

Lecture 20 - March 31, 2023

↓ - Function Approximation in Reinforcement Learning

- Basics of Function Approximations
- Least square Policy Iteration (LSPI)
- Neural Fitted Q-Iterations (NFQI)
- Deep Q-Network (DQN) → Finite action
- Deep Policy Gradients (DPG) → Large/continuous Action

Exam 2 → Tues, April 4

Project 3 → Due April 14

TA's office hour:

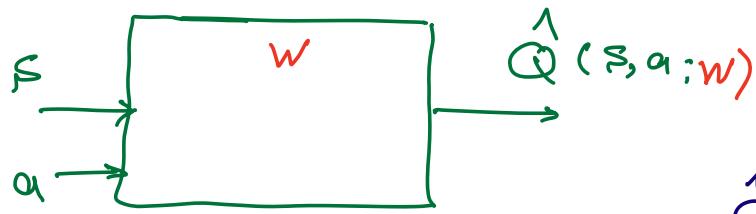
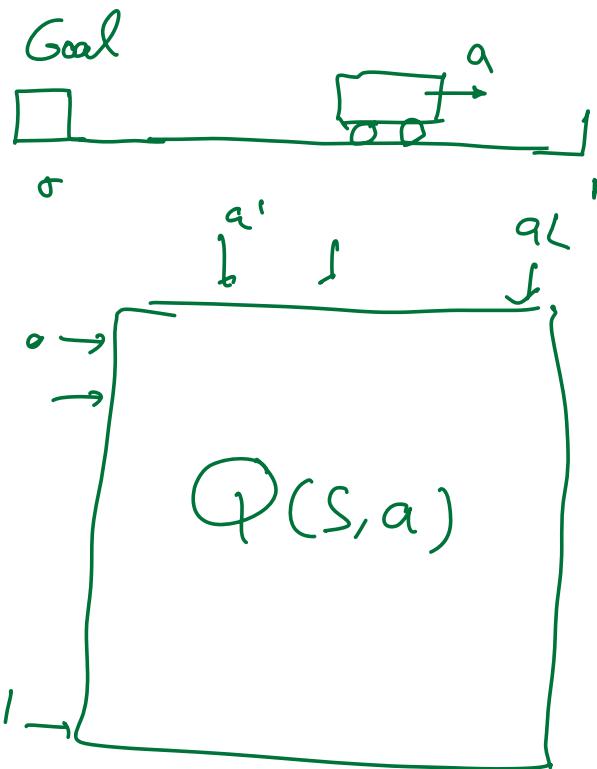
Wednesdays, 2pm-3pm (in-person)

Fridays, 2pm-3pm (virtual)

$$S \in [0, 1]$$

$v_{\min} \leq v \leq v_{\max}$

$$a \in \{-1, 0, +1\}$$



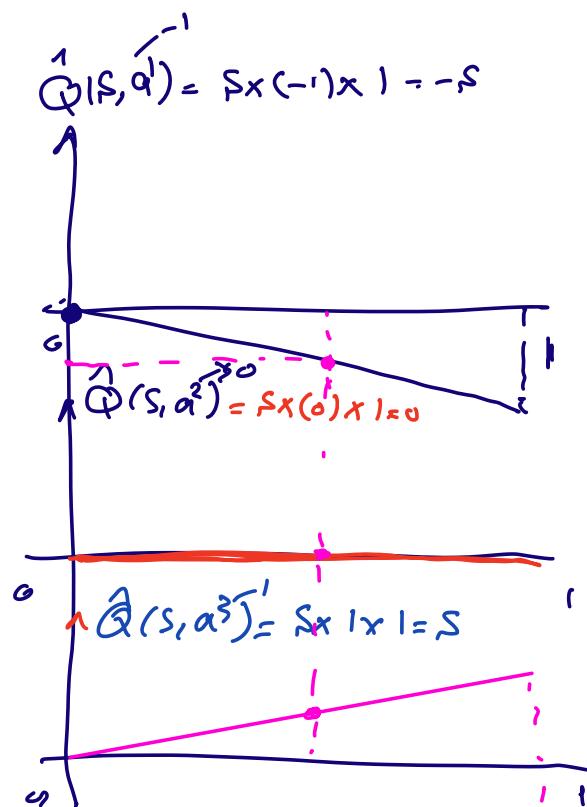
$$\hat{Q}(s, a; w) = s \cdot a \cdot w$$

