

# REINFORCEMENT LEARNING APPROXIMATION AND DEEP LEARNING

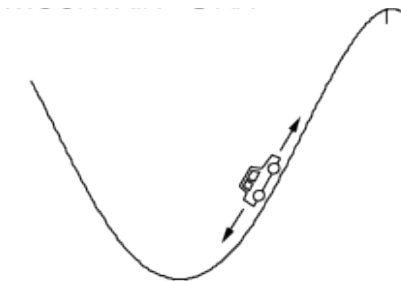
Odalric-Ambrym Maillard

2021-2022

Inria Scool

## CONTINUOUS SPACE AND FUNCTION APPROXIMATION

## EXAMPLE: MOUNTAIN-CAR



States: Position  $\times$  Speed of car (continuous).

Actions: Acceleration of car.

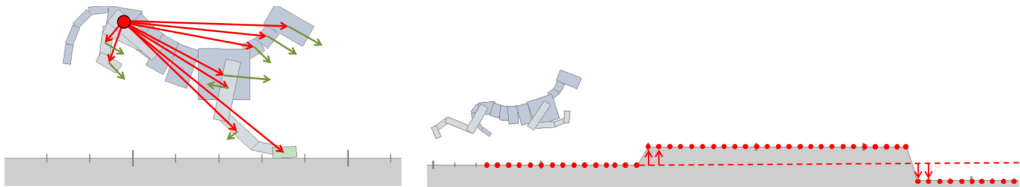
Transitions: Physics.

Rewards: 1 if reach top of mountain, 0 else.

Optimal policy: First go left to get momentum, then go right.

Continuous state-action space. For other classical tasks, check [https://gym.openai.com/envs/#classic\\_control](https://gym.openai.com/envs/#classic_control).

# EXAMPLE: MOTION-PLANNING



States: relative positions (red points) and velocities.

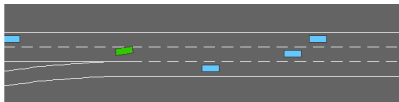
Actions: acceleration at each joint

Rewards: not falling (height of anchor point).

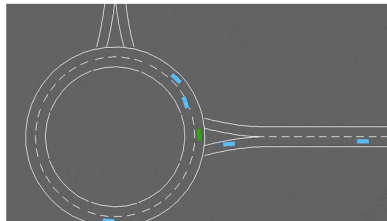
# EXAMPLE: HIGH-WAY

Controlling an ego-vehicle in various road environments

<https://github.com/eleurent/highway-env>



*The merge-v0 environment.*



*The roundabout-v0 environment.*

Actions: acceleration, changing lane, etc.

Rewards: no collision, speed, etc.

Study safety and robustness (here with partially known behavior of other vehicles).

<http://edouardleurent.com/publication/phd-thesis/>

DP requires:

- ▶ the state and action spaces to be small enough;
- ▶ the model to be known.

DP requires:

- ▶ the state and action spaces to be small enough;
- ▶ the model to be known.

Unfortunately:

- ▶ **the state space can be too large** (even continuous) for the value function to be represented exactly,

$$V_{\theta}(s) = \theta^{\top} \varphi(s) = \sum_{i=1}^d \theta_i \varphi_i(s) \quad (1)$$

DP requires:

- ▶ the state and action spaces to be small enough;
- ▶ the model to be known.

Unfortunately:

- ▶ **the state space can be too large** (even continuous) for the value function to be represented exactly,

$$V_{\theta}(s) = \theta^{\top} \varphi(s) = \sum_{i=1}^d \theta_i \varphi_i(s) \quad (1)$$

- ▶ **the model might be unknown** and one has to rely on a dataset

$$\mathcal{D} = \{(s_i, a_i, r_i, s'_i)_{1 \leq i \leq n}\}. \quad (2)$$

- ▶ the dataset can be obtained in multiple ways;



- ▶ Computing a **greedy policy** requires **knowing the model**:

$$\pi \in \mathcal{GV} \Leftrightarrow \forall s \in \mathcal{S}, \quad \pi(s) \in \operatorname{argmax}_{a \in \mathcal{A}} \left( m_a(s) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s') \right).$$

