

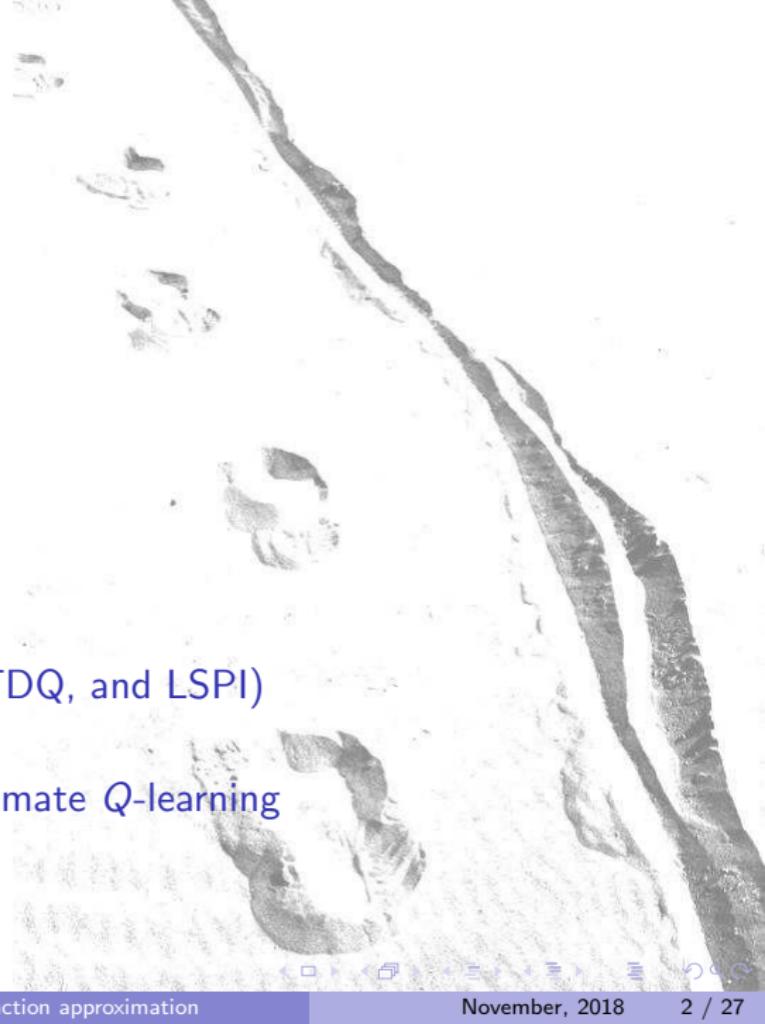
# RL for continuous spaces, function approximation

Sébastien Gros

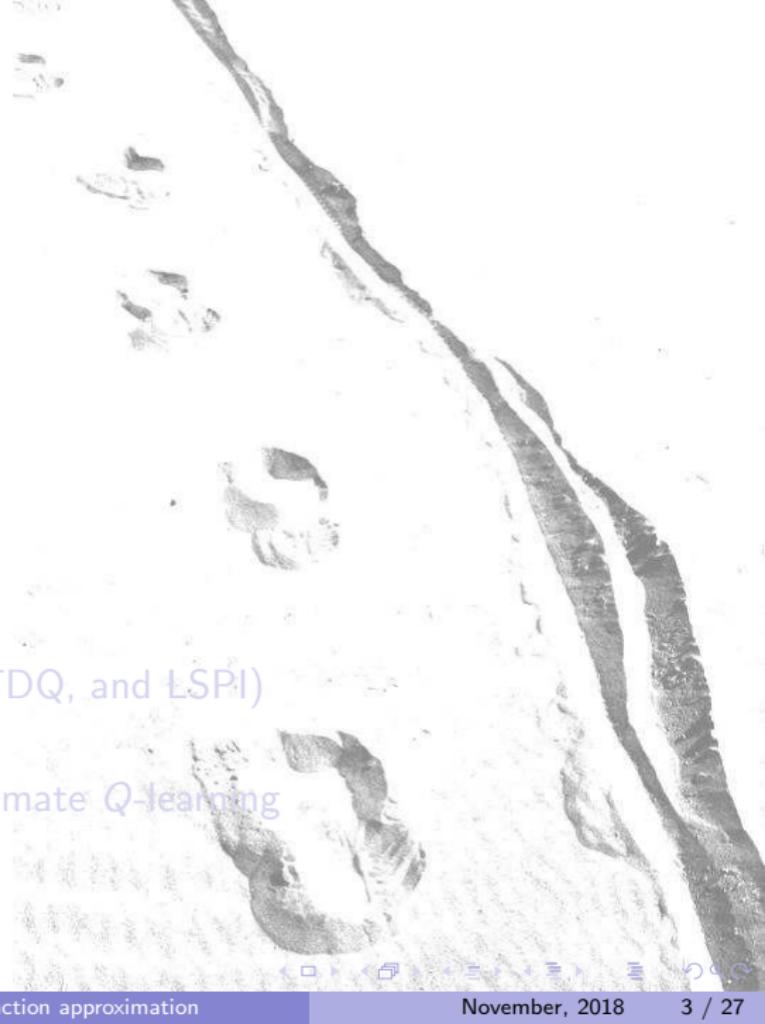
Cybernetic, NTNU  
Elec. Eng., Chalmers

TUM lectures on RL

# Outline

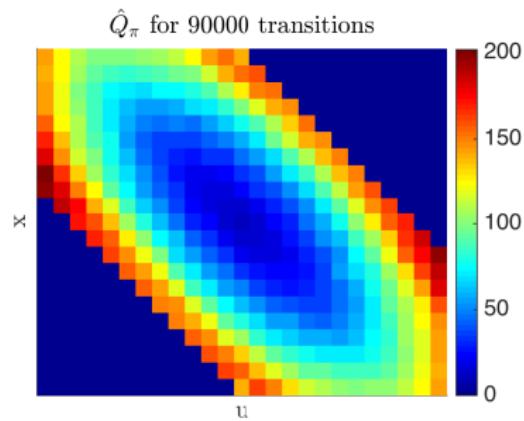
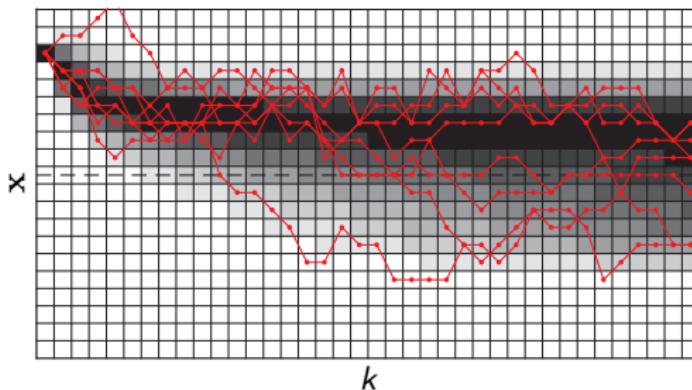
- 
- 1 Introduction
  - 2 Value function fitting
  - 3 Gradient-based approaches
  - 4 Action-value function fitting
  - 5 The “LS” family (LSTD, LSTDQ, and LSPI)
  - 6 Some observations on approximate Q-learning

# Outline

- 
- 1 Introduction
  - 2 Value function fitting
  - 3 Gradient-based approaches
  - 4 Action-value function fitting
  - 5 The “LS” family (LSTD, LSTDQ, and LSPI)
  - 6 Some observations on approximate Q-learning

## Why function approximations?

We have so far looked at...



- Example with small discrete state and input spaces, i.e. e.g.

$$x \in \{-10, -9, \dots, 9, 10\}, \quad u \in \{-12, -11, \dots, 11, 12\}$$

- Value functions  $\hat{V}_\pi$ ,  $\hat{Q}_\pi$ ,  $\hat{V}_*$ ,  $\hat{Q}_*$ , policy  $\pi(x)$ ,  $\pi_*(x)$  can be handled in a **tabular** fashion, e.g.

$$\hat{V}_\pi(-10) = 2.39, \quad \hat{V}_\pi(-9) = 5.43, \dots$$

This approach works only for small discrete spaces, for large and/or continuous spaces, we need to represent all these functions differently

## Approximating the value functions - Examples

- **Polynomials**, e.g.

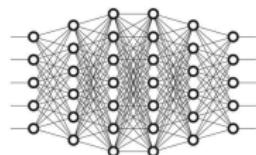
$$\hat{V}_\theta(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top S_\theta \mathbf{x} + \mathbf{f}_\theta^\top \mathbf{x} + c_\theta$$

- **Radial Basis Functions**, e.g.

$$\hat{V}(\mathbf{x}) = \sum_{k=1}^N \theta_k e^{-\frac{1}{2} \|\mathbf{x} - \mathbf{x}_k\|^2}$$

where  $\mathbf{x}_{1,\dots,N}$  is a grid of the state space

- **DNN**



- **What else? Stay tune...**

## Approximating the value functions - Examples

- **Polynomials**, e.g.

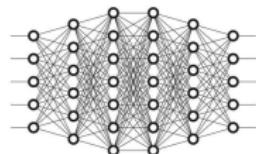
$$\hat{V}_\theta(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top S_\theta \mathbf{x} + \mathbf{f}_\theta^\top \mathbf{x} + c_\theta$$

- **Radial Basis Functions**, e.g.

$$\hat{V}(\mathbf{x}) = \sum_{k=1}^N \theta_k e^{-\frac{1}{2} \|\mathbf{x} - \mathbf{x}_k\|^2}$$

where  $\mathbf{x}_1, \dots, \mathbf{x}_N$  is a grid of the state space

- **DNN**



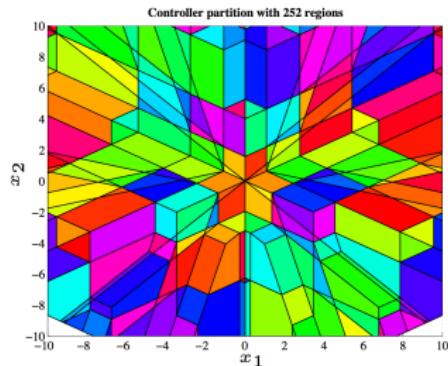
- **What else? Stay tune...**

**Remark:** suppose we have

- Linear, deterministic dynamics
- Quadratic cost
- Affine constraints

then  $V_*$  is piecewise quadratic

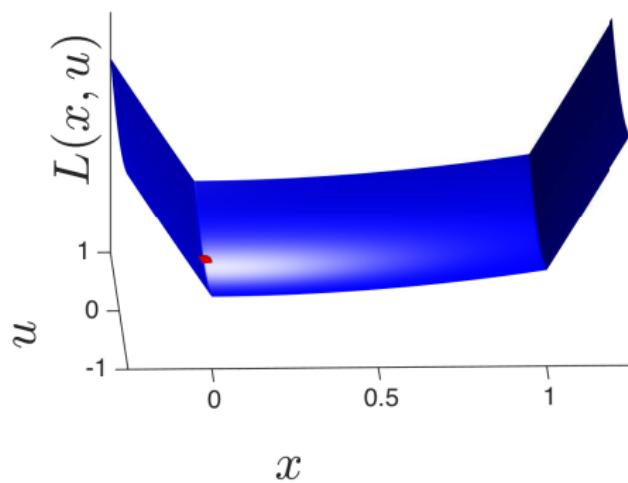
Regions can be arbitrarily complex...



Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

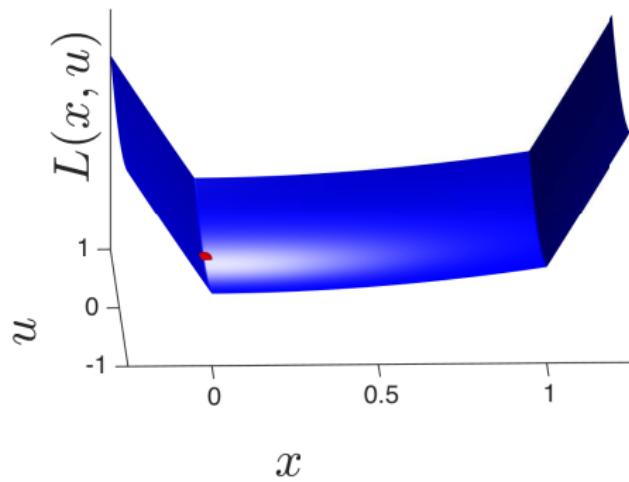
Stage cost



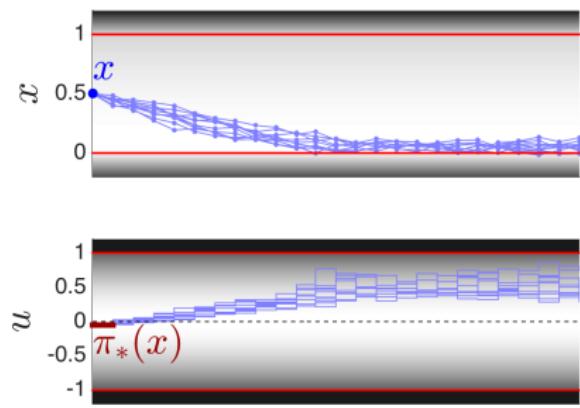
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



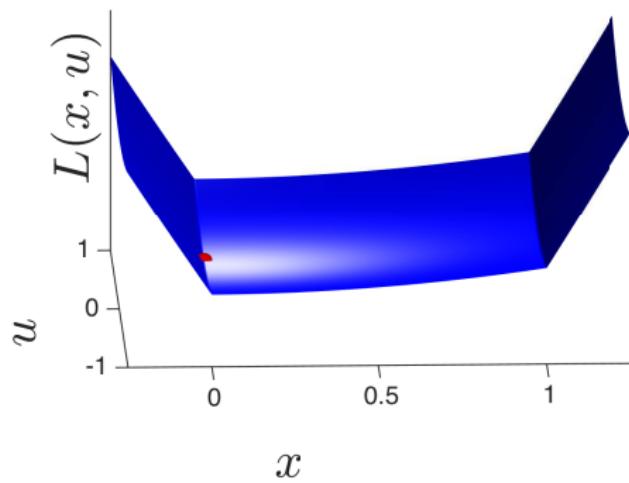
Trajectories



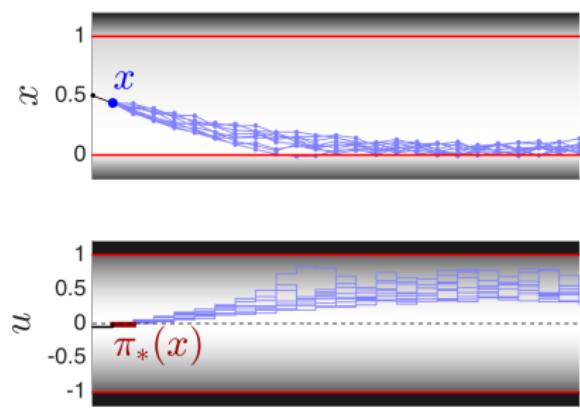
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



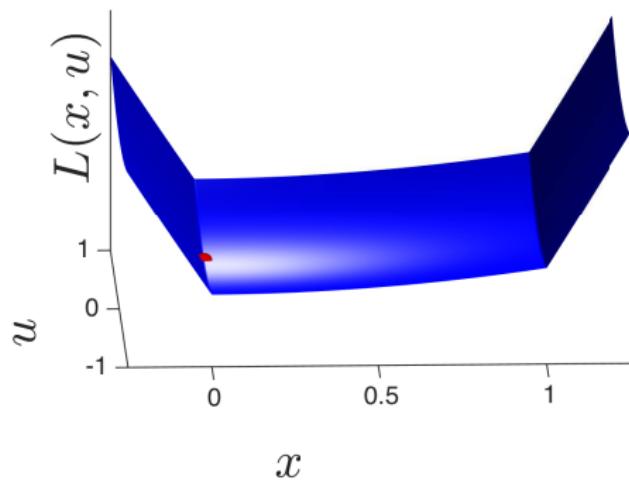
Trajectories



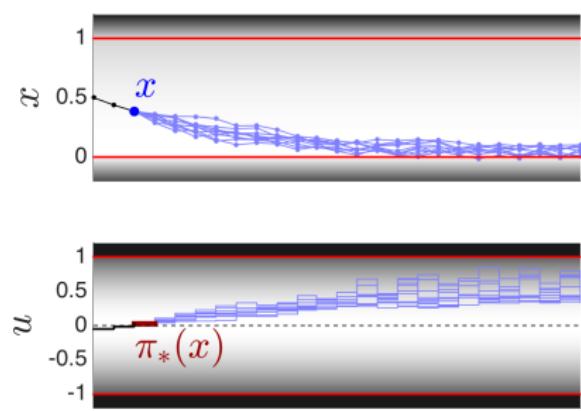
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



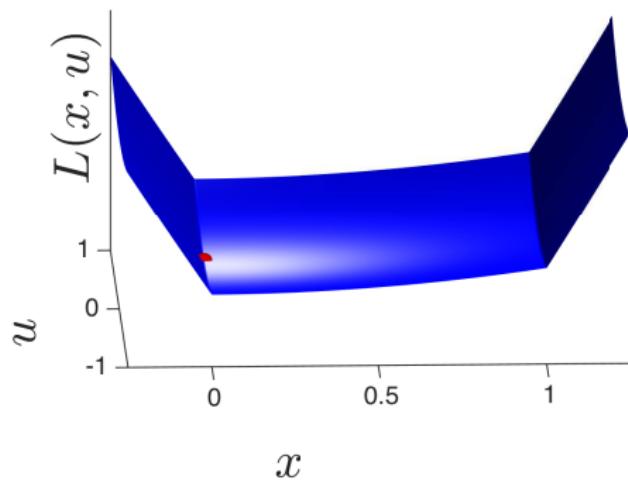
Trajectories



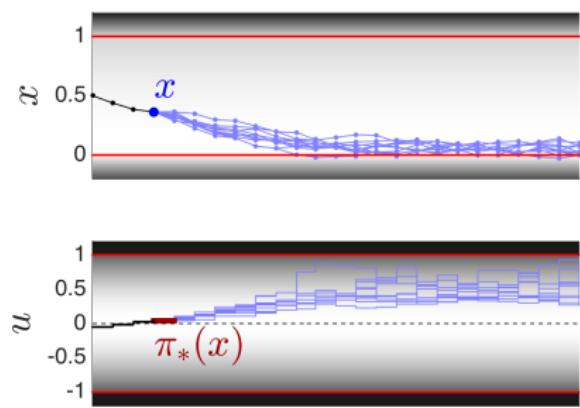
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



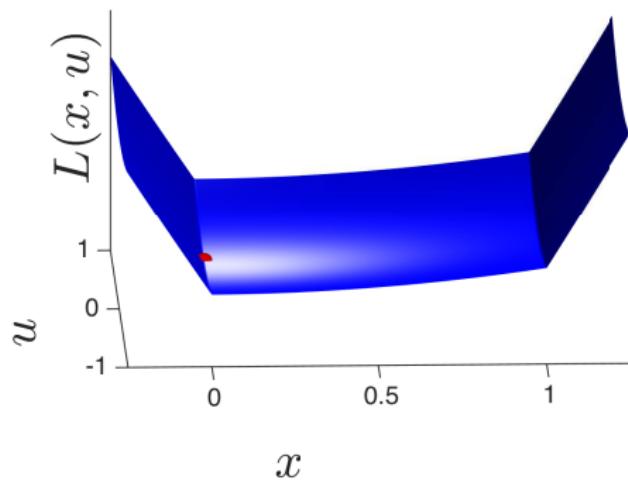
Trajectories



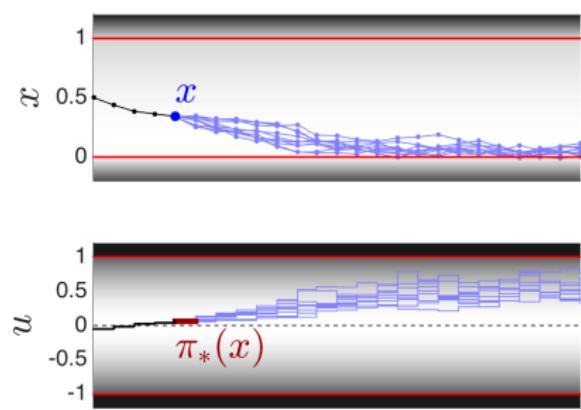
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



