, \boldsymbol{s}')] + \gamma V^{\pi}(\boldsymbol{s}') \right)$$

# Self-Consistency of Value Functions

$$
\begin{aligned}
V^\pi(s) &= \mathbb{E}\big[R|s_0 = s, \pi\big] \\
&= \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|s_0 = s, \pi\big] \\
&= \mathbb{E}\big[r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2}|s_0 = s, \pi\big] \\
&= \sum_{a} p_\pi(a|s) \sum_{s'} p(s'|s, a) \left( \mathbb{E}[r_1] + \gamma \mathbb{E}\big[\sum_{k=0}^{\infty} \gamma^k r_{k+2}|s_1 = s', \pi\big] \right) \\
&= \sum_{a} p_\pi(a|s) \sum_{s'} p(s'|s, a) \left( \mathbb{E}[r_1(s, a, s')] + \gamma V^\pi(s') \right) \\
&= \mathbb{E}[r_1 + \gamma V^\pi(s_1)|s_0 = s, \pi]
\end{aligned}
$$

# Demo: Policy Iteration

Example: Shortest-path problem

# Demo: Policy Iteration

Example: Shortest-path problem

# Demo: Policy Iteration
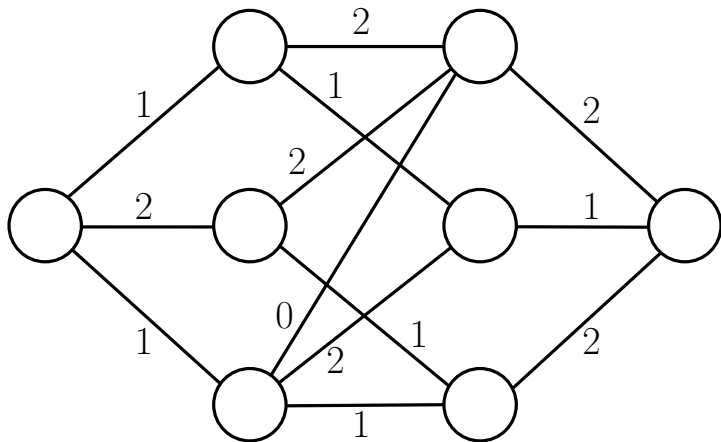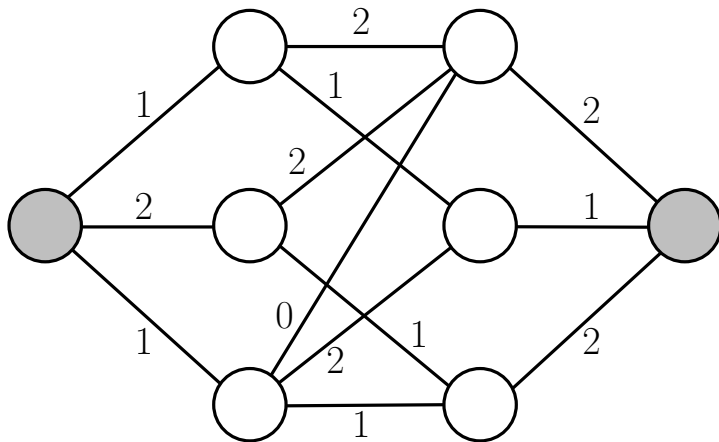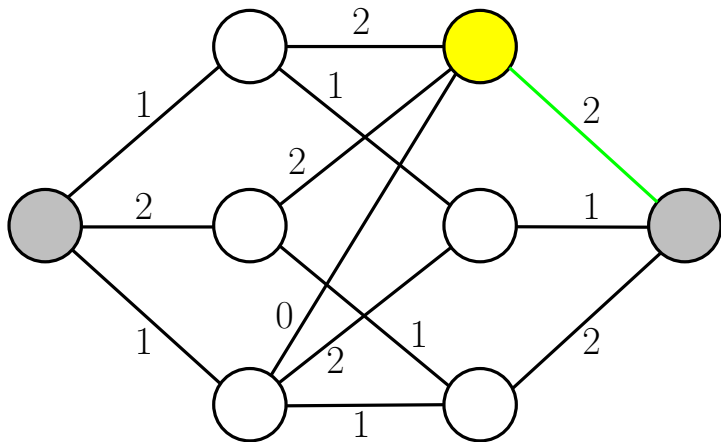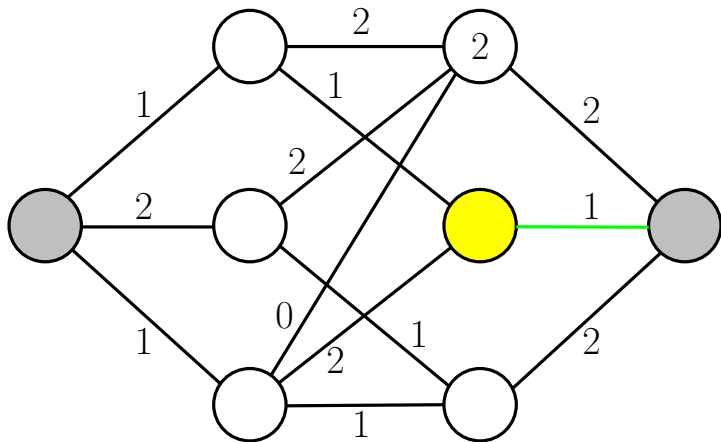
Example: Shortest-path problem