- ▶ Computing a **greedy policy** requires **knowing the model**:

$$\pi \in \mathcal{GV} \Leftrightarrow \forall s \in \mathcal{S}, \quad \pi(s) \in \operatorname{argmax}_{a \in \mathcal{A}} \left( m_a(s) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s') \right).$$

- ▶ **Sampling the optimality operator?**

- ▶ Optimality operator:

$$(\mathcal{T}V)(s) = \max_{a \in \mathcal{A}} \mathbb{E}_{s' \sim P(\cdot|s, a)} [m(s, a) + \gamma V(s')];$$

- ▶ with  $s'_{i,a} \sim P(\cdot|s_i, a)$ , a possible sampled operator:

$$(\hat{\mathcal{T}}V)(s_i) = \max_{a \in \mathcal{A}} (r(s_i, a) + \gamma V(s'_{i,a}));$$

- ▶ it is biased:  $\mathbb{E}[(\hat{\mathcal{T}}V)(s_i)|s_i] \neq (\mathcal{T}V)(s_i)$ .

- ▶ state-action value function (aka Q-function, quality function)

$$Q_{\pi}(s, a) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(S_t, A_t) | S_0 = s, A_0 = a, S_{t+1} \sim \mathbf{p}(\cdot | S_t, A_t), A_{t+1} = \pi(S_{t+1}) \right].$$

- ▶ state-action value function (aka Q-function, quality function)

$$Q_{\pi}(s, a) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(S_t, A_t) | S_0 = s, A_0 = a, S_{t+1} \sim \mathbf{p}(\cdot | S_t, A_t), A_{t+1} = \pi(S_{t+1}) \right].$$

- ▶ Bellman evaluation operator  $T_{\pi} : \mathbb{R}^{S \times A} \rightarrow \mathbb{R}^{S \times A}$

- ▶ definition:  $[T_{\pi} Q](s, a) = r(s, a) + \gamma \sum_{s' \in S} \mathbf{p}(s' | s, a) Q(s', \pi(s'))$ ;
- ▶  $Q_{\pi}$  is its unique fixed point:  $T_{\pi} Q_{\pi} = Q_{\pi}$ ;
- ▶ link to  $v_{\pi}$ :  $v_{\pi}(s) = Q_{\pi}(s, \pi(s))$ .

- ▶ state-action value function (aka Q-function, quality function)

$$Q_\pi(s, a) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(S_t, A_t) \mid S_0 = s, A_0 = a, S_{t+1} \sim \mathbf{p}(\cdot \mid S_t, A_t), A_{t+1} = \pi(S_{t+1}) \right].$$

- ▶ Bellman evaluation operator  $T_\pi : \mathbb{R}^{S \times \mathcal{A}} \rightarrow \mathbb{R}^{S \times \mathcal{A}}$

- ▶ definition:  $[T_\pi Q](s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbf{p}(s' \mid s, a) Q(s', \pi(s'))$ ;
- ▶  $Q_\pi$  is its unique fixed point:  $T_\pi Q_\pi = Q_\pi$ ;
- ▶ link to  $v_\pi$ :  $v_\pi(s) = Q_\pi(s, \pi(s))$ .

- ▶ Bellman optimality operator  $T_* : \mathbb{R}^{S \times \mathcal{A}} \rightarrow \mathbb{R}^{S \times \mathcal{A}}$

- ▶ definition:  $[T_* Q](s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbf{p}(s' \mid s, a) \max_{a' \in \mathcal{A}} Q(s', a')$ ;
- ▶  $Q_*$  is its unique fixed point:  $Q_* = T_* Q_*$ ;
- ▶ link to  $v_*$ :  $v_*(s) = \max_{a \in \mathcal{A}} Q_*(s, a)$ .