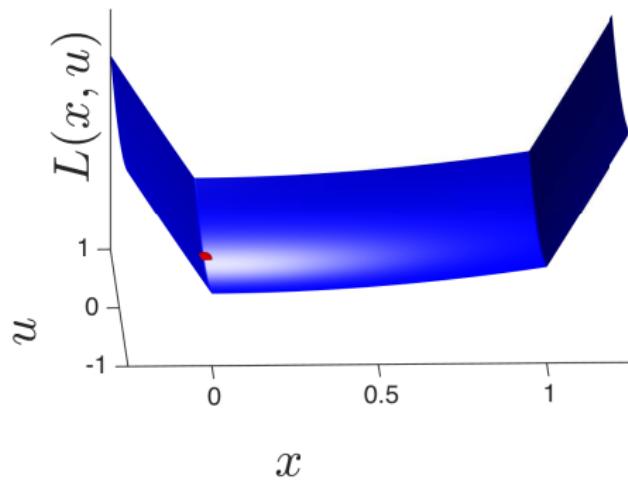
Trajectories



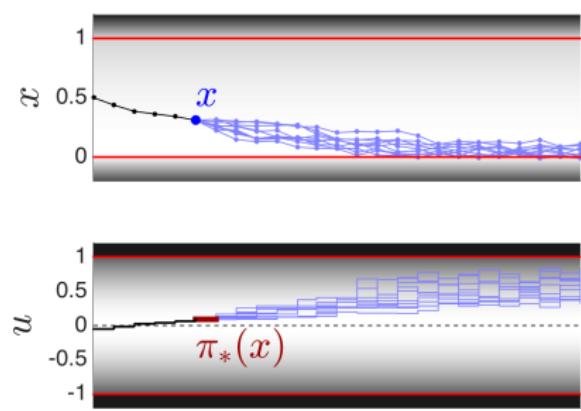
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



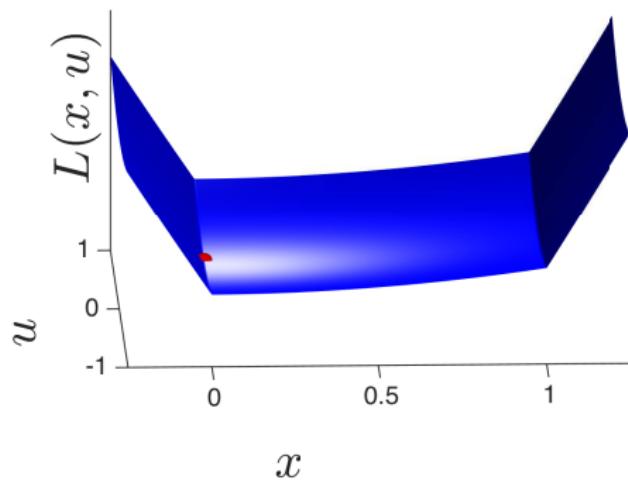
Trajectories



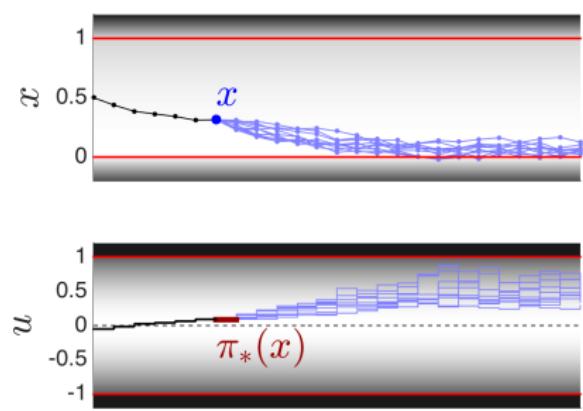
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



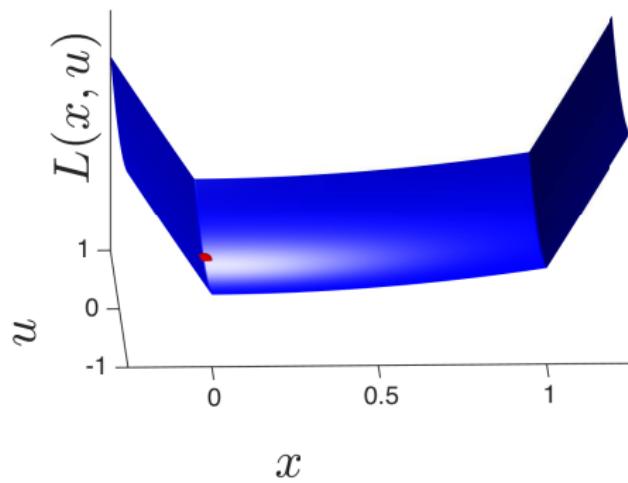
Trajectories



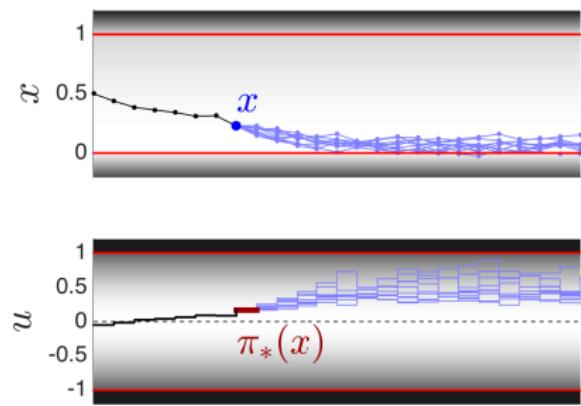
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



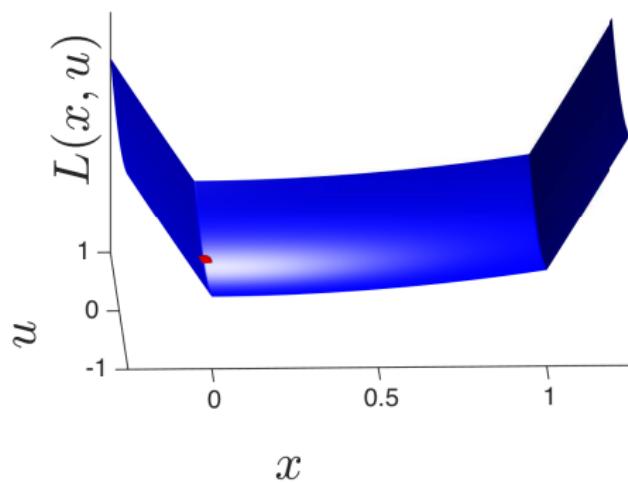
Trajectories



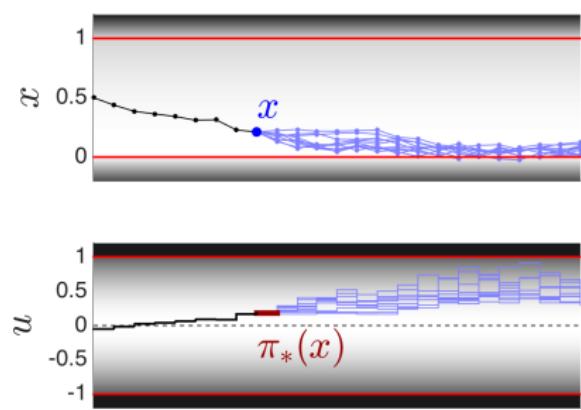
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



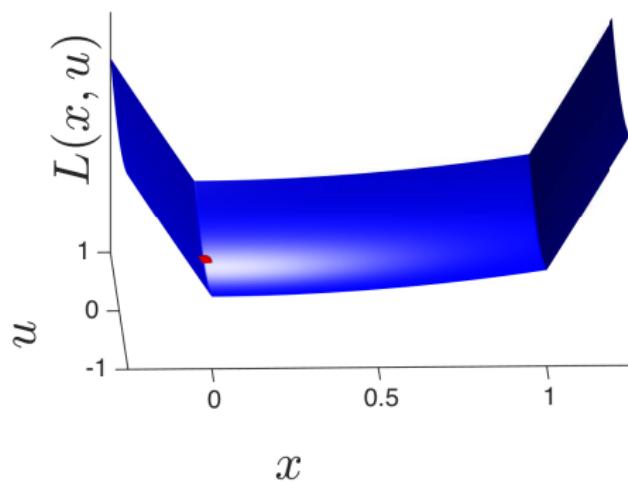
Trajectories



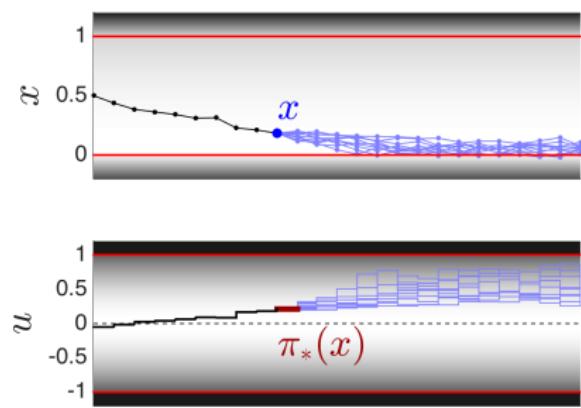
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



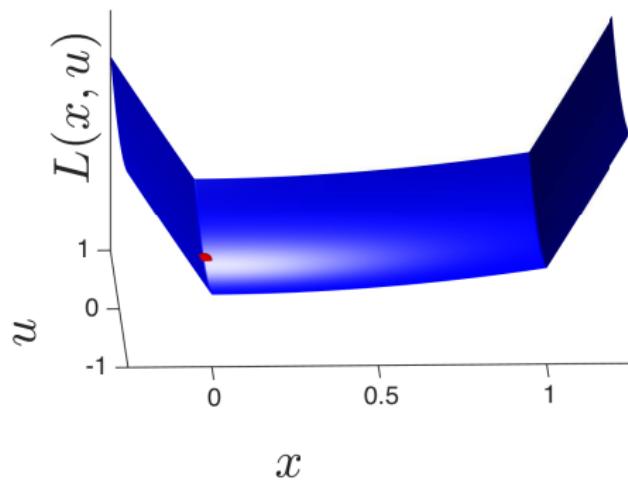
Trajectories



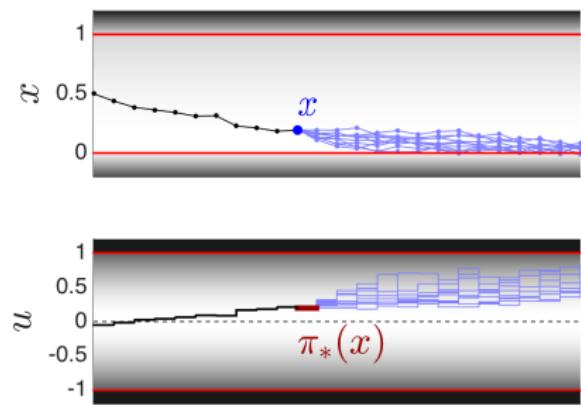
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



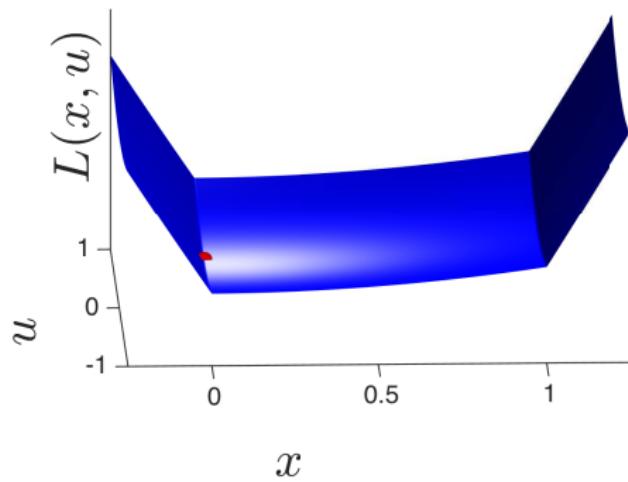
Trajectories



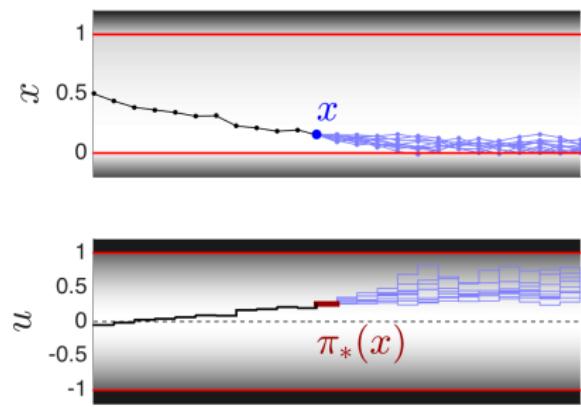
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



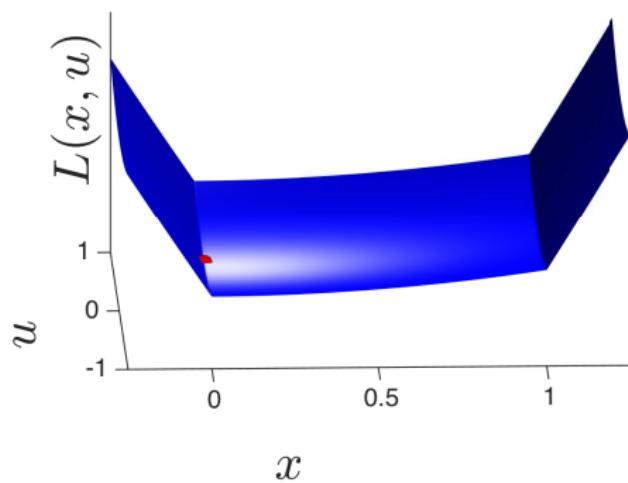
Trajectories



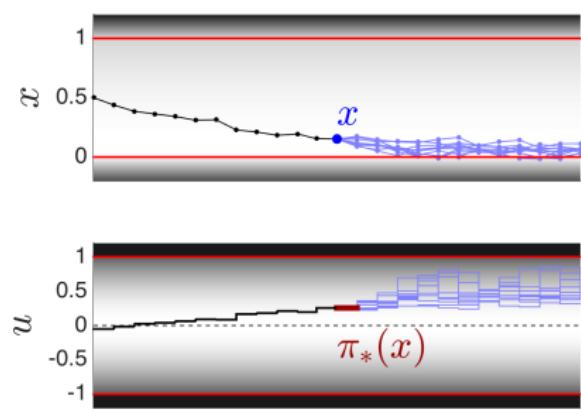
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



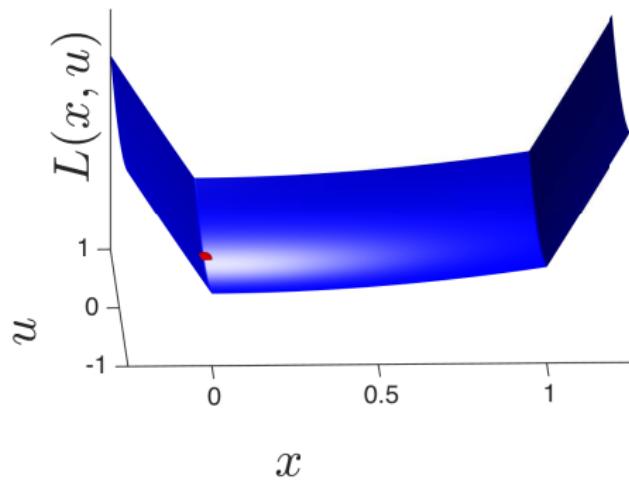
Trajectories



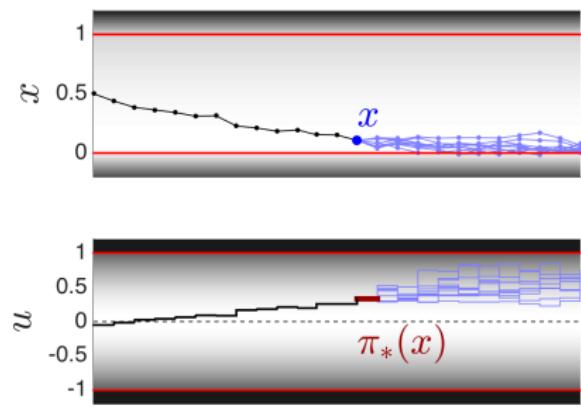
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



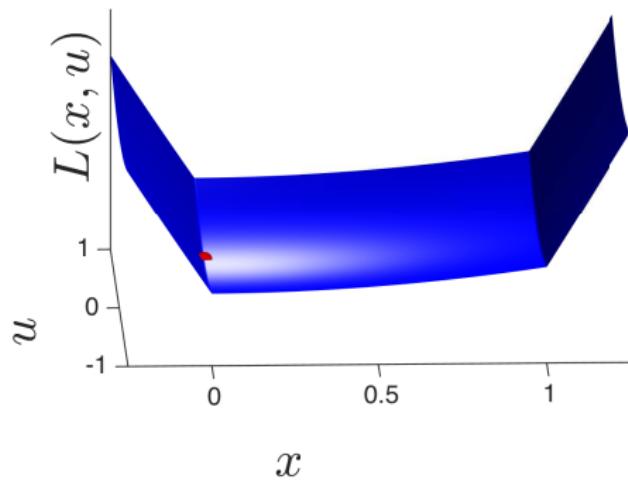
Trajectories



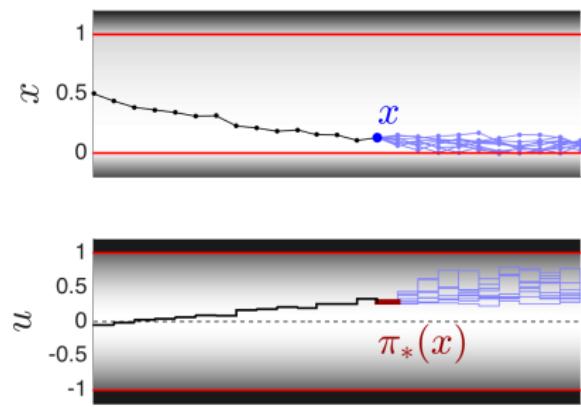
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



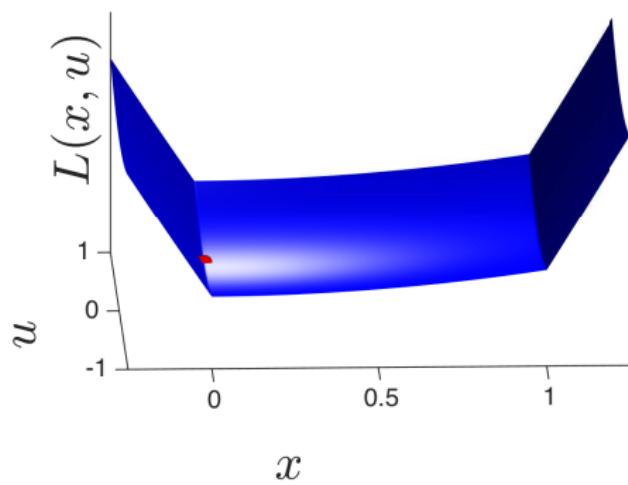
Trajectories



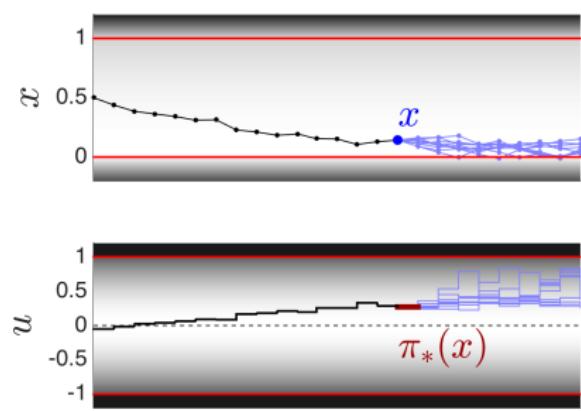
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