$w \neq 1$

$$s=0.5 \rightarrow a^* = \arg \max_{a \in A} \hat{Q}(s, a)$$

$\downarrow$

$$= \arg \max \{-0.5, 0, 0.5\} = a^* = 1$$



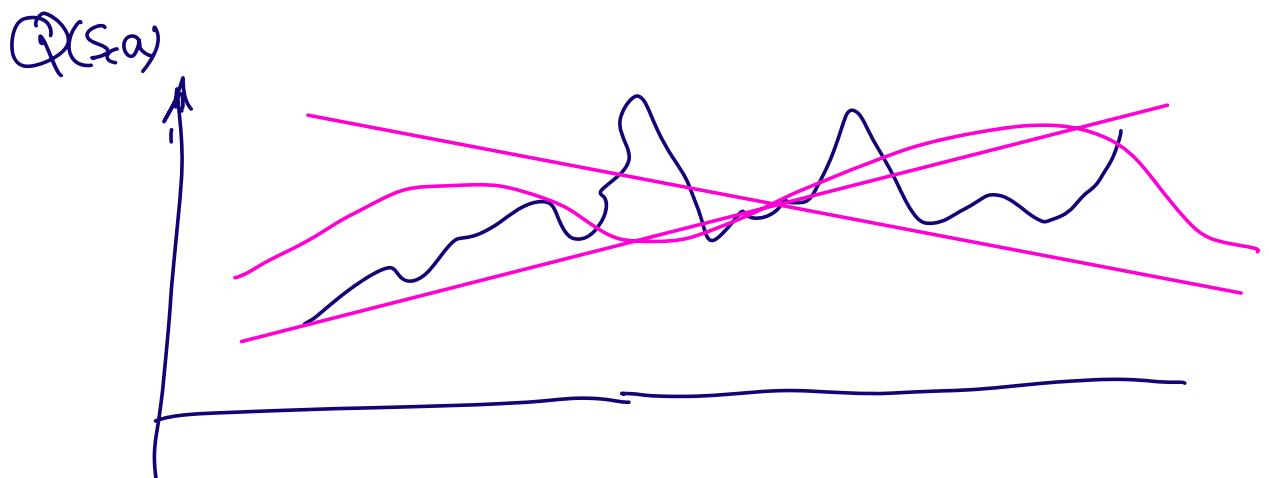
$$\hat{Q}_W(s, a; w) = \phi^T(s, a) w \quad : \text{Linear Function Approx}$$

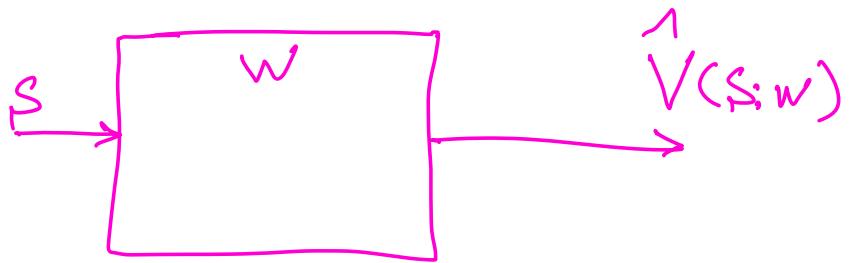
↳ Feature or Basis Function

$$\phi(s, a) = \begin{bmatrix} 1 \\ s \\ sa \\ s^2 a \\ sa + a^2 \end{bmatrix}$$

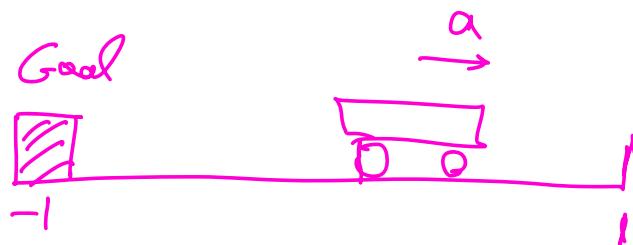
$$\hat{Q}_W(s, a) = [1 \ s \ sa \ s^2 a \ sa + a^2] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \end{bmatrix}$$