► Allows acting greedily:

► resp to  $V^\pi = Q_\pi(s, \pi(s))$ :

$$\begin{aligned}\pi' \in \mathcal{G}(V^\pi) &\Leftrightarrow \forall s \in \mathcal{S}, \quad \pi'(s) \in \operatorname{argmax}_{a \in \mathcal{A}} \left( r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbf{p}(s'|s, a) V^\pi(s') \right) \\ &\Leftrightarrow \forall s \in \mathcal{S}, \quad \pi'(s) \in \operatorname{argmax}_{a \in \mathcal{A}} Q_\pi(s, a).\end{aligned}$$

► resp. to  $V^*$ :

$$\star(s) \in \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a).$$

► resp. to any  $Q \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ :

$$\forall Q \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}, \quad \pi \in \mathcal{G}(Q) \Leftrightarrow \forall s \in \mathcal{S}, \quad \pi(s) \in \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a).$$

- ▶ Allows sampling easily the related operators

- ▶ recall the dataset

$$\mathcal{D} = \{(s_i, a_i, r_i, s'_i)_{1 \leq i \leq n}\}.$$

- ▶ sampled Bellman evaluation operator

$$[\hat{T}_\pi Q](s_i, a_i) = r_i + \gamma Q(s'_i, \pi(s'_i));$$

- ▶ sampled Bellman optimality operator

$$[\hat{T}_* Q](s_i, a_i) = r_i + \gamma \max_{a' \in \mathcal{A}} Q(s'_i, a').$$

⇒ Focus on this key function

- ▶ Allows sampling easily the related operators

- ▶ recall the dataset

$$\mathcal{D} = \{(s_i, a_i, r_i, s'_i)_{1 \leq i \leq n}\}.$$

- ▶ sampled Bellman evaluation operator

$$[\hat{T}_\pi Q](s_i, a_i) = r_i + \gamma Q(s'_i, \pi(s'_i));$$

- ▶ sampled Bellman optimality operator

$$[\hat{T}_* Q](s_i, a_i) = r_i + \gamma \max_{a' \in \mathcal{A}} Q(s'_i, a').$$

⇒ Focus on this key function

- ▶ Features for the Q-function:

- ▶ linear parametrization of the Q-function  $Q \in \mathcal{F}$ :

$$\mathcal{F} = \{Q_\theta(s, a) = \theta^\top \varphi(s, a), \theta \in \mathbb{R}^d\}$$



## CONTINUOUS SPACE AND FUNCTION APPROXIMATION

### **Approximate value iteration**

Approximate policy iteration

Policy representation

When  $\mathcal{S}, \mathcal{A}$  is huge, we may want to represent in Function spaces:  $m, p, V, Q$  or  $\pi$ .

1) Parameterize **Q function**:

$$\mathcal{F} = \{Q_{\theta}(s, a) = \theta^{\top} \varphi(s, a), \theta \in \mathbb{R}^d\}$$

When  $\mathcal{S}, \mathcal{A}$  is huge, we may want to represent in Function spaces:  $m, p, V, Q$  or  $\pi$ .

1) Parameterize **Q function**:

$$\mathcal{F} = \{Q_\theta(s, a) = \theta^\top \varphi(s, a), \theta \in \mathbb{R}^d\}$$

For VI, difficulty is to find  $Q_\theta \simeq \mathcal{T}Q_\theta$ : **Approximate VI**

For PI, difficulty is to find  $Q_\theta \simeq \mathcal{T}_\pi Q_\theta$ : **Approximate PI**

Different strategies:

- ▶ Minimize  $\|Q_\theta - \mathcal{T}_\pi Q_\theta\|$ : **Bellman residual minimization.**
- ▶ Solve  $Q_\theta = \Pi_{\mathcal{F}} \mathcal{T}_\pi Q_\theta$ : **Least-squares Temporal Differences.**

When  $\mathcal{S}, \mathcal{A}$  is huge, we may want to represent in Function spaces:  $m, p, V, Q$  or  $\pi$ .