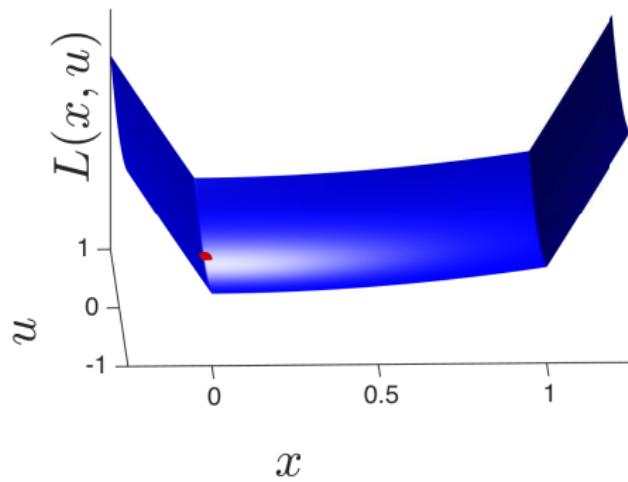
Trajectories



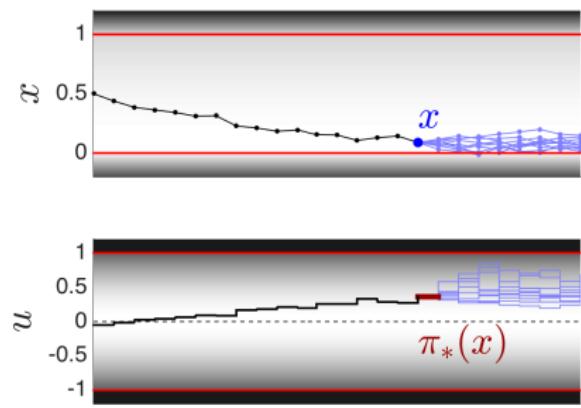
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



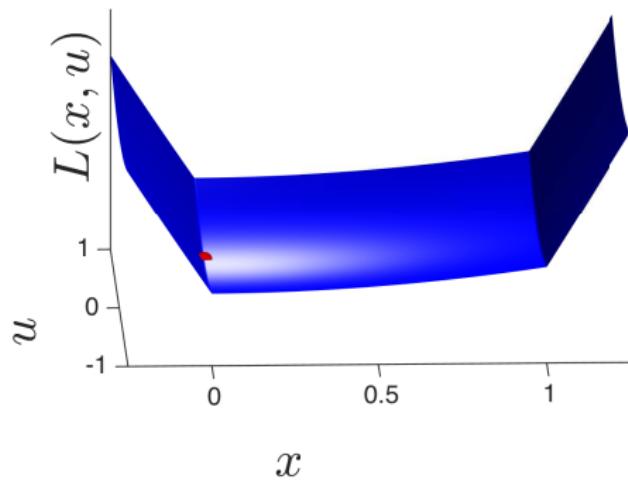
Trajectories



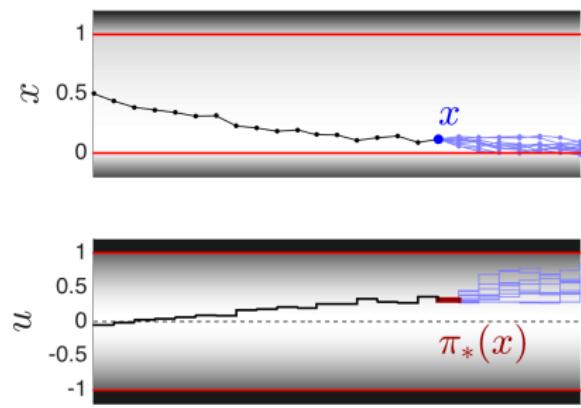
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



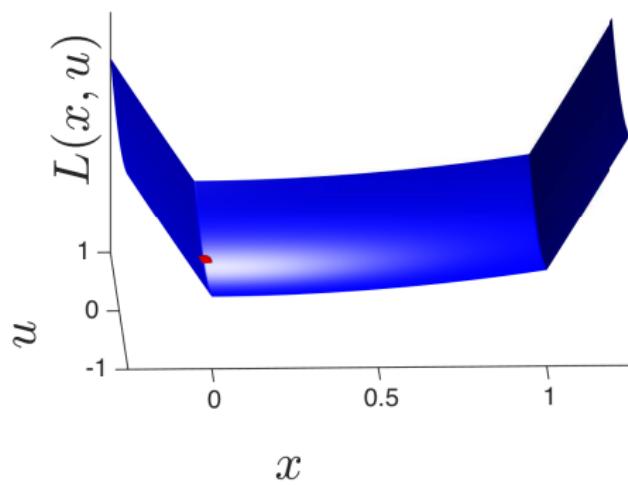
Trajectories



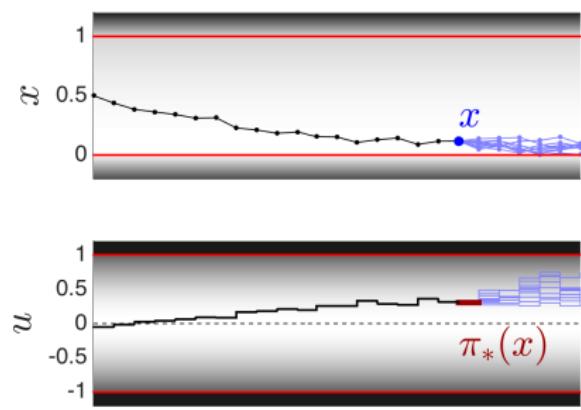
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



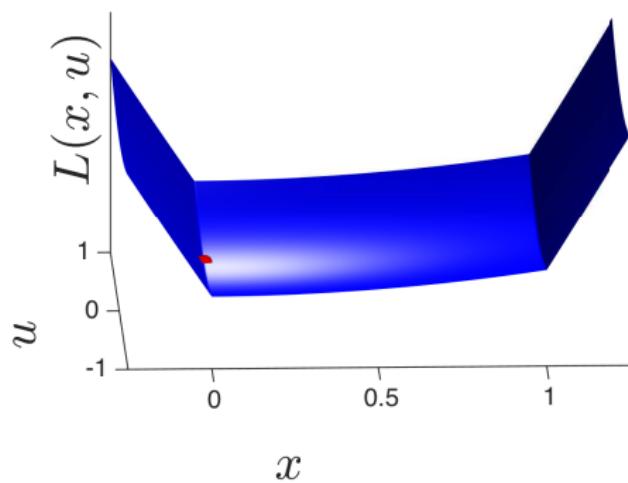
Trajectories



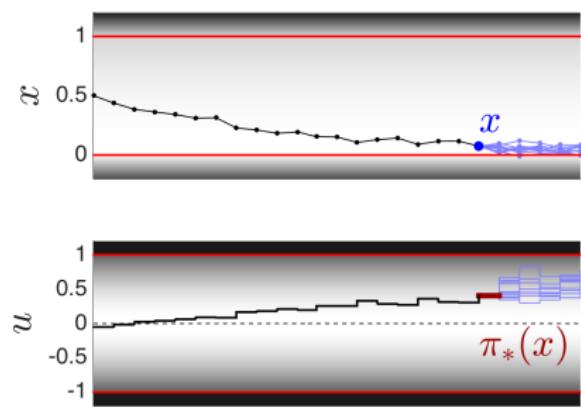
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Stage cost



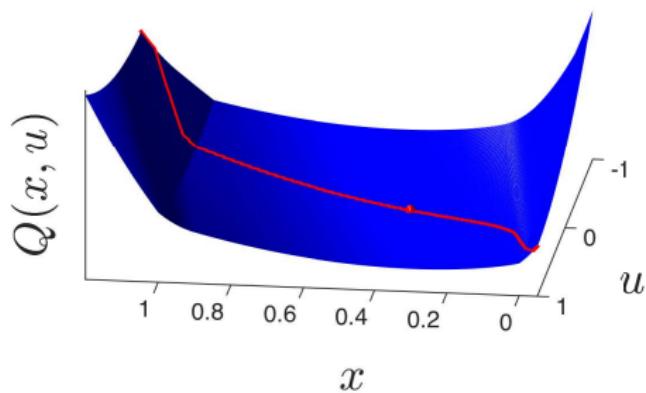
Trajectories



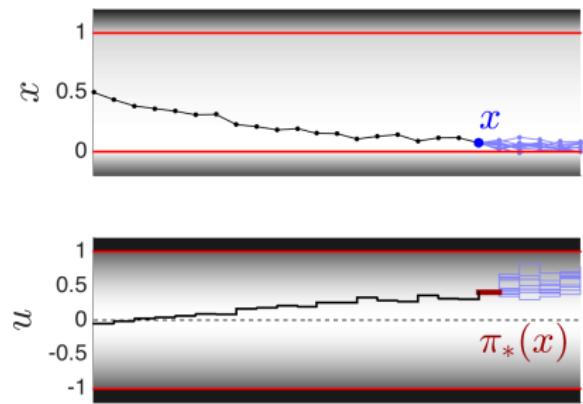
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Action-value function



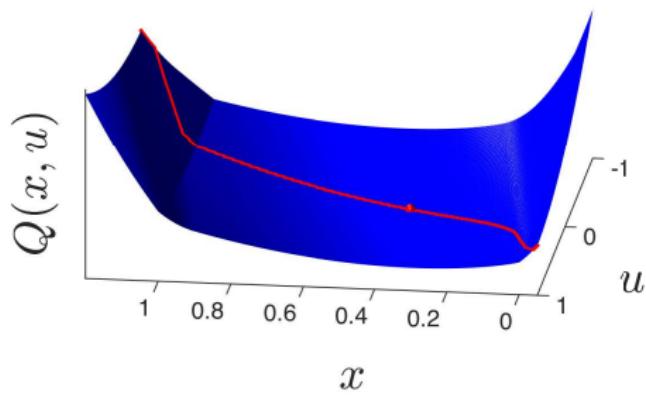
Trajectories



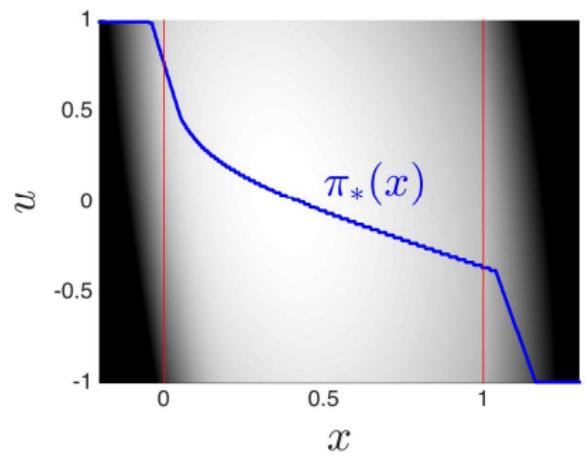
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Action-value function



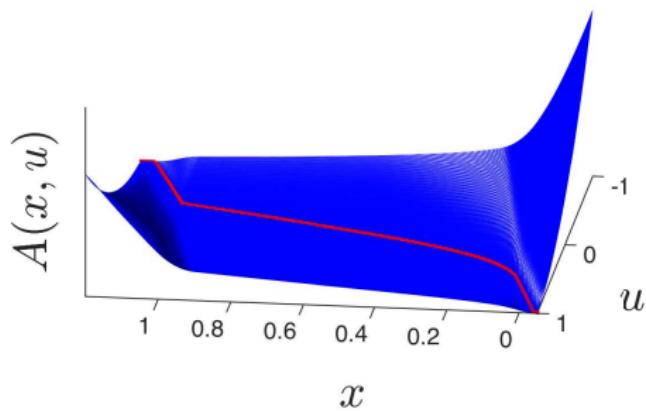
Optimal policy



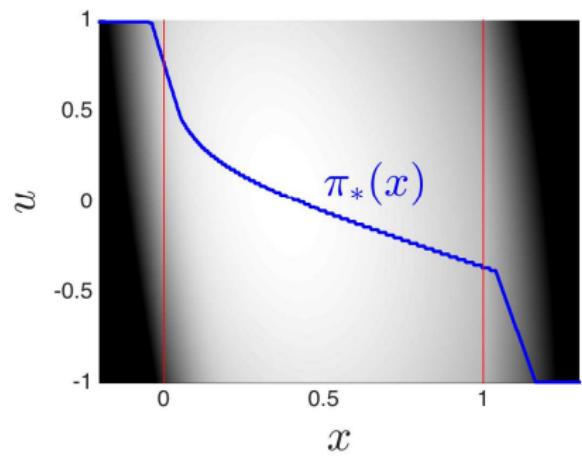
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

Advantage function



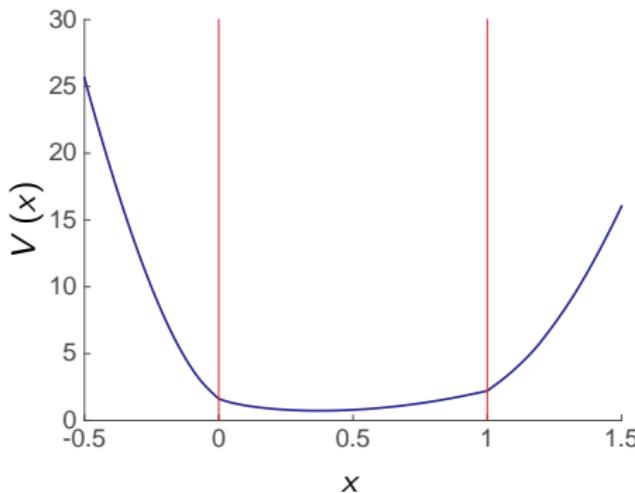
Optimal policy



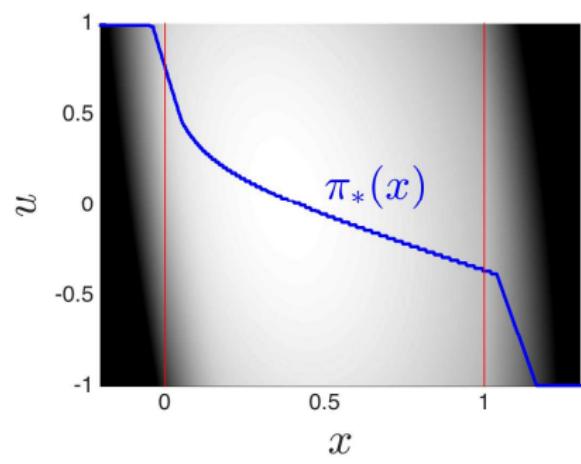
Let's use the following example

- State  $x \in \mathbb{R}$  and input  $u \in [-1, 1]$
- Dynamics:  $x_+ = x + 0.1u + e$  with  $e \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2 + \text{strong penalty for } x \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Noise  $e$  "pushing"  $x$  "against left wall"

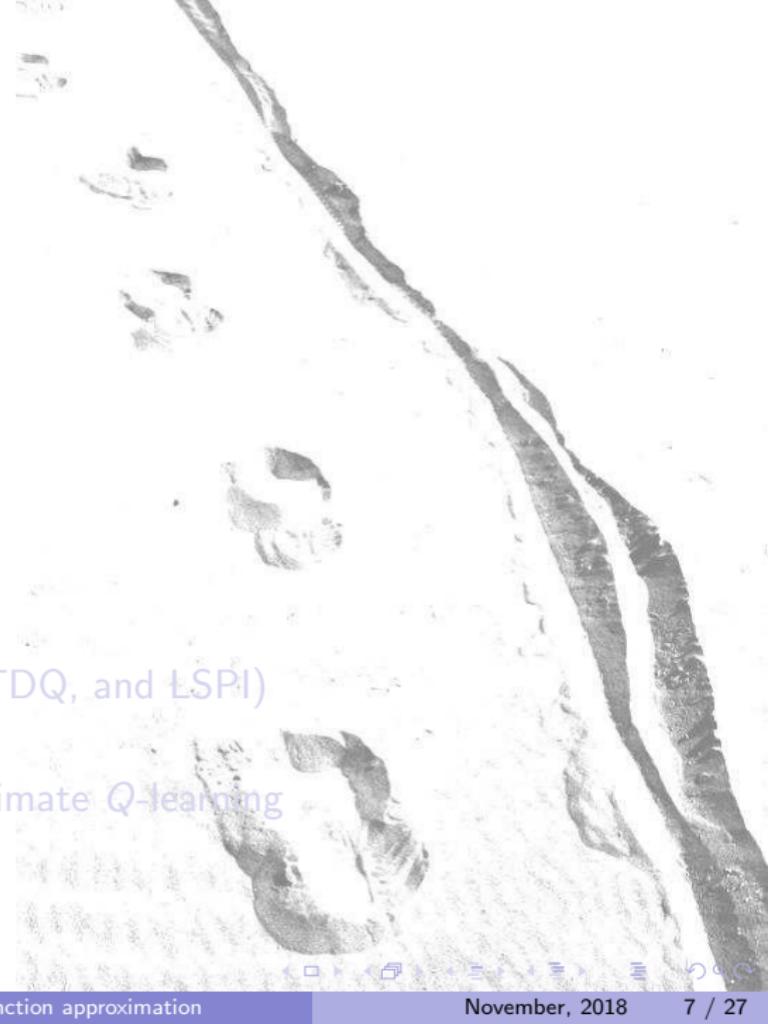
Value function



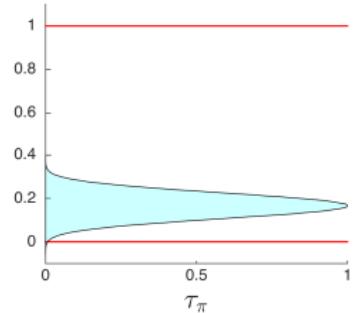
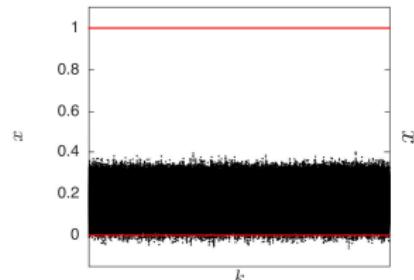
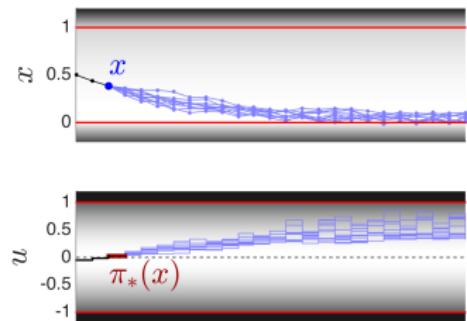
Optimal policy



# Outline

- 
- 1 Introduction
  - 2 Value function fitting
  - 3 Gradient-based approaches
  - 4 Action-value function fitting
  - 5 The “LS” family (LSTD, LSTDQ, and LSPI)
  - 6 Some observations on approximate Q-learning