# Demo: Policy Iteration

Example: Shortest-path problem

# Demo: Policy Iteration

Example: Shortest-path problem

# Demo: Policy Iteration

Example: Shortest-path problem
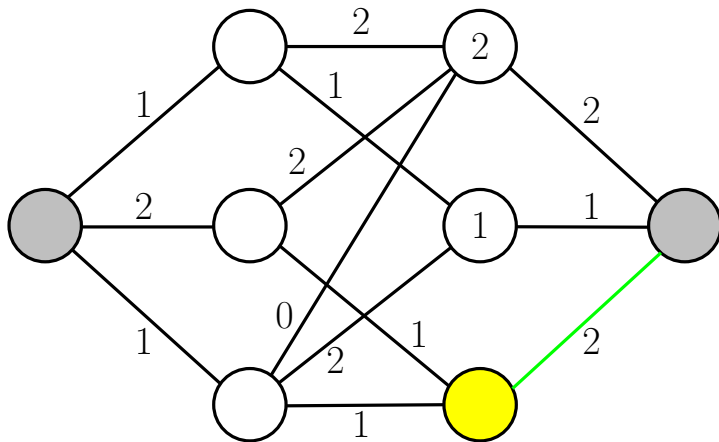
# Demo: Policy Iteration

Example: Shortest-path problem

# Demo: Policy Iteration

Example: Shortest-path problem

# Demo: Policy Iteration

Example: Shortest-path problem

# Demo: Policy Iteration

Example: Shortest-path problem
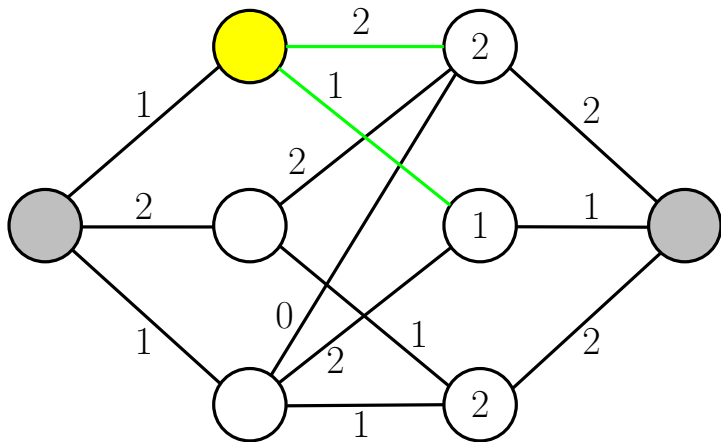
# Demo: Policy Iteration