1) Parameterize **Q function**

$$\mathcal{F} = \{Q_\theta(s, a) = \theta^\top \varphi(s, a), \theta \in \mathbb{R}^d\}$$

For VI, difficulty is to find  $Q_\theta \simeq \mathcal{T}Q_\theta$ : **Approximate VI**

For PI, difficulty is to find  $Q_\theta \simeq \mathcal{T}_\pi Q_\theta$ : **Approximate PI**

Different strategies:

- ▶ Minimize  $\|Q_\theta - \mathcal{T}_\pi Q_\theta\|$ : **Bellman residual minimization.**
- ▶ Solve  $Q_\theta = \Pi_{\mathcal{F}} \mathcal{T}_\pi Q_\theta$ : **Least-squares Temporal Differences.**

2) Parameterize  **$\pi$  function**

- ▶ Solve  $\operatorname{argmin}_{\pi \in \mathcal{F}} (\max_a Q(\cdot, a) - Q(\cdot, \pi))$ : **Direct Policy Search**

- ▶ Difficulty to apply value iteration ( $Q_{k+1} = \mathcal{T}Q_k$ ):
  - ▶ Max operator  $\mathcal{T}$  is **unknown**;
  - ▶  $Q_k \in \mathcal{F}$  may **not** imply  $\mathcal{T}Q_k \in \mathcal{F}$ .

- ▶ Difficulty to apply value iteration ( $Q_{k+1} = \mathcal{T}Q_k$ ):
  - ▶ Max operator  $\mathcal{T}$  is **unknown**;
  - ▶  $Q_k \in \mathcal{F}$  may **not** imply  $\mathcal{T}Q_k \in \mathcal{F}$ .
- ▶ Example of naive **Approximate value iteration**:
  - ▶ writing  $Q_k = Q_{\theta_k}$ , sampled operator:

$$(\hat{\mathcal{T}}Q_k)(s_i, a_i) = r_i + \gamma \max_{a' \in \mathcal{A}} Q_k(s'_i, a')$$

- ▶ search for  $Q \in \mathcal{F}$  being the closest to  $\hat{\mathcal{T}}Q_k$ :

$$Q_{k+1} \in \underset{Q_{\theta} \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \left( Q_{\theta}(s_i, a_i) - (\hat{\mathcal{T}}Q_k)(s_i, a_i) \right)^2.$$

- ▶ summary:  $Q_{k+1} = \Pi_{\mathcal{F}} \hat{\mathcal{T}}Q_k$ .

---

**Algorithm 1** Approximate value iteration

---

**Input:** A dataset  $\mathcal{D} = \{(s_i, a_i, r_i, s'_i)_{1 \leq i \leq n}\}$ , the number  $K$  of iterations, a function approximator, an initial state-action value function  $Q_0$ .

1: **for**  $k = 0$  **to**  $K$  **do**

2:   Apply the sampled Bellman operator to function  $Q_k$ :

$$(\hat{T}Q_k)(s_i, a_i) = r_i + \gamma \max_{a' \in \mathcal{A}} Q_k(s'_i, a').$$

3:   Solve the regression problem with inputs  $(s_i, a_i)$  and outputs  $(\hat{T}Q_k)(s_i, a_i)$  to get the Q-function  $Q_{k+1}$

4: **end for**

5: **return** The greedy policy  $\pi_{K+1} \in \mathcal{G}(Q_{K+1})$ :

$$\forall s \in \mathcal{S}, \quad \pi_{K+1} \in \operatorname{argmax}_{a \in \mathcal{A}} Q_{K+1}(s, a).$$

- ▶ Approximate value iteration:

$$Q_{k+1} = \mathbb{A} \hat{T} Q_k.$$

- ▶  $\mathbb{A}$  is an abstract approximation operator



- ▶ Approximate value iteration:

$$Q_{k+1} = \mathbb{A} \hat{T} Q_k.$$

- ▶  $\mathbb{A}$  is an abstract approximation operator
- ▶ Problem:  $\mathbb{A} \hat{T}$  should be a **contraction**!