## Reminder - Expected value over Markov chains

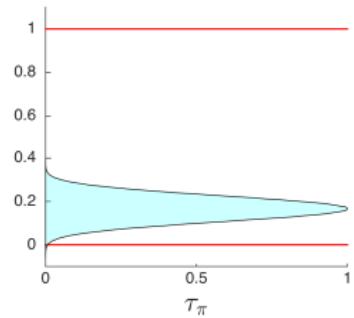
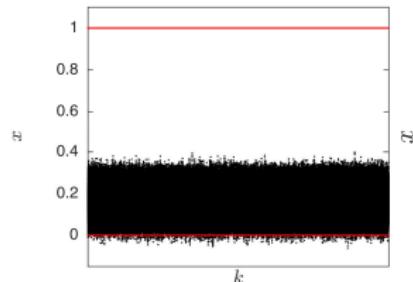
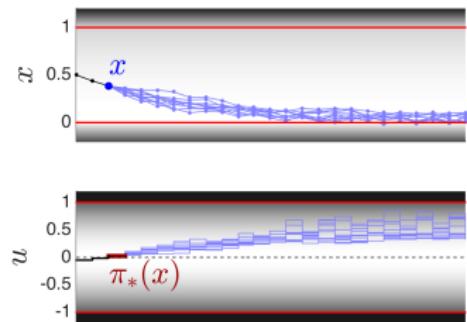


$$\mathbb{P} [\mathbf{x}_k | \mathbf{x}_0] = \int \prod_{i=0}^{k-1} \mathbb{P} [\mathbf{x}_{i+1} | \mathbf{x}_i, \boldsymbol{\pi}(\mathbf{x}_i)] d\mathbf{x}_1 \dots d\mathbf{x}_{k-1}$$

$$V(\mathbf{x}_0) = \sum_{k=0}^{\infty} \int \gamma^k L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) \mathbb{P} [\mathbf{x}_k | \mathbf{x}_0] d\mathbf{x}_k$$

$$J(\boldsymbol{\pi}) = \mathbb{E} [V(\mathbf{x}_0)] = \int V(\mathbf{x}_0) \mathbb{P} [\mathbf{x}_0] d\mathbf{x}_0 = \mathbb{E}_{\tau_{\boldsymbol{\pi}}} [L(\mathbf{x}, \mathbf{u})]$$

## Reminder - Expected value over Markov chains



$$\mathbb{P}[\mathbf{x}_k | \mathbf{x}_0] = \int \prod_{i=0}^{k-1} \mathbb{P}[\mathbf{x}_{i+1} | \mathbf{x}_i, \pi(\mathbf{x}_i)] d\mathbf{x}_1 \dots d\mathbf{x}_{k-1}$$

$$V(\mathbf{x}_0) = \sum_{k=0}^{\infty} \int \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k)) \mathbb{P}[\mathbf{x}_k | \mathbf{x}_0] d\mathbf{x}_k$$

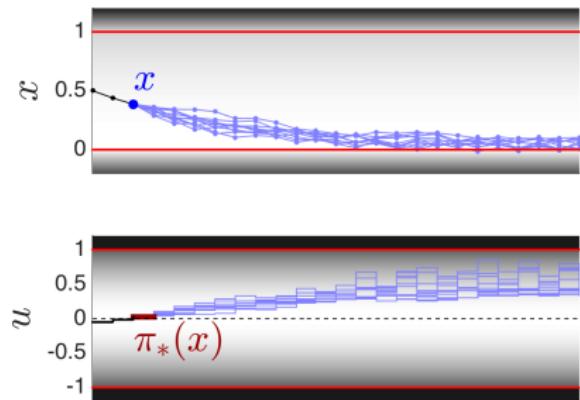
$$J(\pi) = \mathbb{E}[V(\mathbf{x}_0)] = \int V(\mathbf{x}_0) \mathbb{P}[\mathbf{x}_0] d\mathbf{x}_0 = \mathbb{E}_{\tau_\pi}[L(\mathbf{x}, \mathbf{u})]$$

When we approximate the functions ( $V, Q, \pi$ ), "where" we build them (i.e. under which distribution) will matter

## Core principle - Function fitting

**Value function fitting:**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} \left[ (V_{\pi}(x) - V_{\theta}(x))^2 \right]$$



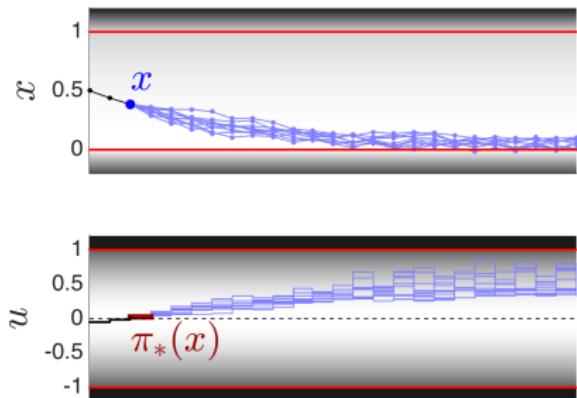
## Core principle - Function fitting

**Value function fitting:**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} \left[ (V_{\pi}(x) - V_{\theta}(x))^2 \right]$$

Optimal parameter  $\hat{\theta}$  satisfies:

$$\mathbb{E}_{\tau_{\pi}} [(V_{\pi}(x) - V_{\theta}(x)) \nabla_{\theta} V_{\theta}(x)] = 0$$



## Core principle - Function fitting

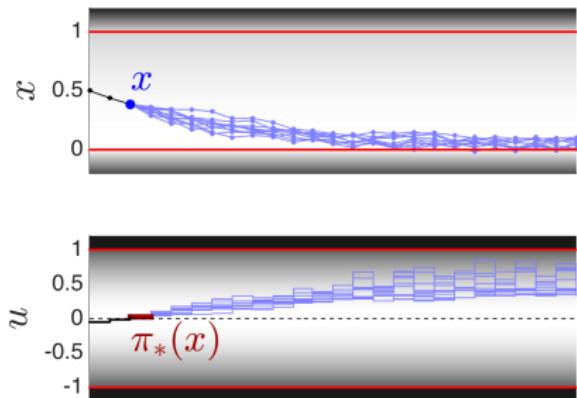
**Value function fitting:**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi} \left[ (V_\pi(x) - V_\theta(x))^2 \right]$$

**Optimal parameter  $\hat{\theta}$  satisfies:**

$$\mathbb{E}_{\tau_\pi} [(V_\pi(x) - V_\theta(x)) \nabla_\theta V_\theta(x)] = 0$$

**Where do we get  $V_\pi$  from?**



## Core principle - Function fitting

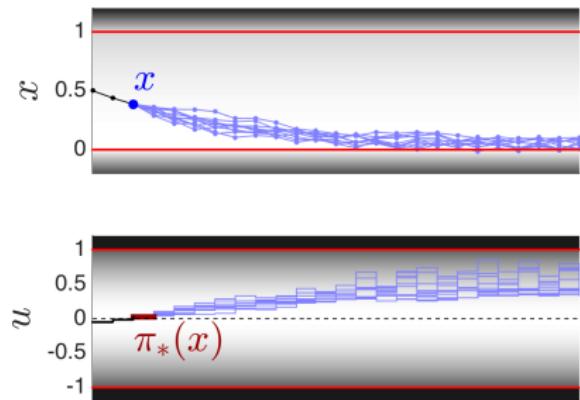
**Value function fitting:**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} \left[ (V_{\pi}(x) - V_{\theta}(x))^2 \right]$$

**Optimal parameter  $\hat{\theta}$  satisfies:**

$$\mathbb{E}_{\tau_{\pi}} [(V_{\pi}(x) - V_{\theta}(x)) \nabla_{\theta} V_{\theta}(x)] = 0$$

**Where do we get  $V_{\pi}$  from?**



$$V_{\pi}(x) = \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k L(x_k, \pi(x_k)) \mid x_0 = x \right] = L(x, \pi(x)) + \gamma \mathbb{E} [V_{\pi}(x_+) \mid x]$$

# Core principle - Function fitting

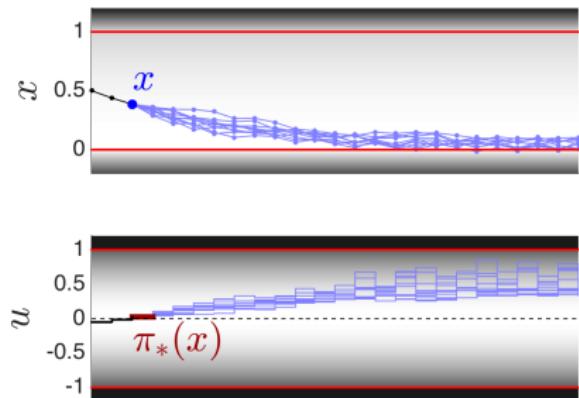
**Value function fitting:**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} \left[ (V_{\pi}(x) - V_{\theta}(x))^2 \right]$$

**Optimal parameter  $\hat{\theta}$  satisfies:**

$$\mathbb{E}_{\tau_{\pi}} [(V_{\pi}(x) - V_{\theta}(x)) \nabla_{\theta} V_{\theta}(x)] = 0$$

**Where do we get  $V_{\pi}$  from?**



$$V_{\pi}(x) = \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k L(x_k, \pi(x_k)) \mid x_0 = x \right] = L(x, \pi(x)) + \gamma \mathbb{E} [V_{\pi}(x_+) \mid x]$$

**MC evaluation**

$$V_{\pi}(x) \approx \sum_{i=0}^{\infty} \gamma^i L(x_i, \pi(x_i))$$

(noisy)

# Core principle - Function fitting

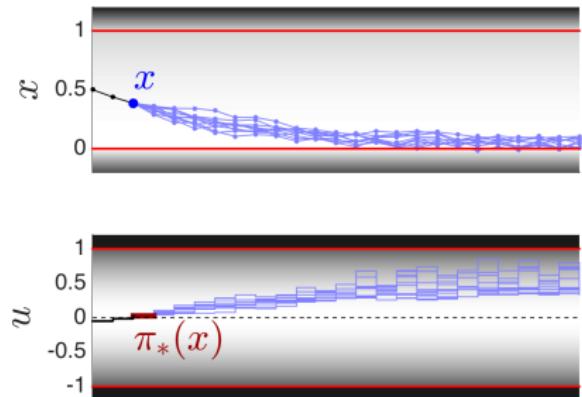
**Value function fitting:**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi} \left[ (V_\pi(x) - V_\theta(x))^2 \right]$$

**Optimal parameter  $\hat{\theta}$  satisfies:**

$$\mathbb{E}_{\tau_\pi} [(V_\pi(x) - V_\theta(x)) \nabla_\theta V_\theta(x)] = 0$$

**Where do we get  $V_\pi$  from?**



$$V_\pi(x) = \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k L(x_k, \pi(x_k)) \mid x_0 = x \right] = L(x, \pi(x)) + \gamma \mathbb{E} [V_\pi(x_+) \mid x]$$

**MC evaluation**

$$V_\pi(x) \approx \sum_{i=0}^{\infty} \gamma^i L(x_i, \pi(x_i))$$

(noisy)

**TD evaluation**

$$\begin{aligned} V_\pi(x) &\approx L(x, \pi(x)) + \gamma V_\pi(x_+) \\ &\approx L(x, \pi(x)) + \gamma V_\theta(x_+) \end{aligned}$$

(biased)

## Core principle - Function fitting

**Value function fitting:**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} \left[ (V_{\pi}(x) - V_{\theta}(x))^2 \right]$$

**Optimal parameter  $\hat{\theta}$  satisfies:**

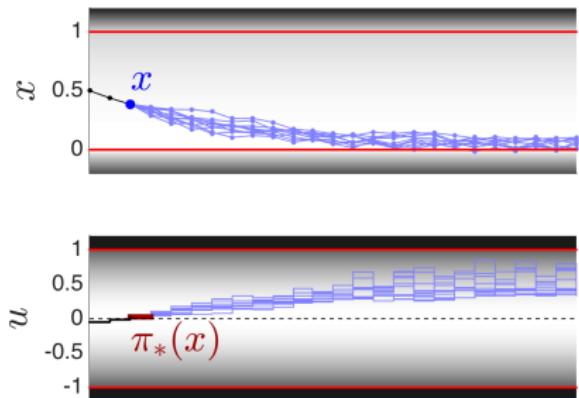
$$\mathbb{E}_{\tau_{\pi}} [(V_{\pi}(x) - V_{\theta}(x)) \nabla_{\theta} V_{\theta}(x)] = 0$$

Where do we get  $V_{\pi}$  from?

Using  $V_{\pi}(x) \approx L(x, \pi(x)) + \gamma V_{\theta}(x_+)$

Compute parameter  $\hat{\theta}$  using:

$$\mathbb{E}_{\tau_{\pi}} [(L(x, \pi(x)) + \gamma V_{\theta}(x_+) - V_{\theta}(x)) \nabla_{\theta} V_{\theta}(x)] = 0$$



# Core principle - Function fitting

**Value function fitting:**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi} \left[ (V_\pi(\mathbf{x}) - V_\theta(\mathbf{x}))^2 \right]$$

**Optimal parameter  $\hat{\theta}$  satisfies:**

$$\mathbb{E}_{\tau_\pi} [(V_\pi(\mathbf{x}) - V_\theta(\mathbf{x})) \nabla_\theta V_\theta(\mathbf{x})] = 0$$

Where do we get  $V_\pi$  from?

Using  $V_\pi(\mathbf{x}) \approx L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_\theta(\mathbf{x}_+)$

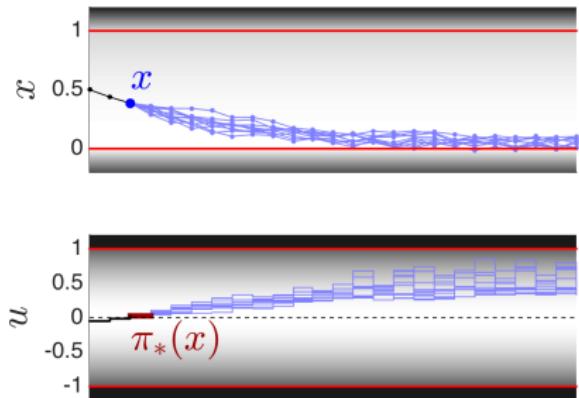
Compute parameter  $\hat{\theta}$  using:

$$\mathbb{E}_{\tau_\pi} [(L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_\theta(\mathbf{x}_+) - V_\theta(\mathbf{x})) \nabla_\theta V_\theta(\mathbf{x})] = 0$$

equivalently:

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_\theta V_\theta] = 0$$

where  $\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_\theta(\mathbf{x}_+) - V_\theta(\mathbf{x})$

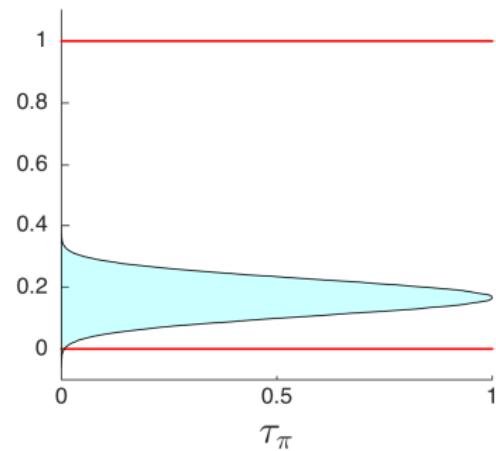
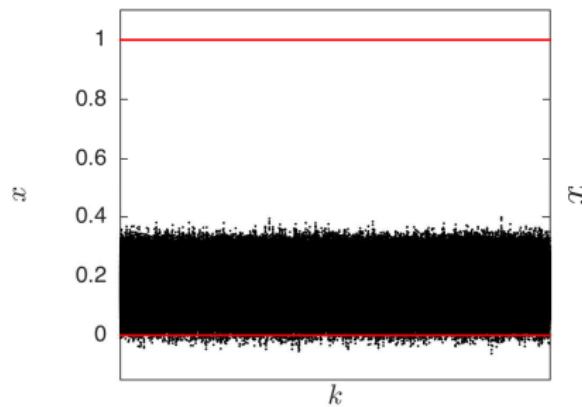


## Function fitting - Example

- Dynamics:  $\mathbf{x}_+ = \mathbf{x} + 0.1\mathbf{u} + \mathbf{e}$  with  $\mathbf{e} \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(\mathbf{x}, \mathbf{u}) = \frac{1}{2}\mathbf{x}^2 + \frac{1}{2}\mathbf{u}^2 + \text{strong penalty for } \mathbf{x} \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Let us e.g. select  $\pi(\mathbf{x}) = -\frac{3}{2}\mathbf{x} + \frac{3}{4}$  clipped to  $[-1, 1]$

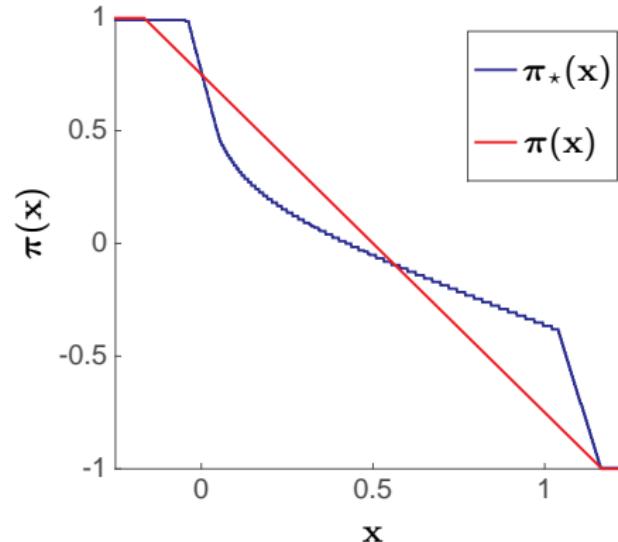
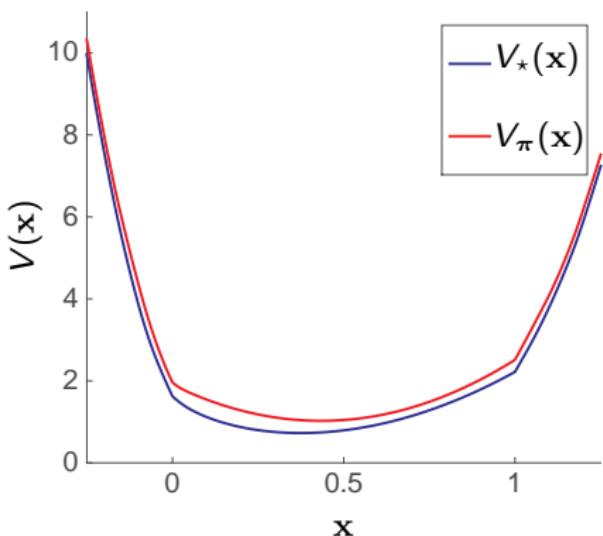
## Function fitting - Example

- Dynamics:  $\mathbf{x}_+ = \mathbf{x} + 0.1\mathbf{u} + \mathbf{e}$  with  $\mathbf{e} \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(\mathbf{x}, \mathbf{u}) = \frac{1}{2}\mathbf{x}^2 + \frac{1}{2}\mathbf{u}^2 + \text{strong penalty for } \mathbf{x} \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Let us e.g. select  $\pi(\mathbf{x}) = -\frac{3}{2}\mathbf{x} + \frac{3}{4}$  clipped to  $[-1, 1]$
- Markov Chain trajectories  $\tau_\pi$