$$= w_1 + s w_2 + sa w_3 + s^2 a w_4 + (sa + a^2) w_5$$





$$\hat{V}(s; w) = \phi^T(s) w$$



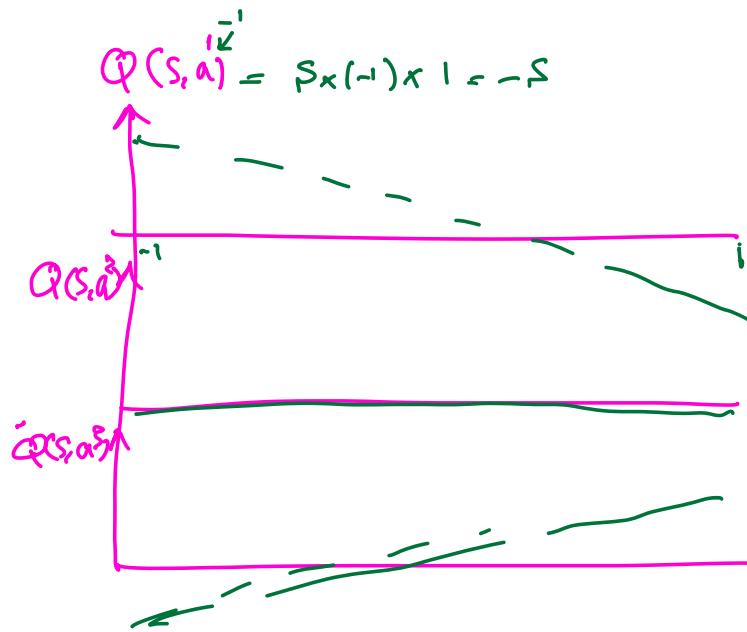
$$\pi^*(s) = -1$$

$$a = \{-1, 0, 1\}$$

$$\hat{Q}(s, a) = s a w$$

$$\hat{Q}(s, a) = |s| a w$$

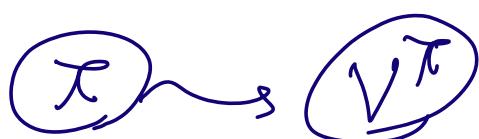
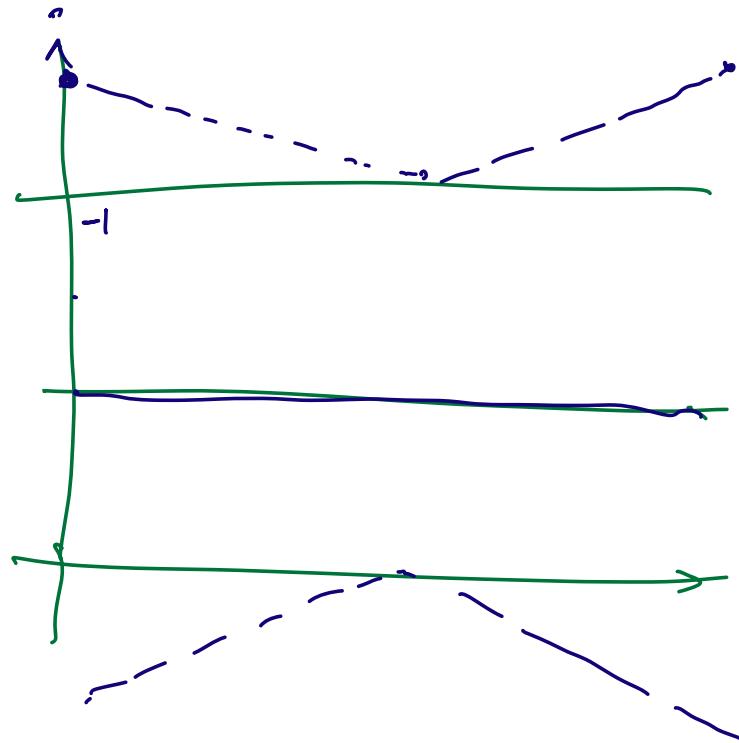
$$\phi(s, a) = |s| a$$