- ▶ Approximate value iteration:

$$Q_{k+1} = \mathbb{A} \hat{T} Q_k.$$

- ▶  $\mathbb{A}$  is an abstract approximation operator
- ▶ Problem:  $\mathbb{A} \hat{T}$  should be a **contraction**!
  - ▶ Otherwise divergence can (will) occur;
  - ▶ Not the case for  $\Pi \hat{T} \dots$
  - ▶ Works if “**averagers**” are used for function approximation, such as
    - ▶ ensemble of trees
    - ▶ notably extremely randomized forests: **fitted-Q**
    - ▶ kernel averagers (Nadaraya-Watson)

Working with AVI is tricky due to **max Bellman operator**  $\mathcal{T}$ .

## CONTINUOUS SPACE AND FUNCTION APPROXIMATION

Approximate value iteration

### **Approximate policy iteration**

Policy representation

▷ Idea: Instead of considering  $\mathcal{T}$ , consider  $\mathcal{T}_\pi$  for a fixed  $\pi$ .

▶ Policy iteration:

- ① policy evaluation: solve the fixed-point equation  $Q_{\pi_k} = T_{\pi_k} Q_{\pi_k}$ ;
- ② policy improvement: compute the greedy policy  $\pi_{k+1} = \mathcal{G}(Q_{\pi_k})$ .

▶ **Approximate policy iteration:**

- ① approximate policy evaluation: find a function  $Q_k \in \mathcal{F}$  such that  $Q_k \approx T_{\pi_k} Q_k$ ;
- ② policy improvement: compute the greedy policy  $\pi_{k+1} = \mathcal{G}(Q_{\pi_k})$ .

---

**Algorithm 2** Approximate policy iteration (generic form)

---

**Input:** An initial  $\pi_0 \in \mathcal{A}^S$  (possibly an initial  $Q_0$  and  $\pi_0 \in \mathcal{G}(Q_0)$ ), number of iterations  $K$

1: **for**  $k = 0$  **to**  $K$  **do**

2:   approximate policy evaluation: find  $Q_k \in \mathcal{F}$  such that  $Q_k \approx T_{\pi_k} Q_k$ .

3:   policy improvement:  $\pi_{k+1} \in \mathcal{G}(Q_k)$ .

4: **end for**

5: **return** the policy  $\pi_{K+1}$

---

▷ How to find an **approximate fixed point** of  $T_\pi$ , that is a function  $Q_\theta \in \mathcal{F}$  such that  $Q_\theta \approx T_\pi Q_\theta$ ?

- ▷ Assume that  $Q_\pi$  is known: simply a **regression** problem. E.g. linear least-squares:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (Q_\pi(s_i, a_i) - Q_\theta(s_i, a_i))^2.$$

- ▷ Assume that  $Q_\pi$  is known: simply a **regression** problem. E.g. linear least-squares:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (Q_\pi(s_i, a_i) - Q_\theta(s_i, a_i))^2.$$

- ▶  $Q_\pi$  is (obviously) unknown...
- ▶ Idea 1: **Monte carlo rollout?**

- ▶ sample a full trajectory starting in  $s_i$  where action  $a_i$  is chosen first, all subsequent states being sampled according to the system dynamics and actions chosen according to  $\pi$ ;
- ▶ let  $q_i$  be the associated discounted cumulative reward;
- ▶ unbiased estimate:  $\mathbb{E}[q_i | s_i, a_i] = Q_\pi(s_i, a_i)$

- ▷ Assume that  $Q_\pi$  is known: simply a **regression** problem. E.g. linear least-squares:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (Q_\pi(s_i, a_i) - Q_\theta(s_i, a_i))^2.$$

- ▶  $Q_\pi$  is (obviously) unknown...
- ▶ Idea 1: **Monte carlo rollout?**

- ▶ sample a full trajectory starting in  $s_i$  where action  $a_i$  is chosen first, all subsequent states being sampled according to the system dynamics and actions chosen according to  $\pi$ ;
- ▶ let  $q_i$  be the associated discounted cumulative reward;
- ▶ unbiased estimate:  $\mathbb{E}[q_i | s_i, a_i] = Q_\pi(s_i, a_i)$
- ▶ replace  $Q_\pi(s_i, a_i)$  by the unbiased estimate  $q_i$ :

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (q_i - Q_\theta(s_i, a_i))^2.$$

- ▷ Drawbacks: requires a simulator, rollouts can be quite noisy.