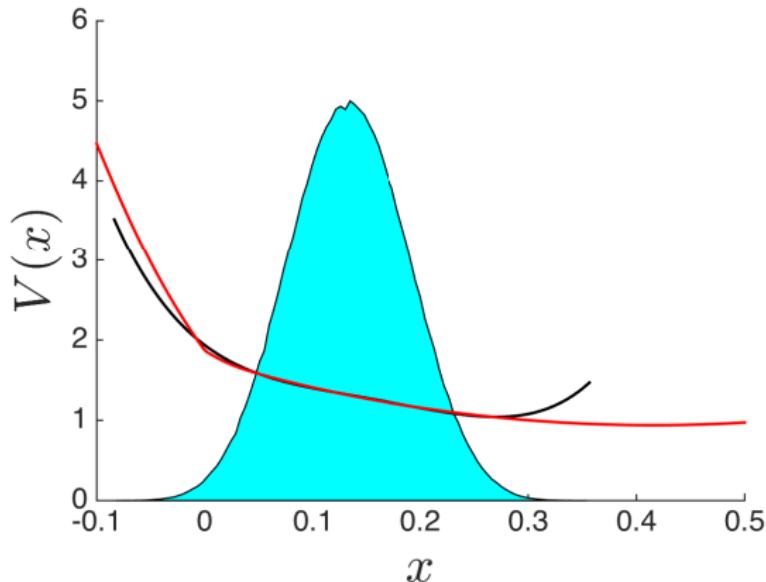


## Function fitting - Example

- Dynamics:  $\mathbf{x}_+ = \mathbf{x} + 0.1\mathbf{u} + \mathbf{e}$  with  $\mathbf{e} \sim \mathcal{U}([-0.1, 0])$
- Stage cost:  $L(\mathbf{x}, \mathbf{u}) = \frac{1}{2}\mathbf{x}^2 + \frac{1}{2}\mathbf{u}^2 + \text{strong penalty for } \mathbf{x} \notin [0, 1]$
- Discount  $\gamma = 0.9$
- Let us e.g. select  $\pi(\mathbf{x}) = -\frac{3}{2}\mathbf{x} + \frac{3}{4}$  clipped to  $[-1, 1]$
- Markov Chain trajectories  $\tau_\pi$
- Resulting value function  $V_\pi(\mathbf{x}) \geq V_*(\mathbf{x})$  for all  $\mathbf{x}$

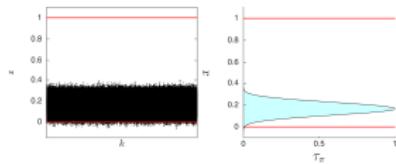


## Function fitting - Example



**Value function fitting:**

$$\hat{\theta} = \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} \left[ \left( V_{\pi} - \hat{V} \right)^2 \right] \iff \mathbb{E}_{\tau_{\pi}} \left[ \left( V_{\pi} - \hat{V} \right) \nabla_{\theta} \hat{V} \right] = 0$$



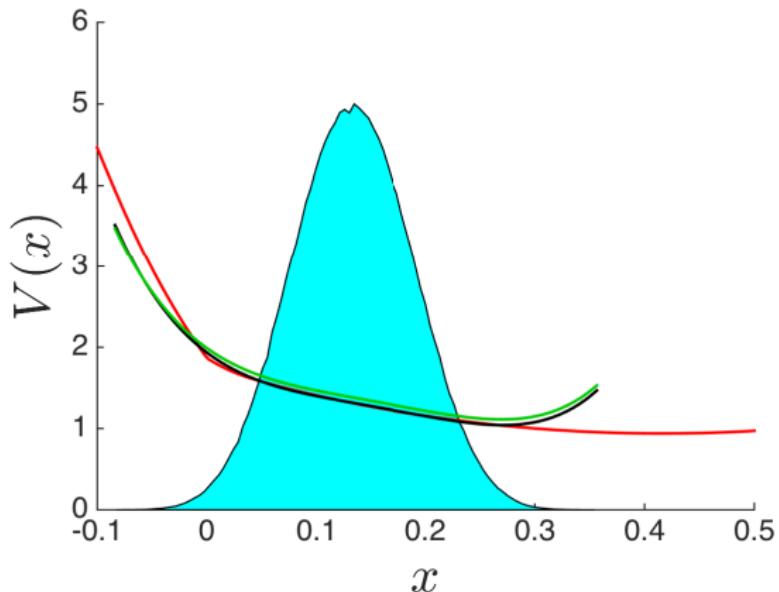
We label

$$V_{\pi}(x), \hat{V}(x)$$

We choose e.g.:

$$\hat{V}(x) = \sum_{k=0}^3 \hat{\theta}_k x^k$$

## Function fitting using target - Example



We label

$$V_\pi(x), \hat{V}(x), V_{TD}(x)$$

We choose e.g.:

$$\hat{V}(x) = \sum_{k=0}^4 \hat{\theta}_k x^k$$

**Value function fitting:**

$$\mathbb{E}_{\tau_\pi} \left[ \left( V_\pi - \hat{V} \right) \nabla_\theta \hat{V} \right] = 0$$

**TD-based fitting:** ( $V_\pi \approx L + \gamma V_{TD}^+$ )

$$\mathbb{E}_{\tau_\pi} \left[ \left( L + \gamma V_{TD}^+ - V_{TD} \right) \nabla_\theta V_{TD} \right] = 0$$

$\hat{V} \neq V_{TD}$  (bias)

## Remark - Fitting towards fixed target?

**Value function fitting:**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} [(V_{\pi} - V_{\theta})^2]$$

using TD approximation

$$V_{\pi}(\mathbf{x}) \approx L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+)$$

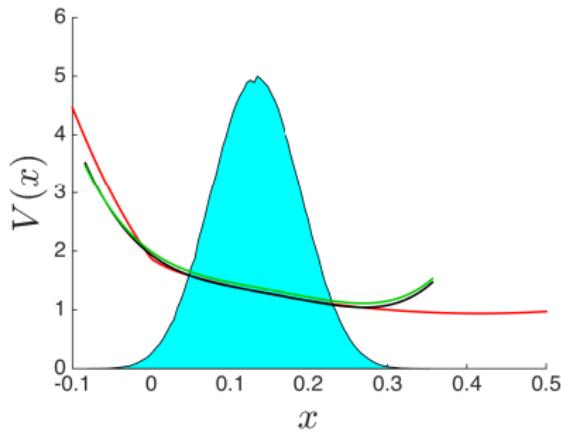
**Compute parameter  $\hat{\theta}$  using:**

$$\mathbb{E}_{\tau_{\pi}} [\delta \nabla_{\theta} V_{\theta}] = 0$$

$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

We label

$$V_*(\mathbf{x}), V_{\pi}(\mathbf{x}), \hat{V}(\mathbf{x}), V_{\text{TD}}(\mathbf{x})$$



## Remark - Fitting towards fixed target?

**Value function fitting:**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} [(V_{\pi} - V_{\theta})^2]$$

using TD approximation

$$V_{\pi}(\mathbf{x}) \approx L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+)$$

Note:  $\delta \approx V_{\pi} - V_{\theta}$  Do min of TD error?

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} [\delta^2]$$

hence solve  $\mathbb{E}_{\tau_{\pi}} [\delta \cdot \nabla_{\theta} \delta] = 0$

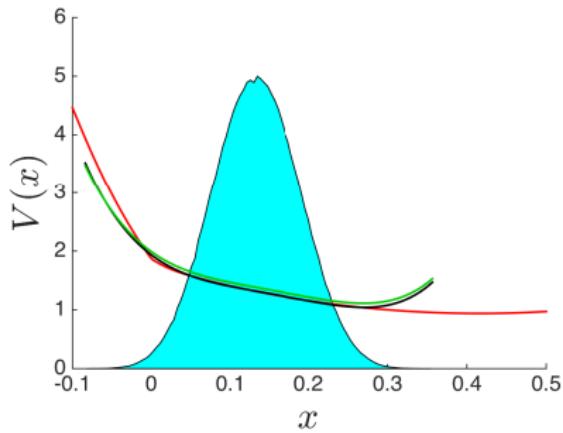
**Compute parameter  $\hat{\theta}$  using:**

$$\mathbb{E}_{\tau_{\pi}} [\delta \nabla_{\theta} V_{\theta}] = 0$$

$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

We label

$$V_*(\mathbf{x}), V_{\pi}(\mathbf{x}), \hat{V}(\mathbf{x}), V_{\text{TD}}(\mathbf{x})$$



## Remark - Fitting towards fixed target?

**Value function fitting:**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} [(V_{\pi} - V_{\theta})^2]$$

using TD approximation

$$V_{\pi}(\mathbf{x}) \approx L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+)$$

Note:  $\delta \approx V_{\pi} - V_{\theta}$  Do min of TD error?

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} [\delta^2]$$

hence solve  $\mathbb{E}_{\tau_{\pi}} [\delta \cdot \nabla_{\theta} \delta] = 0$

- generally **does not work**
- exception: deterministic system and “rich” parametrization (get  $\delta = 0$ )
- true  $\delta$  is not necessarily RMS-minimum

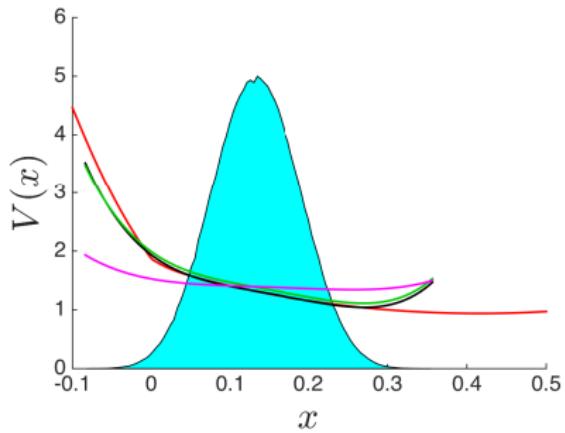
**Compute parameter  $\hat{\theta}$  using:**

$$\mathbb{E}_{\tau_{\pi}} [\delta \nabla_{\theta} V_{\theta}] = 0$$

$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

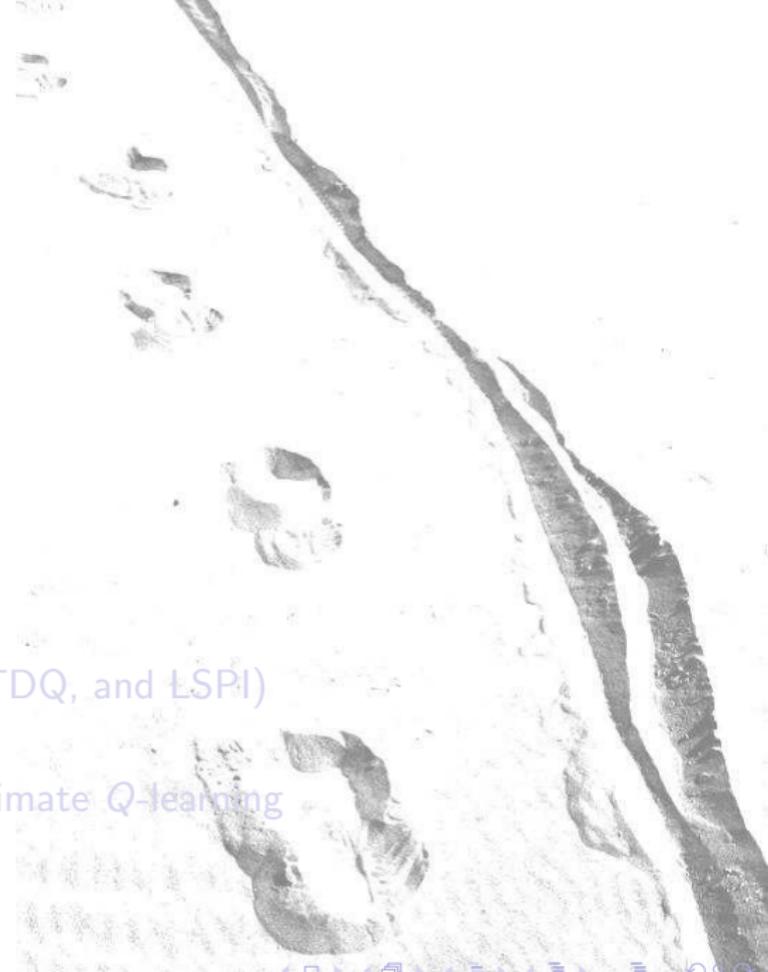
We label

$$V_*(\mathbf{x}), V_{\pi}(\mathbf{x}), \hat{V}(\mathbf{x}), V_{\text{TD}}(\mathbf{x})$$



# Outline

- 1 Introduction
- 2 Value function fitting
- 3 Gradient-based approaches
- 4 Action-value function fitting
- 5 The “LS” family (LSTD, LSTDQ, and LSPI)
- 6 Some observations on approximate Q-learning



Stochastic TD approach - Solving fitting  $\min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi} \left[ (V_\pi - V_\theta)^2 \right]$

**Gradient** approach to solve

$$\hat{\theta} = \arg \min_{\theta} \underbrace{\frac{1}{2} \mathbb{E}_{\tau_\pi} \left[ (V_\pi - V_\theta)^2 \right]}_{:=\phi(\theta)}$$

via gradient steps:

$$\hat{\theta} \leftarrow \hat{\theta} - \alpha \nabla_{\theta} \phi$$

Stochastic TD approach - Solving fitting  $\min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi} \left[ (V_\pi - V_\theta)^2 \right]$

**Gradient** approach to solve

$$\hat{\theta} = \arg \min_{\theta} \underbrace{\frac{1}{2} \mathbb{E}_{\tau_\pi} \left[ (V_\pi - V_\theta)^2 \right]}_{:=\phi(\theta)}$$

via gradient steps:

$$\hat{\theta} \leftarrow \hat{\theta} - \alpha \nabla_{\theta} \phi$$

Iteration reads as:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \mathbb{E}_{\tau_\pi} [(V_\pi - V_\theta) \nabla_{\theta} V_\theta]$$

**Stochastic gradient** approach:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha (V_\pi - V_\theta) \nabla_{\theta} V_\theta$$

Stochastic TD approach - Solving fitting  $\min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} \left[ (V_{\pi} - V_{\theta})^2 \right]$

**Gradient** approach to solve

$$\hat{\theta} = \arg \min_{\theta} \underbrace{\frac{1}{2} \mathbb{E}_{\tau_{\pi}} \left[ (V_{\pi} - V_{\theta})^2 \right]}_{:=\phi(\theta)}$$

via gradient steps:

$$\hat{\theta} \leftarrow \hat{\theta} - \alpha \nabla_{\theta} \phi$$

Iteration reads as:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \mathbb{E}_{\tau_{\pi}} [(V_{\pi} - V_{\theta}) \nabla_{\theta} V_{\theta}]$$

**Stochastic gradient** approach:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha (V_{\pi} - V_{\theta}) \nabla_{\theta} V_{\theta}$$

**TD approximation**

$$V_{\pi}(x) \approx L(x, \pi(x)) + \gamma V_{\theta}(x_+)$$

is used.

Stochastic TD approach - Solving fitting  $\min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_{\pi}} \left[ (V_{\pi} - V_{\theta})^2 \right]$

**Gradient** approach to solve

$$\hat{\theta} = \arg \min_{\theta} \underbrace{\frac{1}{2} \mathbb{E}_{\tau_{\pi}} \left[ (V_{\pi} - V_{\theta})^2 \right]}_{:=\phi(\theta)}$$

via gradient steps:

$$\hat{\theta} \leftarrow \hat{\theta} - \alpha \nabla_{\theta} \phi$$

Iteration reads as:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \mathbb{E}_{\tau_{\pi}} [(V_{\pi} - V_{\theta}) \nabla_{\theta} V_{\theta}]$$

**Stochastic gradient** approach:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha (V_{\pi} - V_{\theta}) \nabla_{\theta} V_{\theta}$$

**TD approximation**

$$V_{\pi}(x) \approx L(x, \pi(x)) + \gamma V_{\theta}(x_+)$$

is used. Yields:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha (L(x, \pi(x)) + \gamma V_{\theta}(x_+) - V_{\theta}(x)) \nabla_{\theta} V_{\theta}$$

Stochastic TD approach - Solving fitting  $\min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi} [(V_\pi - V_\theta)^2]$

**Gradient** approach to solve

$$\hat{\theta} = \arg \min_{\theta} \underbrace{\frac{1}{2} \mathbb{E}_{\tau_\pi} [(V_\pi - V_\theta)^2]}_{:=\phi(\theta)}$$

via gradient steps:

$$\hat{\theta} \leftarrow \hat{\theta} - \alpha \nabla_{\theta} \phi$$

Iteration reads as:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \mathbb{E}_{\tau_\pi} [(V_\pi - V_\theta) \nabla_{\theta} V_\theta]$$

**TD approximation**

$$V_\pi(x) \approx L(x, \pi(x)) + \gamma V_\theta(x_+)$$

is used. Yields:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha (L(x, \pi(x)) + \gamma V_\theta(x_+) - V_\theta(x)) \nabla_{\theta} V_\theta$$

or equivalently

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta \cdot \nabla_{\theta} V_\theta$$

$$\delta = L(x, \pi(x)) + \gamma V_\theta(x_+) - V_\theta(x)$$

**Stochastic gradient** approach:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha (V_\pi - V_\theta) \nabla_{\theta} V_\theta$$

## Stochastic TD approach - Solving fitting $\min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi} [(V_\pi - V_\theta)^2]$

**Gradient** approach to solve

$$\hat{\theta} = \arg \min_{\theta} \underbrace{\frac{1}{2} \mathbb{E}_{\tau_\pi} [(V_\pi - V_\theta)^2]}_{:=\phi(\theta)}$$

via gradient steps:

$$\hat{\theta} \leftarrow \hat{\theta} - \alpha \nabla_{\theta} \phi$$

Iteration reads as:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \mathbb{E}_{\tau_\pi} [(V_\pi - V_\theta) \nabla_{\theta} V_\theta]$$

**Stochastic gradient** approach:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha (V_\pi - V_\theta) \nabla_{\theta} V_\theta$$

**TD approximation**

$$V_\pi(x) \approx L(x, \pi(x)) + \gamma V_\theta(x_+)$$

is used. Yields:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha (L(x, \pi(x)) + \gamma V_\theta(x_+) - V_\theta(x)) \nabla_{\theta} V_\theta$$

or equivalently

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta \cdot \nabla_{\theta} V_\theta$$

$$\delta = L(x, \pi(x)) + \gamma V_\theta(x_+) - V_\theta(x)$$

note that  $\delta \nabla_{\theta} V_\theta$  is not the gradient of anything...