Example: Shortest-path problem

# Demo: Policy Iteration

Example: Shortest-path problem

# Demo: Policy Iteration

Example: Shortest-path problem

# Demo: Policy Iteration

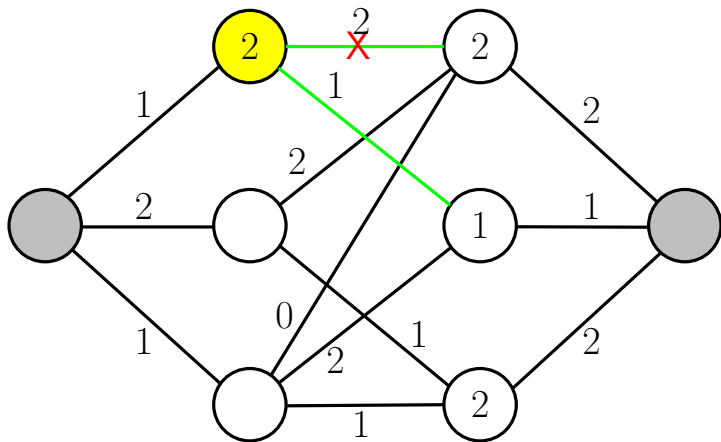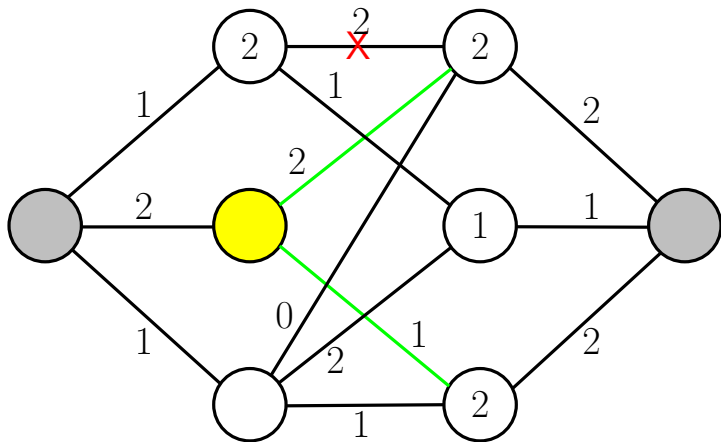Example: Shortest-path problem

# Demo: Value Iteration

Example: Shortest-path problem

# Demo: Value Iteration

Example: Shortest-path problem

# Demo: Value Iteration

Example: Shortest-path problem
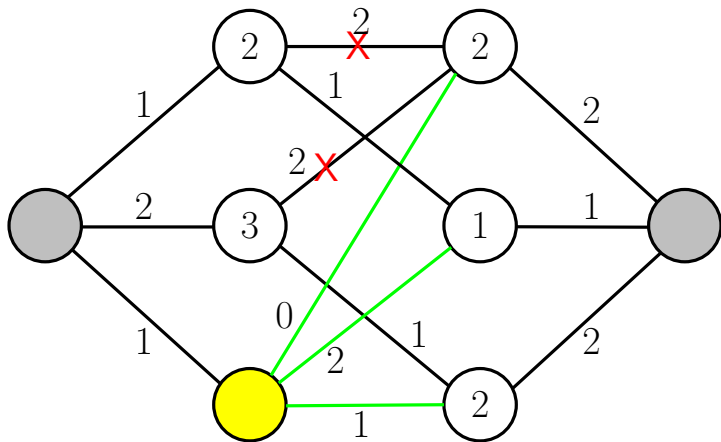
# Demo: Value Iteration

Example: Shortest-path problem

# Demo: Value Iteration

Example: Shortest-path problem

# Demo: Value Iteration

Example: Shortest-path problem

# Demo: Value Iteration

Example: Shortest-path problem

# Demo: Value Iteration

Example: Shortest-path problem

# Demo: Value Iteration

Example: Shortest-path problem

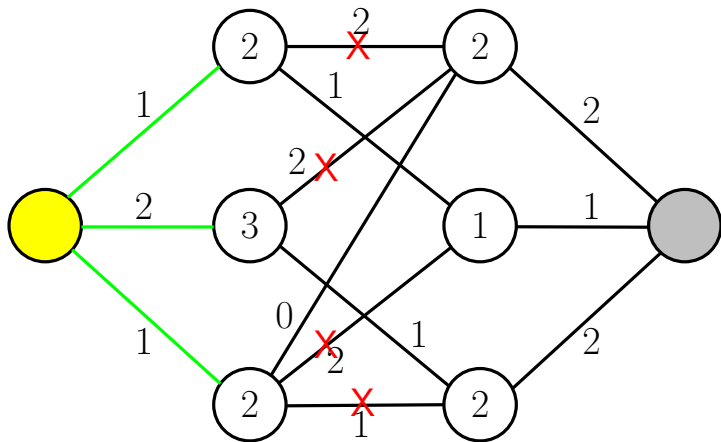# Demo: Value Iteration

Example: Shortest-path problem

# Demo: Value Iteration

Example: Shortest-path problem