- ▶ Idea 2: **minimize the residual**  $\|Q_\theta - T_\pi Q_\theta\|$  for some norm.
- ▶ With an  $\ell_2$ -loss, parametric representation and sampled operator:

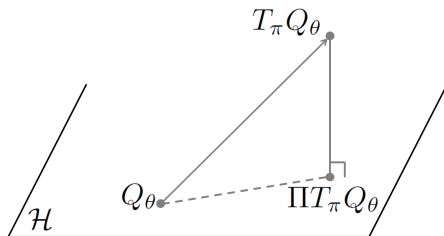
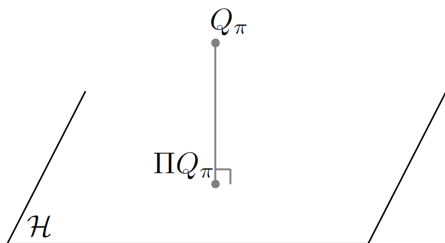
$$\begin{aligned} & \min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \left( (\hat{T}_\pi Q_\theta)(s_i, a_i) - Q_\theta(s_i, a_i) \right)^2 \\ &= \min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \left( r_i + \gamma Q_\theta(s'_i, \pi(s'_i)) - Q_\theta(s_i, a_i) \right)^2. \end{aligned}$$

- ▶ Idea 2: **minimize the residual**  $\|Q_\theta - T_\pi Q_\theta\|$  for some norm.
- ▶ With an  $\ell_2$ -loss, parametric representation and sampled operator:

$$\begin{aligned} & \min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \left( (\hat{T}_\pi Q_\theta)(s_i, a_i) - Q_\theta(s_i, a_i) \right)^2 \\ &= \min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \left( r_i + \gamma Q_\theta(s'_i, \pi(s'_i)) - Q_\theta(s_i, a_i) \right)^2. \end{aligned}$$

- ▶ However, there is a **bias problem**:

$$\begin{aligned} & \mathbb{E}[(\hat{T}_\pi Q_\theta)(s_i, a_i) - Q_\theta(s_i, a_i)]^2 | s_i, a_i] \\ &= ((\mathcal{T}_\pi Q_\theta)(s_i, a_i) - Q_\theta(s_i, a_i))^2 + \mathbb{V}((\hat{T}_\pi Q_\theta)(s_i, a_i) | s_i, a_i) \\ &\neq ((\mathcal{T}_\pi Q_\theta)(s_i, a_i) - Q_\theta(s_i, a_i))^2. \end{aligned}$$



► Idea 3: Solve  $Q_\theta = \Pi T_\pi Q_\theta$ . as a **nested optimization** problem:

$$w_\theta = \operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (r_i + \gamma Q_\theta(s'_i, \pi(s'_i)) - Q_w(s_i, a_i))^2 \quad (3)$$

$$\theta_n = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (Q_\theta(s_i, a_i) - Q_{w_\theta}(s_i, a_i))^2. \quad (4)$$

- (3) is ordinary least-squares in  $w$ . Solution:

$$w_{\theta} = \left( \sum_{i=1}^n \varphi(s_i, a_i) \varphi(s_i, a_i)^{\top} \right)^{-1} \sum_{i=1}^n \varphi(s_i, a_i) (r_i + \gamma \theta^{\top} \varphi(s'_i, \pi(s'_i))).$$

- (3) is ordinary least-squares in  $w$ . Solution:

$$w_\theta = \left( \sum_{i=1}^n \varphi(s_i, a_i) \varphi(s_i, a_i)^\top \right)^{-1} \sum_{i=1}^n \varphi(s_i, a_i) (r_i + \gamma \theta^\top \varphi(s'_i, \pi(s'_i))).$$