Stochastic TD approach - Solving fitting  $\min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi} \left[ (V_\pi - V_\theta)^2 \right]$

**Gradient** approach to solve

$$\hat{\theta} = \arg \min_{\theta} \underbrace{\frac{1}{2} \mathbb{E}_{\tau_\pi} \left[ (V_\pi - V_\theta)^2 \right]}_{:=\phi(\theta)}$$

via gradient steps:

$$\hat{\theta} \leftarrow \hat{\theta} - \alpha \nabla_{\theta} \phi$$

Iteration reads as:

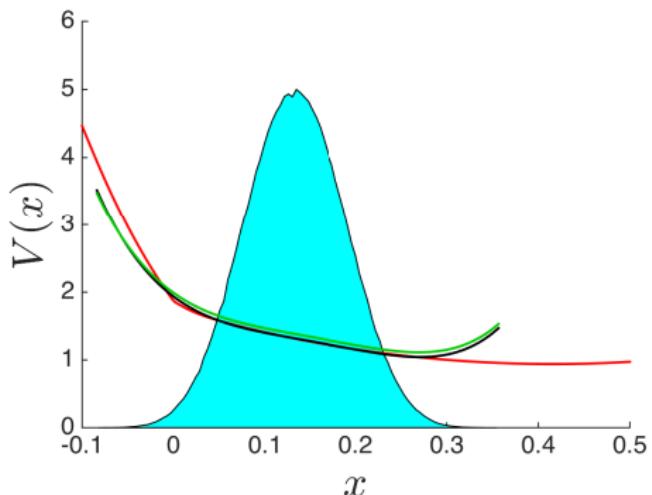
$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \mathbb{E}_{\tau_\pi} [(V_\pi - V_\theta) \nabla_{\theta} V_\theta]$$

**Stochastic gradient** approach:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha (V_\pi - V_\theta) \nabla_{\theta} V_\theta$$

$V_\pi(x)$ ,  $\hat{V}(x)$ ,  $V_{\text{TD}}(x)$

$$\begin{aligned}\hat{\theta} &\leftarrow \hat{\theta} + \alpha \delta \cdot \nabla_{\theta} V_\theta \\ \delta &= L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_\theta(\mathbf{x}_+) - V_\theta(\mathbf{x})\end{aligned}$$



Stochastic TD approach - Solving fitting  $\min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi} [(V_\pi - V_\theta)^2]$

**Gradient** approach to solve

$$\hat{\theta} = \arg \min_{\theta} \underbrace{\frac{1}{2} \mathbb{E}_{\tau_\pi} [(V_\pi - V_\theta)^2]}_{:=\phi(\theta)}$$

$$\begin{aligned}\hat{\theta} &\leftarrow \hat{\theta} + \alpha \delta \cdot \nabla_{\theta} V_{\theta} \\ \delta &= L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})\end{aligned}$$

via gradient steps:

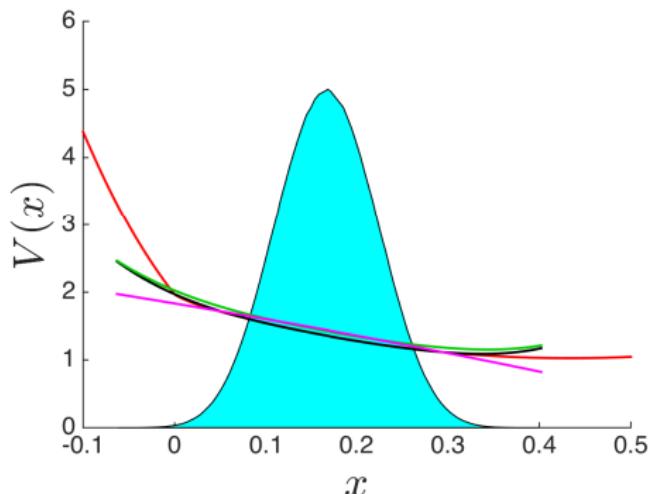
$$\hat{\theta} \leftarrow \hat{\theta} - \alpha \nabla_{\theta} \phi$$

Iteration reads as:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \mathbb{E}_{\tau_\pi} [(V_\pi - V_\theta) \nabla_{\theta} V_{\theta}]$$

**Stochastic gradient** approach:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha (V_\pi - V_\theta) \nabla_{\theta} V_{\theta}$$



$V_\pi(\mathbf{x})$ ,  $\hat{V}(\mathbf{x})$ ,  $V_{\text{TD}}(\mathbf{x})$ ,  $V_{\text{Gradient}}(\mathbf{x})$

## TD( $\lambda$ ) approach

**Stochastic gradient with TD:**

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta \nabla_{\theta} V_{\theta}$$

$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

**Stochastic gradient with TD( $\lambda$ ):**

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta E$$

$$E \leftarrow \gamma \lambda E + \nabla_{\theta} V_{\theta}$$

$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

## TD( $\lambda$ ) approach

**Stochastic gradient with TD:**

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta \nabla_{\theta} V_{\theta}$$

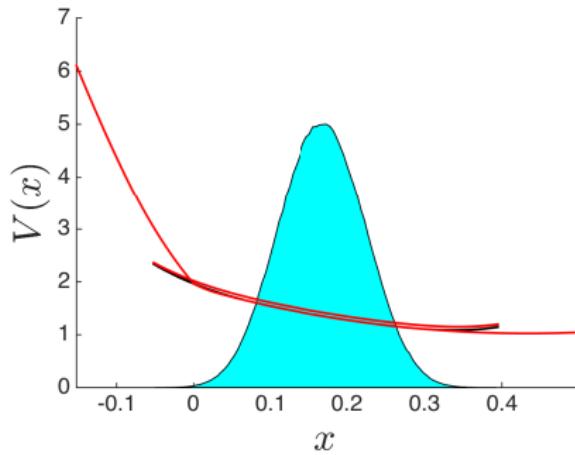
$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

**Stochastic gradient with TD( $\lambda$ ):**

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta E$$

$$E \leftarrow \gamma \lambda E + \nabla_{\theta} V_{\theta}$$

$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$



# Outline

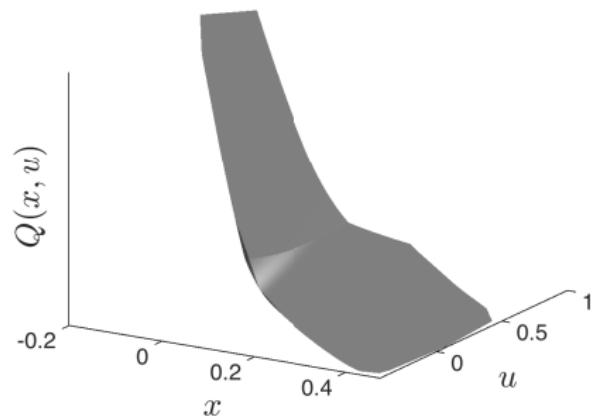
- 1 Introduction
- 2 Value function fitting
- 3 Gradient-based approaches
- 4 Action-value function fitting
- 5 The “LS” family (LSTD, LSTDQ, and LSPI)
- 6 Some observations on approximate Q-learning

## Function fitting for $Q_\theta$ - Temporal Difference

**Value function fitting:**

$$\hat{\theta} = \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi^d} \left[ (Q_\pi(x, u) - Q_\theta(x, u))^2 \right]$$

Where  $\tau_\pi^d$  is the Markov chain distribution under disturbed policy  $\pi + d$



## Function fitting for $Q_\theta$ - Temporal Difference

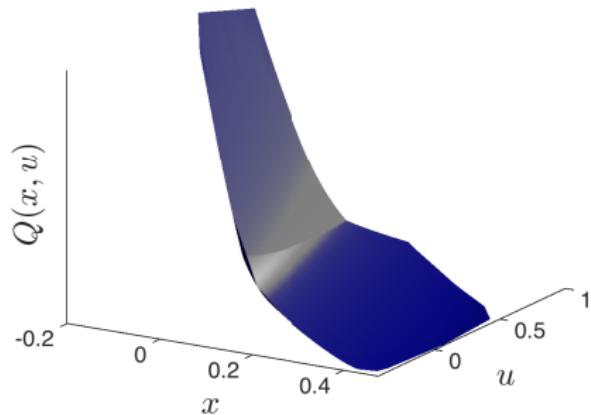
**Value function fitting:**

$$\hat{\theta} = \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi^d} \left[ (Q_\pi(x, u) - Q_\theta(x, u))^2 \right]$$

**Optimal parameter  $\hat{\theta}$  satisfies:**

$$\mathbb{E}_{\tau_\pi^d} [(Q_\pi - Q_\theta) \nabla_\theta Q_\theta] = 0$$

Where  $\tau_\pi^d$  is the Markov chain distribution under disturbed policy  $\pi + d$



## Function fitting for $Q_\theta$ - Temporal Difference

**Value function fitting:**

$$\hat{\theta} = \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi^d} \left[ (Q_\pi(\mathbf{x}, \mathbf{u}) - Q_\theta(\mathbf{x}, \mathbf{u}))^2 \right]$$

Optimal parameter  $\hat{\theta}$  satisfies:

$$\mathbb{E}_{\tau_\pi^d} [(Q_\pi - Q_\theta) \nabla_\theta Q_\theta] = 0$$

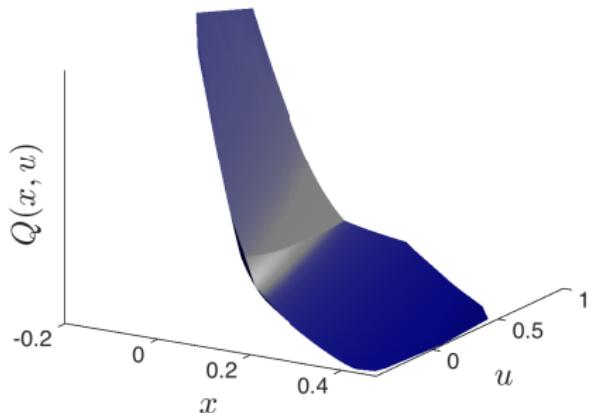
Where do we get  $Q_\pi$  from? Recall that

$$Q_\pi(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[Q_\pi(\mathbf{x}_+, \pi(\mathbf{x}_+)) | \mathbf{x}, \mathbf{u}]$$

### TD evaluation

$$\begin{aligned} Q_\pi(\mathbf{x}, \mathbf{u}) &\approx L(\mathbf{x}, \mathbf{u}) + \gamma Q_\pi(\mathbf{x}_+, \pi(\mathbf{x})) \\ &\approx L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi(\mathbf{x})) \end{aligned}$$

Where  $\tau_\pi^d$  is the Markov chain distribution under disturbed policy  $\pi + d$



## Function fitting for $Q_\theta$ - Temporal Difference

**Value function fitting:**

$$\hat{\theta} = \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi^d} \left[ (Q_\pi(\mathbf{x}, \mathbf{u}) - Q_\theta(\mathbf{x}, \mathbf{u}))^2 \right]$$

**Optimal parameter  $\hat{\theta}$  satisfies:**

$$\mathbb{E}_{\tau_\pi^d} [(Q_\pi - Q_\theta) \nabla_\theta Q_\theta] = 0$$

Using

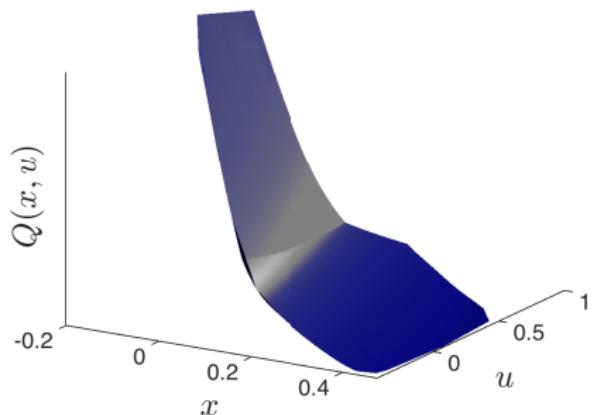
$$Q_\pi(\mathbf{x}, \mathbf{u}) \approx L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi(\mathbf{x}_+))$$

**Compute parameter  $\hat{\theta}$  using:**

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi(\mathbf{x}_+)) - Q_\theta(\mathbf{x}, \mathbf{u})$$

Where  $\tau_\pi^d$  is the Markov chain distribution under disturbed policy  $\pi + \mathbf{d}$



## Function fitting for $Q_\theta$ - Temporal Difference

**Value function fitting:**

$$\hat{\theta} = \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi^d} \left[ (Q_\pi(\mathbf{x}, \mathbf{u}) - Q_\theta(\mathbf{x}, \mathbf{u}))^2 \right]$$

**Optimal parameter  $\hat{\theta}$  satisfies:**

$$\mathbb{E}_{\tau_\pi^d} [(Q_\pi - Q_\theta) \nabla_\theta Q_\theta] = 0$$

Using

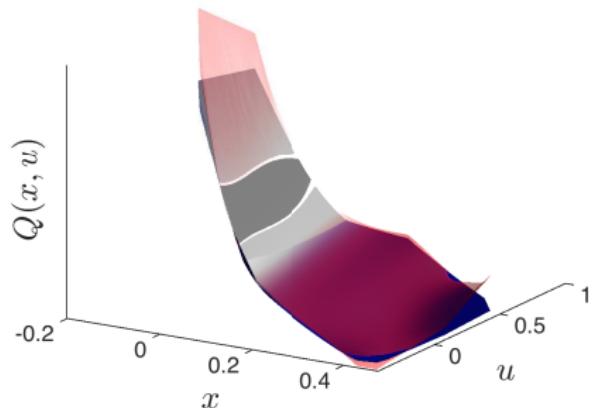
$$Q_\pi(\mathbf{x}, \mathbf{u}) \approx L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi(\mathbf{x}_+))$$

**Compute parameter  $\hat{\theta}$  using:**

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi(\mathbf{x}_+)) - Q_\theta(\mathbf{x}, \mathbf{u})$$

Where  $\tau_\pi^d$  is the Markov chain distribution under disturbed policy  $\pi + d$



## Function fitting for $Q_\theta$ - Temporal Difference

**Value function fitting:**

$$\hat{\theta} = \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi^d} \left[ (Q_\pi(\mathbf{x}, \mathbf{u}) - Q_\theta(\mathbf{x}, \mathbf{u}))^2 \right]$$

**Optimal parameter  $\hat{\theta}$  satisfies:**

$$\mathbb{E}_{\tau_\pi^d} [(Q_\pi - Q_\theta) \nabla_\theta Q_\theta] = 0$$

Using

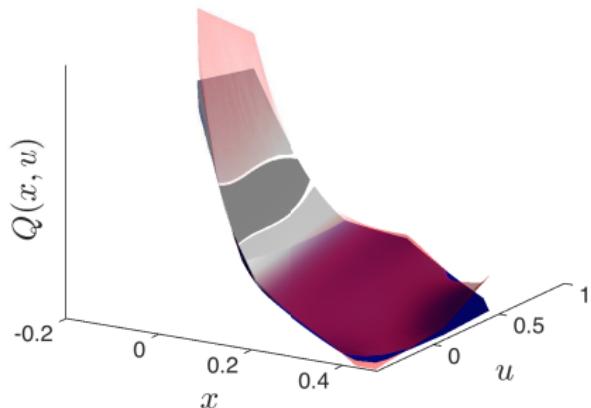
$$Q_\pi(\mathbf{x}, \mathbf{u}) \approx L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi(\mathbf{x}_+))$$

**Compute parameter  $\hat{\theta}$  using:**

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi(\mathbf{x}_+)) - Q_\theta(\mathbf{x}, \mathbf{u})$$

Where  $\tau_\pi^d$  is the Markov chain distribution under disturbed policy  $\pi + \mathbf{d}$



**Stochastic gradient:**

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_\theta Q_\theta$$

## Function fitting for $Q_\theta$ - Temporal Difference

**Value function fitting:**

$$\hat{\theta} = \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi^d} \left[ (Q_\pi(\mathbf{x}, \mathbf{u}) - Q_\theta(\mathbf{x}, \mathbf{u}))^2 \right]$$

**Optimal parameter  $\hat{\theta}$  satisfies:**

$$\mathbb{E}_{\tau_\pi^d} [(Q_\pi - Q_\theta) \nabla_\theta Q_\theta] = 0$$

Using

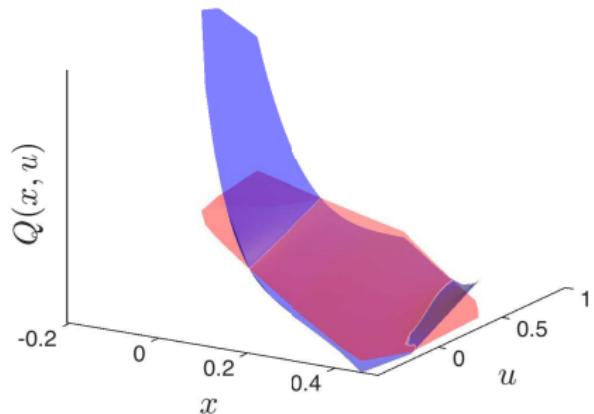
$$Q_\pi(\mathbf{x}, \mathbf{u}) \approx L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi(\mathbf{x}_+))$$

**Compute parameter  $\hat{\theta}$  using:**

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi(\mathbf{x}_+)) - Q_\theta(\mathbf{x}, \mathbf{u})$$

Where  $\tau_\pi^d$  is the Markov chain distribution under disturbed policy  $\pi + \mathbf{d}$



**Stochastic gradient:**

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_\theta Q_\theta$$

# Off/On-policy Q-Learning - SARSA( $\lambda$ )

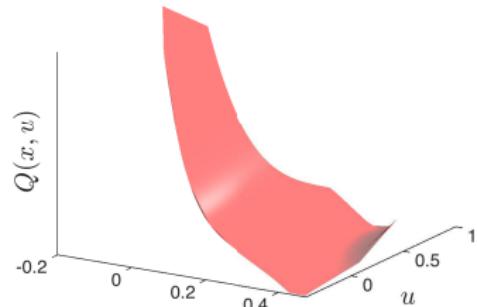
## Off-policy TD error

$$\delta_Q^{\text{off}} = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

## Stochastic gradient

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q^{\text{off}} \nabla_{\theta} Q_{\theta}$$

where  $\mathbf{u}_k$  can be arbitrary



$$\mathbb{E}_{\tau_{\pi}^{\text{d}}} \left[ \delta_Q^{\text{on/off}} \nabla_{\theta} Q_{\theta} \right] = 0$$

# Off/On-policy Q-Learning - SARSA( $\lambda$ )

## Off-policy TD error

$$\delta_Q^{\text{off}} = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

## Stochastic gradient

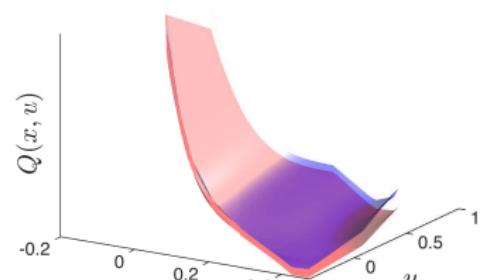
$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q^{\text{off}} \nabla_{\theta} Q_{\theta}$$

where  $\mathbf{u}_k$  can be arbitrary

## On-policy TD error

$$\delta_Q^{\text{on}} = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

where  $\mathbf{u}_k$  is drawn from  $\pi(\mathbf{x}_k) + \text{disturbance}$



$$\mathbb{E}_{\tau_{\pi}^{\text{d}}} [\delta_Q^{\text{on/off}} \nabla_{\theta} Q_{\theta}] = 0$$

# Off/On-policy Q-Learning - SARSA( $\lambda$ )

## Off-policy TD error

$$\delta_Q^{\text{off}} = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

## Stochastic gradient

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q^{\text{off}} \nabla_{\theta} Q_{\theta}$$

where  $\mathbf{u}_k$  can be arbitrary

## On-policy TD error

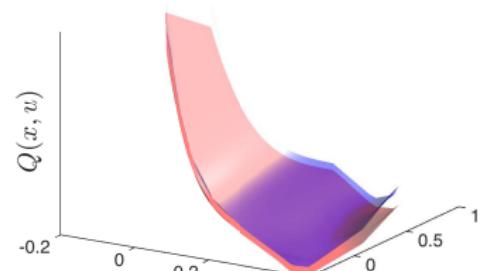
$$\delta_Q^{\text{on}} = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

## SARSA( $\lambda$ ):

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q^{\text{on}} E$$

$$E \leftarrow \gamma \lambda E + \nabla_{\theta} Q_{\theta}$$

where  $\mathbf{u}_k$  is drawn from  $\pi(\mathbf{x}_k) + \text{disturbance}$



$$\mathbb{E}_{\tau_{\pi}^{\text{d}}} [\delta_Q^{\text{on/off}} \nabla_{\theta} Q_{\theta}] = 0$$