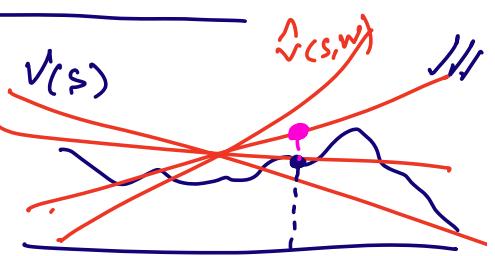
$$\hat{Q}(s, a, w) = |s| a w$$

$$w = -1$$

$$\hat{Q}(s, a, w) = -|s| a$$



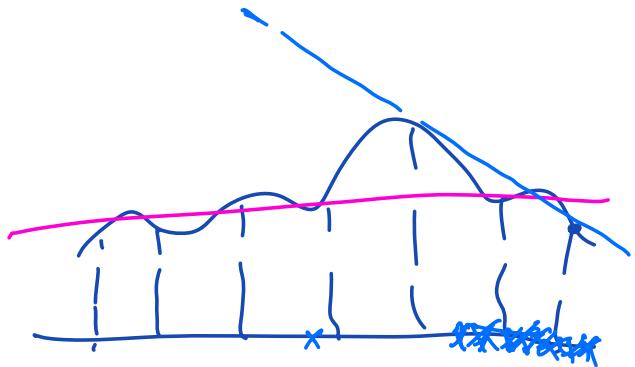
$$V_w(s) = \phi^T(s) w$$



$$J(w) = E_{\pi} \left[ (V_{(s)}^{\pi} - \hat{V}_{(s; w)})^2 \right]$$

$$\Rightarrow \text{Learn } w^* = \arg \min_{w \in W} J(w)$$

① ← Degree of freedom



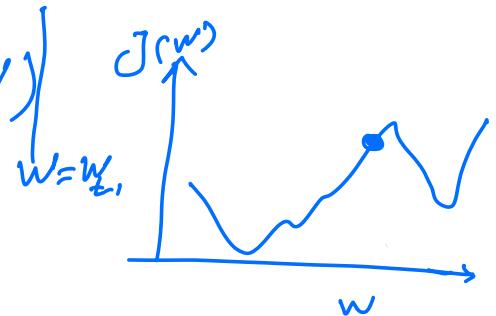
$$J(w) = \underset{\pi}{E} \left[ (\hat{V}_{(s)}^{\pi} - \hat{V}_{(s;w)})^2 \right]$$

distribution of states  
under policy  $\pi$

$$= \sum_{s \in S} \mu_{(s)} (\hat{V}_{(s)}^{\pi} - \hat{V}_{(s;w)})^2$$

Gradient descent  $\Rightarrow$  learn  $w^* = \underset{w \in W}{\operatorname{argmin}} J(w)$

$$W_t \leftarrow W_{t-1} - \frac{1}{2} \alpha \nabla_w J(w)$$



$$W_t \leftarrow W_{t-1} + \frac{1}{2} \alpha 2 (\hat{V}_{(s)}^{\pi} - \hat{V}_{(s;w)}) \nabla_w \hat{V}_{(s;w)}$$