- (4) is minimized for  $\theta = w_\theta$ :

$$\begin{aligned} \theta_n = w_{\theta_n} &\Leftrightarrow \theta_n = \left( \sum_{i=1}^n \varphi(s_i, a_i) \varphi(s_i, a_i)^\top \right)^{-1} \sum_{i=1}^n \varphi(s_i, a_i) (r_i + \gamma \theta_n^\top \varphi(s'_i, \pi(s'_i))) \\ &\Leftrightarrow \theta_n = \left( \sum_{i=1}^n \varphi(s_i, a_i) (\varphi(s_i, a_i) - \gamma \varphi(s'_i, \pi(s'_i)))^\top \right)^{-1} \sum_{i=1}^n \varphi(s_i, a_i) r_i. \end{aligned}$$

- This API strategy with LSTD is called **LSPi (Least-Squares Policy Iteration)**.

---

## Algorithm 3 Least-squares policy iteration

---

**Input:** An initial  $\pi_0 \in \mathcal{A}^S$  (possibly an initial  $Q_0$  and  $\pi_0 \in \mathcal{G}(Q_0)$ ), number of iterations  $K$

1: **for**  $k = 0$  **to**  $K$  **do**

2:   approximate policy evaluation:

$$\theta_k = \left( \sum_{i=1}^n \varphi(s_i, a_i) (\varphi(s_i, a_i) - \gamma \varphi(s'_i, \pi_k(s'_i)))^\top \right)^{-1} \sum_{i=1}^n \varphi(s_i, a_i) r_i.$$

3:   policy improvement:

$$\pi_{k+1} \in \mathcal{G}(Q_{\theta_k}).$$

4: **end for**

5: **return** the policy  $\pi_{K+1}$

---

## Proposition (Williams and Baird, 93)

$$\|V_\pi - \hat{V}_{BR}\|_\infty \leq \frac{1+\gamma}{1-\gamma} \|V_\pi - \hat{V}_{\text{best}}\|_\infty$$

## Proposition (Tsitsiklis and Van Roy, 97)

Let  $\rho$  be the stationary distribution of  $P_\pi$ , then

$$\|V_\pi - \hat{V}_{LSTD}\|_{2,\rho} \leq \frac{1}{1-\gamma} \|V_\pi - \hat{V}_{\text{best}}\|_{2,\rho}$$

Approximate PI:  $\pi_k = \mathcal{G}[v_{k-1}]$  and  $v_k = v_{\pi_k} + \varepsilon_k$

## Proposition

Assume  $\|\varepsilon_k\|_\infty \leq \varepsilon$ . The loss due to running  $\pi_k$  instead of  $\pi_\star$  satisfies

$$\limsup_{k \rightarrow \infty} \|V_\star - v_{\pi_k}\|_\infty \leq \frac{2\gamma}{(1-\gamma)^2} \varepsilon.$$

## CONTINUOUS SPACE AND FUNCTION APPROXIMATION

Approximate value iteration

Approximate policy iteration

---

**Policy representation**

---



- ▶ Idea: represent the **policy** instead of **Q-function**

- ▶ Idea: represent the **policy** instead of **Q-function**
- ▶ Motivation: might be easier to learn (we control what it is)
- ▶ Let  $\mathcal{F} \subset \mathcal{A}^S$  be an hypothesis space of policies. At iteration  $k$ , assume  $Q_{\pi_k}(s_i, a)$  are known, and solve the **classification problem**

$$\pi_{k+1} \in \operatorname{argmin}_{\pi \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \left( \max_{a \in \mathcal{A}} Q_{\pi_k}(s_i, a) - Q_{\pi_k}(s_i, \pi(s_i)) \right).$$

- ▶ Practically, replace  $Q_{\pi_k}(s_i, a)$  by a Monte Carlo rollout.
- ▶ Often called **DPI** for **Direct Policy Iteration**.

*“The more applied you go, the stronger theory you need”*

# MERCI

[odalricambrym.maillard@inria.fr](mailto:odalricambrym.maillard@inria.fr)

[odalricambrymmaillard.wordpress.com](https://odalricambrymmaillard.wordpress.com)