# Off/On-policy Q-Learning - SARSA( $\lambda$ )

## Off-policy TD error

$$\delta_Q^{\text{off}} = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

## Stochastic gradient

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q^{\text{off}} \nabla_{\theta} Q_{\theta}$$

where  $\mathbf{u}_k$  can be arbitrary

## On-policy TD error

$$\delta_Q^{\text{on}} = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

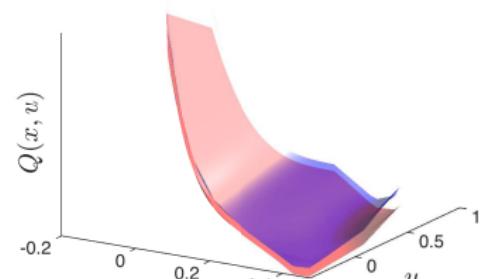
## SARSA( $\lambda$ ):

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q^{\text{on}} E$$

$$E \leftarrow \gamma \lambda E + \nabla_{\theta} Q_{\theta}$$

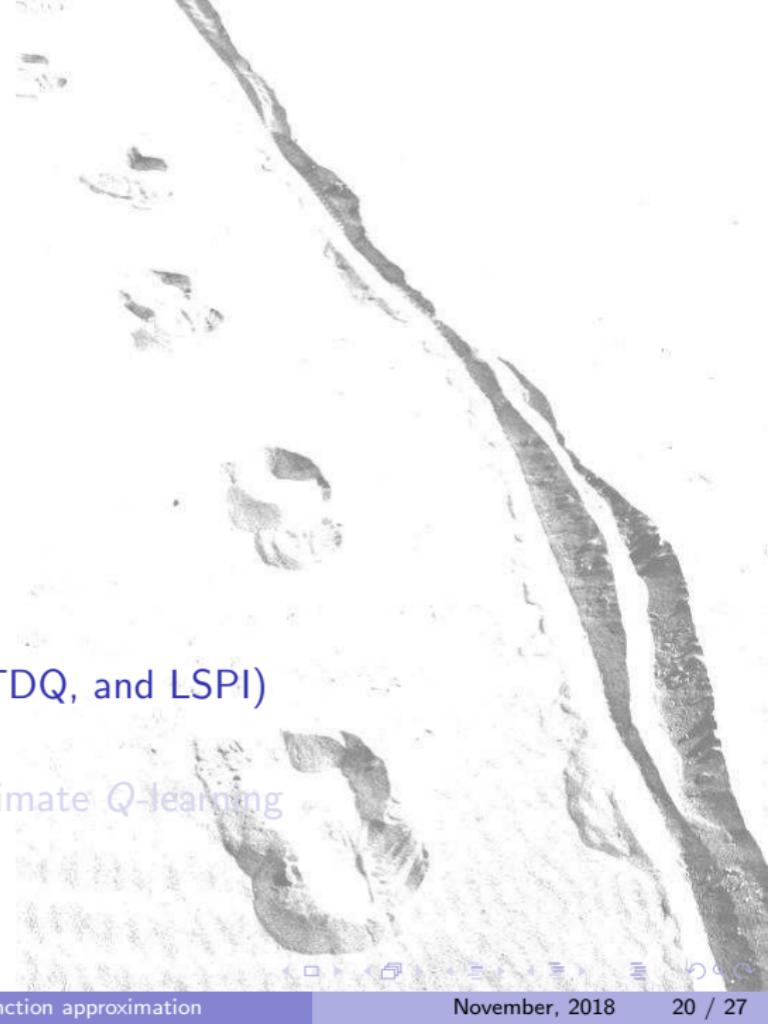
where  $\mathbf{u}_k$  is drawn from  $\pi(\mathbf{x}_k) + \text{disturbance}$

**SARSA( $\lambda$ )** gets  $Q_{\theta} \approx Q_{\pi}$  if disturbance vanishes over time, while exploring the state-input state " $\infty$ "-many times



$$\mathbb{E}_{\tau_{\pi}^{\text{d}}} \left[ \delta_Q^{\text{on/off}} \nabla_{\theta} Q_{\theta} \right] = 0$$

# Outline

- 
- 1 Introduction
  - 2 Value function fitting
  - 3 Gradient-based approaches
  - 4 Action-value function fitting
  - 5 The “LS” family (LSTD, LSTDQ, and LSPI)
  - 6 Some observations on approximate Q-learning

## LSTD - Efficient evaluation of $V$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_\theta V_\theta] = 0$$

$$\delta = L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_\theta(\mathbf{x}_+) - V_\theta(\mathbf{x})$$

for  $\theta$

**Stochastic gradient:**

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta \nabla_\theta V_\theta$$

gets us there (or close enough...)

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- – inefficient (many many iterations)

## LSTD - Efficient evaluation of $V$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_\theta V_\theta] = 0$$

$$\delta = L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_\theta(\mathbf{x}_+) - V_\theta(\mathbf{x})$$

for  $\theta$

Given trajectory, **LSTD** solves:

$$\frac{1}{N} \sum_{k=0}^{N-1} \delta_k(\mathbf{x}_k, \mathbf{x}_{k+1}) \nabla_\theta V_\theta(\mathbf{x}_k) = 0$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma V_\theta(\mathbf{x}_{k+1}) - V_\theta(\mathbf{x}_k)$$

for observed  $\mathbf{x}_0, \dots, \mathbf{x}_N$

**Stochastic gradient:**

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta \nabla_\theta V_\theta$$

gets us there (or close enough...)

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- – inefficient (many many iterations)

## LSTD - Efficient evaluation of $V$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_\theta V_\theta] = 0$$

$$\delta = L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_\theta(\mathbf{x}_+) - V_\theta(\mathbf{x})$$

for  $\theta$

Given trajectory, **LSTD** solves:

$$\frac{1}{N} \sum_{k=0}^{N-1} \delta_k(\mathbf{x}_k, \mathbf{x}_{k+1}) \nabla_\theta V_\theta(\mathbf{x}_k) = 0$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma V_\theta(\mathbf{x}_{k+1}) - V_\theta(\mathbf{x}_k)$$

for observed  $\mathbf{x}_{0,\dots,N}$

**Stochastic gradient:**

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta \nabla_\theta V_\theta$$

gets us there (or close enough...)

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- – inefficient (many many iterations)

- If  $V_\theta$  is linear in  $\theta$  then solution is explicit, global solution
- If nonlinear, Newton required
- Nonlinear  $\rightarrow$  possibly local solution / initial guess matters

## LSTD - Efficient evaluation of $V$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_\theta V_\theta] = 0$$

$$\delta = L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_\theta(\mathbf{x}_+) - V_\theta(\mathbf{x})$$

for  $\theta$

Given trajectory, **LSTD** solves:

$$\frac{1}{N} \sum_{k=0}^{N-1} \delta_k(\mathbf{x}_k, \mathbf{x}_{k+1}) \nabla_\theta V_\theta(\mathbf{x}_k) = 0$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma V_\theta(\mathbf{x}_{k+1}) - V_\theta(\mathbf{x}_k)$$

for observed  $\mathbf{x}_0, \dots, N$

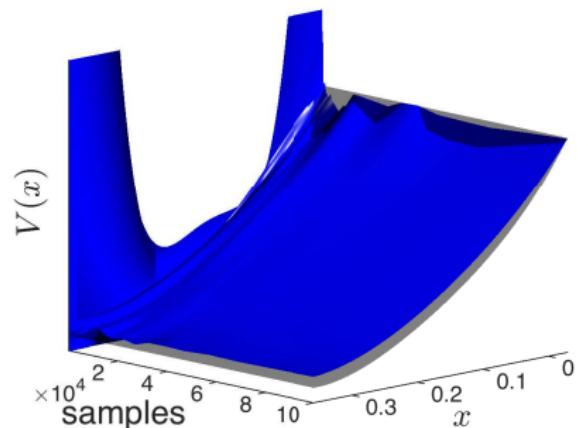
- If  $V_\theta$  is linear in  $\theta$  then solution is explicit, global solution
- If nonlinear, Newton required
- Nonlinear  $\rightarrow$  possibly local solution / initial guess matters

## Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta \nabla_\theta V_\theta$$

gets us there (or close enough...)

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)



Recall that we want to solve

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_{\theta} Q_{\theta}] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_{\theta}(\mathbf{x}_+, \pi_+) - Q_{\theta}(\mathbf{x}, \mathbf{u})$$

for  $\theta$

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_{\theta} Q_{\theta}$$

gets us there.

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)

## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \boldsymbol{\pi}_+) - Q_\theta(\mathbf{x}, \mathbf{u})$$

for  $\theta$

### Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_\theta Q_\theta$$

gets us there.

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)

**LSTDQ** solves:

$$\sum_{k=0}^{N-1} \delta_k \nabla_\theta Q_\theta(\mathbf{x}_k, \mathbf{u}_k) = 0$$

$$\delta_k = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_\theta(\mathbf{x}_{k+1}, \boldsymbol{\pi}(\mathbf{x}_{k+1})) - Q_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

for observed trajectory  $(\mathbf{x}_0, \mathbf{u}_0), \dots, (\mathbf{x}_N, \mathbf{u}_N)$

- If  $Q_\theta$  is linear in  $\theta$  then solution is explicit, global solution
- If nonlinear, Newton required
- Nonlinear  $\rightarrow$  possibly local solution / initial guess matters

## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_{\theta} Q_{\theta}] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_{\theta}(\mathbf{x}_+, \pi_+) - Q_{\theta}(\mathbf{x}, \mathbf{u})$$

for  $\theta$

### Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_{\theta} Q_{\theta}$$

gets us there.

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)

---

### Algorithm: LSTDQ (prototype)

---

**Input:** Initial  $Q$  parameters  $\theta$

Collect  $\mathcal{D} = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 0, \dots, N\}$

**while**  $\|\zeta_{\theta}\| > \text{Tol}$  **do**

Form

$$\zeta_{\theta} = \sum_{\mathcal{D}} \delta_k \nabla_{\theta} Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

and gradient

$$\text{Newton step: } \theta \leftarrow \theta + \alpha \nabla_{\theta} \zeta_{\theta}^{-1} \zeta_{\theta}$$

---

## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi_+) - Q_\theta(\mathbf{x}, \mathbf{u})$$

for  $\theta$

### Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_\theta Q_\theta$$

gets us there.

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)

---

### Algorithm: LSTDQ (prototype)

---

**Input:** Initial  $Q$  parameters  $\theta$

Collect  $\mathcal{D} = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 0, \dots, N\}$

**while**  $\|\zeta_\theta\| > \text{Tol}$  **do**

Form

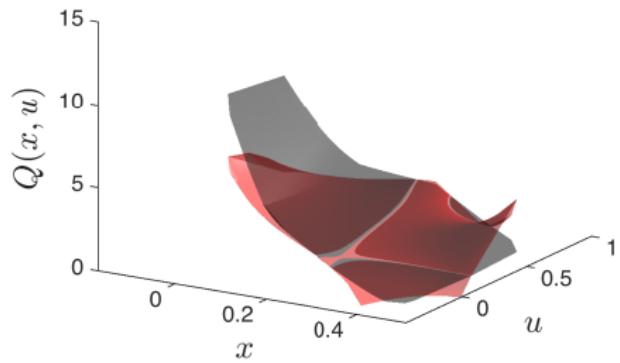
$$\zeta_\theta = \sum_{\mathcal{D}} \delta_k \nabla_\theta Q_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

and gradient

$$\text{Newton step: } \theta \leftarrow \theta + \alpha \nabla_\theta \zeta_\theta^{-1} \zeta_\theta$$

---

# samples = 10



## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi_+) - Q_\theta(\mathbf{x}, \mathbf{u})$$

for  $\theta$

### Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_\theta Q_\theta$$

gets us there.

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)

---

### Algorithm: LSTDQ (prototype)

---

**Input:** Initial  $Q$  parameters  $\theta$

Collect  $\mathcal{D} = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 0, \dots, N\}$

**while**  $\|\zeta_\theta\| > \text{Tol}$  **do**

Form

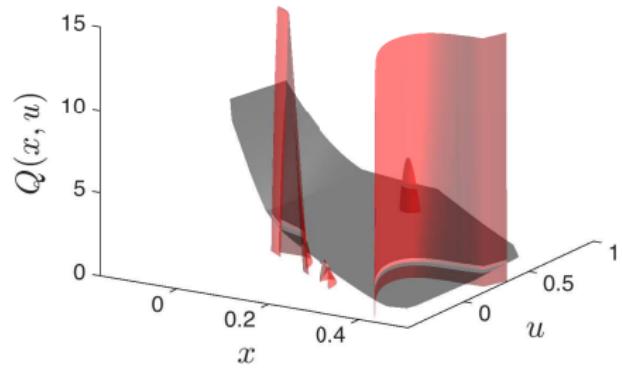
$$\zeta_\theta = \sum_{\mathcal{D}} \delta_k \nabla_\theta Q_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

and gradient

$$\text{Newton step: } \theta \leftarrow \theta + \alpha \nabla_\theta \zeta_\theta^{-1} \zeta_\theta$$

---

# samples = 35



## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi_+) - Q_\theta(\mathbf{x}, \mathbf{u})$$

for  $\theta$

**Stochastic gradient:**

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_\theta Q_\theta$$

gets us there.

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)

---

**Algorithm:** LSTDQ (prototype)

---

**Input:** Initial  $Q$  parameters  $\theta$

Collect  $\mathcal{D} = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 0, \dots, N\}$

**while**  $\|\zeta_\theta\| > \text{Tol}$  **do**

Form

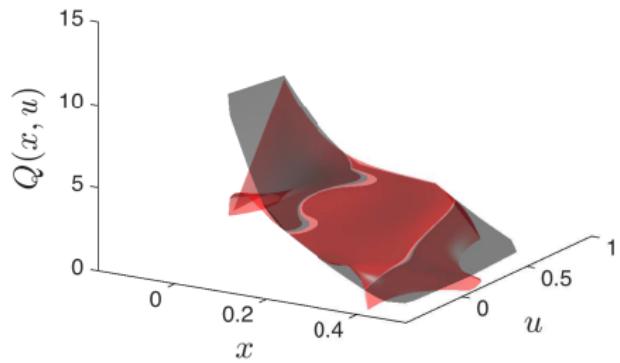
$$\zeta_\theta = \sum_{\mathcal{D}} \delta_k \nabla_\theta Q_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

and gradient

$$\text{Newton step: } \theta \leftarrow \theta + \alpha \nabla_\theta \zeta_\theta^{-1} \zeta_\theta$$

---

# samples = 129



## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi_+) - Q_\theta(\mathbf{x}, \mathbf{u})$$

for  $\theta$

### Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_\theta Q_\theta$$

gets us there.

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)

---

### Algorithm: LSTDQ (prototype)

---

**Input:** Initial  $Q$  parameters  $\theta$

Collect  $\mathcal{D} = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 0, \dots, N\}$

**while**  $\|\zeta_\theta\| > \text{Tol}$  **do**

Form

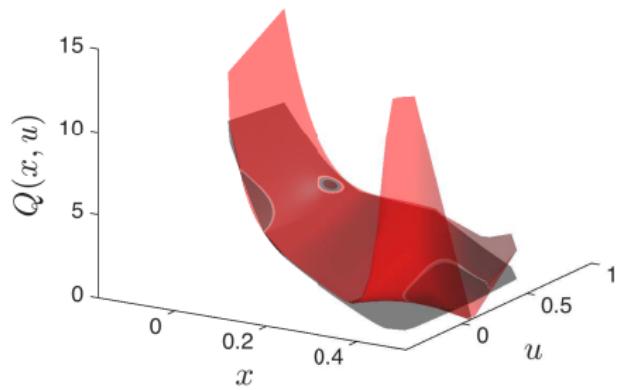
$$\zeta_\theta = \sum_{\mathcal{D}} \delta_k \nabla_\theta Q_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

and gradient

$$\text{Newton step: } \theta \leftarrow \theta + \alpha \nabla_\theta \zeta_\theta^{-1} \zeta_\theta$$

---

# samples = 464



## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi_+) - Q_\theta(\mathbf{x}, \mathbf{u})$$

for  $\theta$

### Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_\theta Q_\theta$$

gets us there.

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)

---

### Algorithm: LSTDQ (prototype)

---

**Input:** Initial  $Q$  parameters  $\theta$

Collect  $\mathcal{D} = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 0, \dots, N\}$

**while**  $\|\zeta_\theta\| > \text{Tol}$  **do**

Form

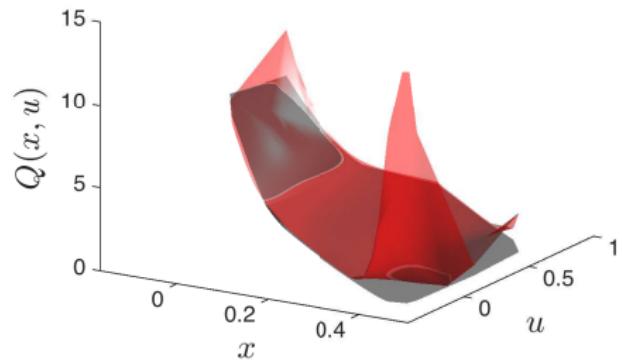
$$\zeta_\theta = \sum_{\mathcal{D}} \delta_k \nabla_\theta Q_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

and gradient

$$\text{Newton step: } \theta \leftarrow \theta + \alpha \nabla_\theta \zeta_\theta^{-1} \zeta_\theta$$

---

# samples = 1668



## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_{\underline{\pi}}^d} [\delta \nabla_{\theta} Q_{\theta}] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi_+) - Q_\theta(\mathbf{x}, \mathbf{u}) \quad \text{for } \theta$$

### Algorithm: LSTDQ (prototype)

**Input:** Initial  $Q$  parameters  $\theta$

Collect  $\mathcal{D} = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 0, \dots, N\}$

**while**  $\|\zeta_\theta\| > \text{Tol}$  **do**

## Form

$$\zeta_{\theta} = \sum_{\mathcal{P}} \delta_k \nabla_{\theta} Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

and gradient

Newton step:  $\theta \leftarrow \theta + \alpha \nabla_{\theta} \zeta_{\theta}^{-1} \zeta_{\theta}$

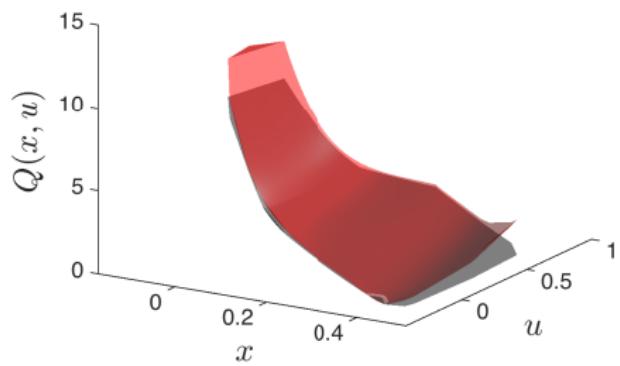
### Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_{\theta} Q_{\theta}$$

gets us there.

- + trivial to code
  - + computationally cheap (for  $\theta$  large)
  - – inefficient (many many iterations)

# samples = 5994



## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi_+) - Q_\theta(\mathbf{x}, \mathbf{u})$$

for  $\theta$

### Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_\theta Q_\theta$$

gets us there.