9	10	11	12
8			
7			
6	5		
4	3	2	1

■ Wall    
 ■ Bump    
 ■ Goal

$$V^\pi(s) \Rightarrow \hat{V}(s; w)$$

$$\hat{V}(s; w) = \underbrace{\begin{bmatrix} 1_{s=s^1} & 1_{s=s^2} & \dots & 1_{s=s^9} \end{bmatrix}}_{\Phi^T(s)} \begin{bmatrix} w_1 \\ \vdots \\ w_9 \end{bmatrix}$$

$$s \in S \rightarrow \hat{V}(s; w) = \underbrace{\Phi^T(s)}_{w_1} = w_1$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & \dots \end{bmatrix}^T \begin{bmatrix} w_1 \\ \vdots \\ w_9 \end{bmatrix}$$

$$s \in S \rightarrow \hat{V}(s; w) = w_5$$

Monte-Carlo  $\rightarrow$  Large/Continuous space

$$\pi \rightarrow \sqrt{\lambda} V(s; w)$$

$$s_0, G(s_0) \xrightarrow{s_0 \xrightarrow{\pi} s_1, r_1} \sim V(s_0) = E[G_{s_0} | \dots]$$

$$s_1, G(s_1) \xrightarrow{s_1 \xrightarrow{\pi} s_2, r_2} G(s_0) = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots + \gamma^{T-1} R_T$$

$$G(s_1) = r_2 + \gamma r_3 + \gamma^2 r_4 + \dots + \gamma^{T-2} R_T$$

$$s_2 \xrightarrow{\pi} s_3, r_3$$

:

$$s_{T-1}, G(s_{T-1}) \xrightarrow[s_{T-1} \xrightarrow{\pi} s_T, r_T]{\text{---}} G(s_T) = R_T$$

$$J(w) = E_{\pi}[(V(s) - \hat{V}(s; w))^2]$$

$$= \frac{1}{|D|} \sum_{s \in D} (G_t - \hat{V}(s; w))^2$$

$$w_t = w_{t-1} - \frac{1}{2} \alpha \nabla J(w)|_{w_{t-1}}$$

$$w_t = w_{t-1} + \frac{1}{2} \times 2\alpha (G_t - \hat{V}(s; w)) \nabla \hat{V}(s; w)|_{w_{t-1}}$$

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

Input: a policy  $\pi$  to be evaluated

Initialize:

$V(s) \in \mathbb{R}$ , arbitrarily, for all  $s \in \mathcal{S}$

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

Loop forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

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

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

Append  $G$  to  $Returns(S_t)$

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

### Gradient Monte Carlo Algorithm for Estimating $\hat{v} \approx v_\pi$

Input: the policy  $\pi$  to be evaluated

Input: a differentiable function  $\hat{v} : \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameter: step size  $\alpha > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Loop forever (for each episode):

Generate an episode  $S_0, A_0, R_1, S_1, A_1, \dots, R_T, S_T$  using  $\pi$

Loop for each step of episode,  $t = 0, 1, \dots, T-1$ :

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [G_t - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$

$$\pi \rightarrow V^\pi \implies \hat{V}(s; w) \quad \text{TD-version}$$

$s_0, \pi(s_0), s_1, r_1$

$$V(s) = V(s) + \alpha [R + \gamma V(s') - V(s)]$$

$$J(w) = E_\pi \left[ \underbrace{(V^\pi(s) - \hat{V}(s; w))^2}_{R + \gamma \hat{V}(s'; w)} \right]$$

$$w_t = w_{t-1} - \frac{1}{2} \alpha \nabla J(w)|_{w_{t-1}}$$

$$= w_{t-1} + \frac{1}{2} \alpha \left( R + \gamma \hat{V}(s_t; w) - \hat{V}(s_{t-1}; w) \right) \nabla \hat{V}(s_t; w)$$

$$\nabla \hat{V}(s; w) = \nabla (\phi^\top(s) w) = \phi(s)$$

### Semi-gradient TD(0) for estimating $\hat{v} \approx v_\pi$

Input: the policy  $\pi$  to be evaluated

Input: a differentiable function  $\hat{v} : \mathcal{S}^+ \times \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $\hat{v}(\text{terminal}, \cdot) = 0$

Algorithm parameter: step size  $\alpha > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

        Choose  $A \sim \pi(\cdot | S)$

        Take action  $A$ , observe  $R, S'$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})] \nabla \hat{v}(S, \mathbf{w})$

$S \leftarrow S'$

    until  $S$  is terminal

Fixed Policy  $\pi \rightarrow \mathcal{V}^\pi := \hat{V}(s; \mathbf{w})$

$$Q(s, a) = Q(s, a) + \alpha [R + \gamma Q(s', a) - Q(s, a)]$$

$s_0 \xrightarrow{a \sim \pi} s_1, R \rightarrow a' \sim \pi'$

$$J(w) = E_{\pi} \left[ \underbrace{(Q(s, a) - \hat{Q}(s, a; w))^2}_{R + \gamma \frac{\hat{Q}(s', a'; w)}{\hat{\phi}^T(s', a')w}} \right]$$

$$\begin{aligned} w_t &= w_{t-1} - \frac{1}{2} \alpha J(w) \\ &= w_{t-1} + \frac{1}{2} \alpha (R + \gamma \hat{Q}(s', a'; w) - \hat{Q}(s, a; w)) \nabla \hat{Q}(s, a; w) \end{aligned}$$

### Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_*$

Input: a differentiable function  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Initialize value-function weights  $w \in \mathbb{R}^d$  arbitrarily (e.g.,  $w = 0$ )

Repeat (for each episode):

$S, A \leftarrow$  initial state and action of episode (e.g.,  $\varepsilon$ -greedy)

Repeat (for each step of episode):

Take action  $A$ , observe  $R, S'$   $\rightarrow \pi(s) = \begin{cases} \arg \max Q(s, a) & 1-\varepsilon \\ \arg \max_{a \in \mathcal{A}} \phi(s, a) w & \varepsilon \\ \text{Random} & \varepsilon \end{cases}$   
If  $S'$  is terminal:

$$w \leftarrow w + \alpha [R - \hat{q}(S, A, w)] \nabla \hat{q}(S, A, w)$$

Go to next episode

Choose  $A'$  as a function of  $\hat{q}(S', \cdot, w)$  (e.g.,  $\varepsilon$ -greedy)

$$w \leftarrow w + \alpha [R + \gamma \hat{q}(S', A', w) - \hat{q}(S, A, w)] \nabla \hat{q}(S, A, w)$$

$$S \leftarrow S'$$

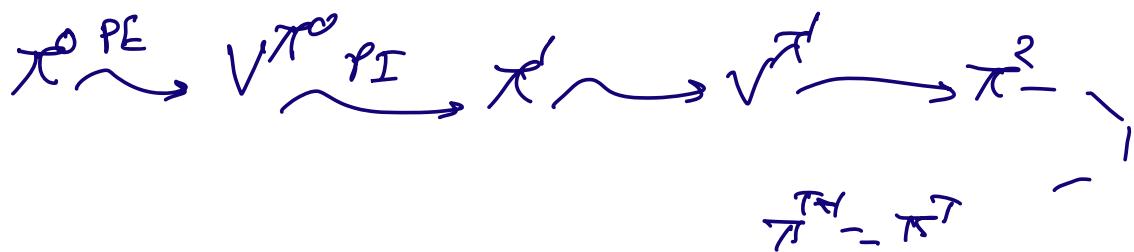
$$A \leftarrow A'$$

## Batch Learning

$$D = \{(s_0, a_0, r_0, s_1), \dots\}$$

$\Downarrow$   
 $\pi^*$

LSPIT



Policy Evaluation:

$$V_{k+1}(s) = \sum_{s'} P(s'|s, \pi(s)) [R + \gamma V_k(s')]$$

Policy Improve-

$$\pi'(s) = \underset{a \in A}{\operatorname{argmax}} \sum_{s'} P(s'|s, a) [R + \gamma V_{\pi}(s')]$$

$$Q^{\pi}(s, a) = \sum_{s'} P(s'|s, a) [R + \gamma V_{\pi}(s')]$$

$$= \sum_{s'} P(s'|s, a) [R + \gamma Q_{\pi}(s', \pi(s'))]$$

$$Q^{\pi} = R + \gamma M Q^{\pi}$$

↓

$$\begin{bmatrix} Q^T(s_1, a_1) & \dots & Q^T(s_1, a^{(A)}) \\ Q^T(s_1, a_1) & \dots & Q^T(s_1, a^{(A)}) \end{bmatrix}$$

LSPI