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)

---

### Algorithm: LSTDQ (prototype)

---

**Input:** Initial  $Q$  parameters  $\theta$

Collect  $\mathcal{D} = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 0, \dots, N\}$

**while**  $\|\zeta_\theta\| > \text{Tol}$  **do**

Form

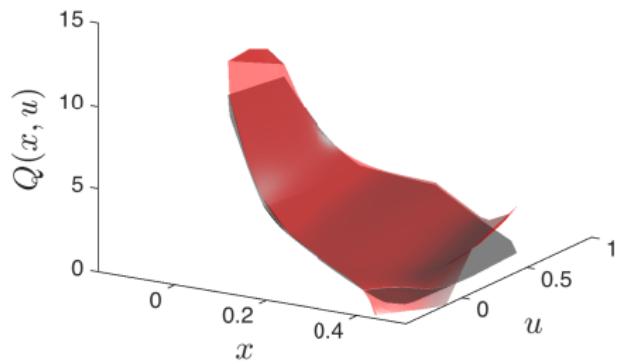
$$\zeta_\theta = \sum_{\mathcal{D}} \delta_k \nabla_\theta Q_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

and gradient

$$\text{Newton step: } \theta \leftarrow \theta + \alpha \nabla_\theta \zeta_\theta^{-1} \zeta_\theta$$

---

# samples = 21544



## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_{\underline{\pi}}^d} [\delta \nabla_{\theta} Q_{\theta}] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi_+) - Q_\theta(\mathbf{x}, \mathbf{u}) \quad \text{for } \theta$$

**Algorithm:** LSTDQ (prototype)

**Input:** Initial  $Q$  parameters  $\theta$

Collect  $\mathcal{D} = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 0, \dots, N\}$

**while**  $\|\zeta_\theta\| > \text{Tol}$  **do**

## Form

$$\zeta_{\theta} = \sum_{\mathcal{D}} \delta_k \nabla_{\theta} Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

and gradient

Newton step:  $\theta \leftarrow \theta + \alpha \nabla_{\theta} \zeta_{\theta}^{-1} \zeta_{\theta}$

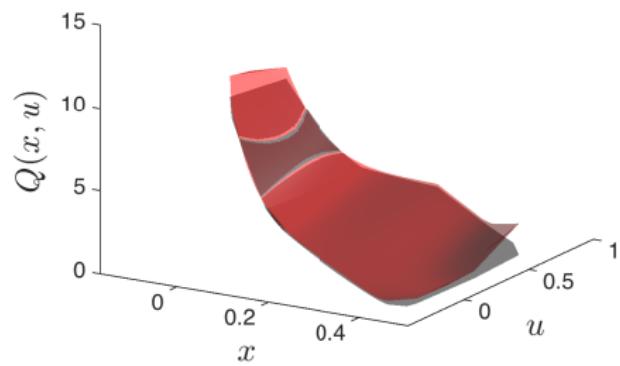
## Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_{\theta} Q_{\theta}$$

gets us there.

- + trivial to code
  - + computationally cheap (for  $\theta$  large)
  - – inefficient (many many iterations)

# samples = 77426



## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi_+) - Q_\theta(\mathbf{x}, \mathbf{u})$$

for  $\theta$

### Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_\theta Q_\theta$$

gets us there.

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)

---

### Algorithm: LSTDQ (prototype)

---

**Input:** Initial  $Q$  parameters  $\theta$

Collect  $\mathcal{D} = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 0, \dots, N\}$

**while**  $\|\zeta_\theta\| > \text{Tol}$  **do**

Form

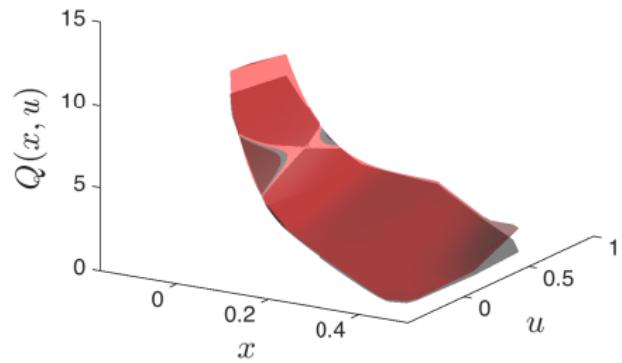
$$\zeta_\theta = \sum_{\mathcal{D}} \delta_k \nabla_\theta Q_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

and gradient

$$\text{Newton step: } \theta \leftarrow \theta + \alpha \nabla_\theta \zeta_\theta^{-1} \zeta_\theta$$

---

# samples = 278255



## LSTDQ - Efficient evaluation of $Q_\pi$

Recall that we want to solve

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi_+) - Q_\theta(\mathbf{x}, \mathbf{u})$$

for  $\theta$

### Stochastic gradient:

$$\hat{\theta} \leftarrow \hat{\theta} + \alpha \delta_Q \nabla_\theta Q_\theta$$

gets us there.

- + trivial to code
- + computationally cheap (for  $\theta$  large)
- - inefficient (many many iterations)

---

### Algorithm: LSTDQ (prototype)

---

**Input:** Initial  $Q$  parameters  $\theta$

Collect  $\mathcal{D} = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 0, \dots, N\}$

**while**  $\|\zeta_\theta\| > \text{Tol}$  **do**

Form

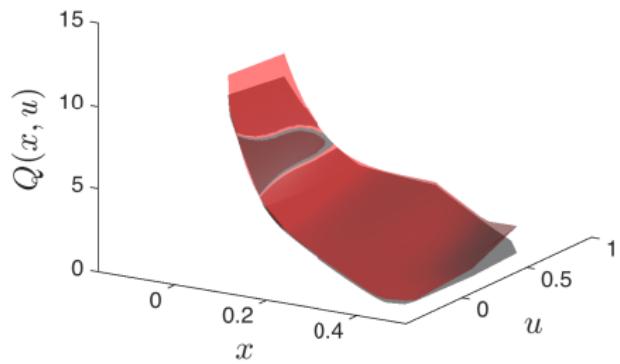
$$\zeta_\theta = \sum_{\mathcal{D}} \delta_k \nabla_\theta Q_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

and gradient

$$\text{Newton step: } \theta \leftarrow \theta + \alpha \nabla_\theta \zeta_\theta^{-1} \zeta_\theta$$

---

# samples = 1000000



# LSPI - Find the optimal policy

Recall

$$\delta_k = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

---

## Algorithm: LSPI (prototype)

---

**Input:** Initial  $Q$  parameters  $\theta$

**for**  $i = 1 \rightarrow \infty$  **do**

    Set policy  $\pi(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}, \mathbf{u})$

    Collect

$D = \{(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1}), k = 1, \dots, N\}$

**while**  $\|\zeta_{\theta}\| > \text{Tol}$  **do**

        Form

$$\zeta_{\theta} = \sum_{\mathcal{D}} \delta_k \nabla_{\theta} Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

        and gradient

$$\text{Newton step: } \theta \leftarrow \theta + \alpha \nabla_{\theta} \zeta_{\theta}^{-1} \zeta_{\theta}$$

---

**return**  $Q$  parameters  $\theta$  for  $\mathcal{D}$

---

- ➊ Set policy  $\pi$  from  $Q_{\theta}$
- ➋ Collect data ( $N$  transitions)
- ➌ Update  $Q_{\theta}$  using LSTDQ
- ➍ Back to 1

# LSPI - Find the optimal policy

Recall

$$\delta_k = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

---

**Algorithm:** LSPI on-the-fly

---

**Input:** Initial  $Q$  parameters  $\theta$

**for**  $k = 1 \rightarrow \infty$  **do**

    Set  $\mathbf{u}_k = \arg \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}_k, \mathbf{u})$

    Observe transition  $(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1})$ , build

$\delta_k =$

$L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}_{k+1}, \mathbf{u}) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$

    Form (for some forgetting factor  $\eta < 1$ )

$$\zeta_{\theta} = \sum_{i=0}^k \eta^i \delta_i \nabla_{\theta} Q_{\theta}(\mathbf{x}_i, \mathbf{u}_i)$$

    and gradient

    Newton step:  $\theta \leftarrow \theta + \alpha \nabla_{\theta} \zeta_{\theta}^{-1} \zeta_{\theta}$

---

- 1 Set policy  $\pi$  from  $Q_{\theta}$
- 2 Observe 1 transition, add to data
- 3 Update  $Q_{\theta}$  using LSTDQ (possibly discounted)
- 4 Back to 1

# LSPI - Find the optimal policy

Recall

$$\delta_k = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

---

## Algorithm: LSPI on-the-fly

---

**Input:** Initial  $Q$  parameters  $\theta$

**for**  $k = 1 \rightarrow \infty$  **do**

    Set  $\mathbf{u}_k = \arg \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}_k, \mathbf{u})$

    Observe transition  $(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1})$ , build

$\delta_k =$

$$L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}_{k+1}, \mathbf{u}) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

    Form (for some forgetting factor  $\eta < 1$ )

$$\zeta_{\theta} = \sum_{i=0}^k \eta^i \delta_i \nabla_{\theta} Q_{\theta}(\mathbf{x}_i, \mathbf{u}_i)$$

    and gradient

$$\text{Newton step: } \theta \leftarrow \theta + \alpha \nabla_{\theta} \zeta_{\theta}^{-1} \zeta_{\theta}$$

---

- 1 Set policy  $\pi$  from  $Q_{\theta}$
- 2 Observe 1 transition, add to data
- 3 Update  $Q_{\theta}$  using LSTDQ (possibly discounted)
- 4 Back to 1

Old samples should be progressively ignored, as they are generated from an old MDP distribution  $\tau_{\pi}$  (policy has changed!)

## Experience replay

Consider the (possibly huge) data set (state transitions):

$$\mathcal{D} = \{(\mathbf{x}_1, \mathbf{u}_1 \rightarrow \mathbf{x}_2), \dots (\mathbf{x}_N, \mathbf{u}_N \rightarrow \mathbf{x}_{N+1})\}$$

Full LSTD update

$$\zeta_{\theta} = \sum_{\mathcal{D}} \delta_k \nabla_{\theta} Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

with

$$\delta_k = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

## Experience replay

Consider the (possibly huge) data set (state transitions):

$$\mathcal{D} = \{(\mathbf{x}_1, \mathbf{u}_1 \rightarrow \mathbf{x}_2), \dots (\mathbf{x}_N, \mathbf{u}_N \rightarrow \mathbf{x}_{N+1})\}$$

Full LSTD update

$$\zeta_{\theta} = \sum_{\mathcal{D}} \delta_k \nabla_{\theta} Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

with

$$\delta_k = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

May not be good:

- we need to re-evaluate  $\delta_k$ ,  $Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1}))$  and  $Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$  for all  $k$  whenever we change  $\theta$ , becomes very expensive for  $N$  large

## Experience replay

Consider the (possibly huge) data set (state transitions):

$$\mathcal{D} = \{(\mathbf{x}_1, \mathbf{u}_1 \rightarrow \mathbf{x}_2), \dots (\mathbf{x}_N, \mathbf{u}_N \rightarrow \mathbf{x}_{N+1})\}$$

Full LSTD update

$$\zeta_{\theta} = \sum_{\mathcal{D}} \delta_k \nabla_{\theta} Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

with

$$\delta_k = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

May not be good:

- we need to re-evaluate  $\delta_k$ ,  $Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1}))$  and  $Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$  for all  $k$  whenever we change  $\theta$ , becomes very expensive for  $N$  large

**Experience replay:** form  $\zeta_{\theta}$  on a sub-data set  $\mathcal{D}_S \subset \mathcal{D}$  of limited cardinality

- how to select the samples  $\mathcal{D}_S$ ? In principle random, but should match MDP distribution under new policy

## Experience replay

Consider the (possibly huge) data set (state transitions):

$$\mathcal{D} = \{(\mathbf{x}_1, \mathbf{u}_1 \rightarrow \mathbf{x}_2), \dots (\mathbf{x}_N, \mathbf{u}_N \rightarrow \mathbf{x}_{N+1})\}$$

Full LSTD update

$$\zeta_{\theta} = \sum_{\mathcal{D}} \delta_k \nabla_{\theta} Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

with

$$\delta_k = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

May not be good:

- we need to re-evaluate  $\delta_k$ ,  $Q_{\theta}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1}))$  and  $Q_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$  for all  $k$  whenever we change  $\theta$ , becomes very expensive for  $N$  large

**Experience replay:** form  $\zeta_{\theta}$  on a sub-data set  $\mathcal{D}_S \subset \mathcal{D}$  of limited cardinality

- how to select the samples  $\mathcal{D}_S$ ? In principle random, but should match MDP distribution under new policy
- many tricks to make a “smart” subset selection

# Outline

- 1 Introduction
- 2 Value function fitting
- 3 Gradient-based approaches
- 4 Action-value function fitting
- 5 The “LS” family (LSTD, LSTDQ, and LSPI)
- 6 Some observations on approximate Q-learning

# Is LSTD(Q) right?

## LSTD

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_{\theta} V_{\theta}] = 0$$

$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

## LSTDQ

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_{\theta} Q_{\theta}] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_{\theta}(\mathbf{x}_+, \pi(\mathbf{x})) - Q_{\theta}(\mathbf{x}, \mathbf{u})$$

**Remark:** if there is a  $\theta_*$  such that  $V_{\theta_*} = V_{\pi}$  for all  $\mathbf{x}$  then

# Is LSTD(Q) right?

## LSTD

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_{\theta} V_{\theta}] = 0$$

$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

## LSTDQ

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_{\theta} Q_{\theta}] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_{\theta}(\mathbf{x}_+, \pi(\mathbf{x})) - Q_{\theta}(\mathbf{x}, \mathbf{u})$$

**Remark:** if there is a  $\theta_*$  such that  $V_{\theta_*} = V_\pi$  for all  $\mathbf{x}$  then

$$V_{\theta_*}(\mathbf{x}) = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V_{\theta_*}(\mathbf{x}_+) | \mathbf{x}] \quad \text{holds } \forall \mathbf{x}$$

# Is LSTD(Q) right?

## LSTD

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_{\theta} V_{\theta}] = 0$$

$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

## LSTDQ

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_{\theta} Q_{\theta}] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_{\theta}(\mathbf{x}_+, \pi(\mathbf{x})) - Q_{\theta}(\mathbf{x}, \mathbf{u})$$

**Remark:** if there is a  $\theta_*$  such that  $V_{\theta_*} = V_\pi$  for all  $\mathbf{x}$  then

$$L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V_{\theta_*}(\mathbf{x}_+) | \mathbf{x}] - V_{\theta_*}(\mathbf{x}) = 0 \quad \text{holds } \forall \mathbf{x}$$

# Is LSTD(Q) right?

## LSTD

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_{\theta} V_{\theta}] = 0$$

$$\delta = L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

## LSTDQ

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_{\theta} Q_{\theta}] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_{\theta}(\mathbf{x}_+, \boldsymbol{\pi}(\mathbf{x})) - Q_{\theta}(\mathbf{x}, \mathbf{u})$$

**Remark:** if there is a  $\theta_*$  such that  $V_{\theta_*} = V_\pi$  for all  $\mathbf{x}$  then

$$\mathbb{E} [L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_{\theta_*}(\mathbf{x}_+) - V_{\theta_*}(\mathbf{x}) \mid \mathbf{x}] = 0 \quad \text{holds } \forall \mathbf{x}$$

# Is LSTD(Q) right?

## LSTD

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_{\theta} V_{\theta}] = 0$$

$$\delta = L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

## LSTDQ

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_{\theta} Q_{\theta}] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_{\theta}(\mathbf{x}_+, \boldsymbol{\pi}(\mathbf{x})) - Q_{\theta}(\mathbf{x}, \mathbf{u})$$

**Remark:** if there is a  $\theta_*$  such that  $V_{\theta_*} = V_\pi$  for all  $\mathbf{x}$  then

$$\bar{\delta}(\mathbf{x}) := \mathbb{E} [L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_{\theta_*}(\mathbf{x}_+) - V_{\theta_*}(\mathbf{x}) \mid \mathbf{x}] = 0 \quad \text{holds} \quad \forall \mathbf{x}$$

# Is LSTD(Q) right?

## LSTD

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_{\theta} V_{\theta}] = 0$$

$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

## LSTDQ

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_{\theta} Q_{\theta}] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_{\theta}(\mathbf{x}_+, \pi(\mathbf{x})) - Q_{\theta}(\mathbf{x}, \mathbf{u})$$

**Remark:** if there is a  $\theta_*$  such that  $V_{\theta_*} = V_\pi$  for all  $\mathbf{x}$  then

$$\mathbb{E}_{\tau_\pi} [\bar{\delta} \nabla_{\theta} V_{\theta_*}] = 0 \quad \text{with} \quad \bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) | \mathbf{x}]$$

*how does this relate to LSTD?*

# Is LSTD(Q) right?

## LSTD

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_{\theta} V_{\theta}] = 0$$

$$\delta = L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_{\theta}(\mathbf{x}_+) - V_{\theta}(\mathbf{x})$$

## LSTDQ

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_{\theta} Q_{\theta}] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_{\theta}(\mathbf{x}_+, \boldsymbol{\pi}(\mathbf{x})) - Q_{\theta}(\mathbf{x}, \mathbf{u})$$

**Remark:** if there is a  $\theta_*$  such that  $V_{\theta_*} = V_{\pi}$  for all  $\mathbf{x}$  then

$$\mathbb{E}_{\tau_\pi} [\bar{\delta} \nabla_{\theta} V_{\theta_*}] = 0 \quad \text{with} \quad \bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) | \mathbf{x}]$$

how does this relate to LSTD?

If  $\mathbb{E}[\zeta(\mathbf{x}, \mathbf{x}_+) | \mathbf{x}] = 0$  for all  $\mathbf{x}$ , then  $\mathbb{E}_{\tau_\pi} [\zeta(\mathbf{x}, \mathbf{x}_+) \xi(\mathbf{x})] = 0$  for any function  $\xi$

Note:  $\mathbb{E}[\cdot | \mathbf{x}]$  is an expectation over state transitions  $\mathbf{x}, \boldsymbol{\pi}(\mathbf{x}) \rightarrow \mathbf{x}_+$ , i.e. it matches  $\tau_\pi$

# Is LSTD(Q) right?

## LSTD

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_\theta V_\theta] = 0$$

$$\delta = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma V_\theta(\mathbf{x}_+) - V_\theta(\mathbf{x})$$

## LSTDQ

$$\mathbb{E}_{\tau_\pi^d} [\delta_Q \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \pi(\mathbf{x})) - Q_\theta(\mathbf{x}, \mathbf{u})$$

**Remark:** if there is a  $\theta_*$  such that  $V_{\theta_*} = V_\pi$  for all  $\mathbf{x}$  then

$$\mathbb{E}_{\tau_\pi} [\bar{\delta} \nabla_\theta V_{\theta_*}] = 0 \quad \text{with} \quad \bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) | \mathbf{x}]$$

how does this relate to LSTD?

If  $\mathbb{E}[\zeta(\mathbf{x}, \mathbf{x}_+) | \mathbf{x}] = 0$  for all  $\mathbf{x}$ , then  $\mathbb{E}_{\tau_\pi} [\zeta(\mathbf{x}, \mathbf{x}_+) \xi(\mathbf{x})] = 0$  for any function  $\xi$

Note:  $\mathbb{E}[\cdot | \mathbf{x}]$  is an expectation over state transitions  $\mathbf{x}, \pi(\mathbf{x}) \rightarrow \mathbf{x}_+$ , i.e. it matches  $\tau_\pi$

Consequence:  $\bar{\delta}(\mathbf{x}) = 0$  for all  $\mathbf{x}$  implies  $\mathbb{E}_{\tau_\pi} [\delta \nabla_\theta V_{\theta_*}] = 0$

Hence LSTD holds if  $V_\theta = V_\pi$

Similarly LSTDQ holds if  $Q_\theta = Q_\pi$

## What does not hold?

### LSTD

$$\mathbb{E}_{\tau_\pi} [\delta \nabla_\theta V_\theta] = 0$$

$$\delta = L(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \gamma V_\theta(\mathbf{x}_+) - V_\theta(\mathbf{x})$$

### LSTDQ

$$\mathbb{E}_{\tau_\pi^d} [\delta \nabla_\theta Q_\theta] = 0$$

$$\delta_Q = L(\mathbf{x}, \mathbf{u}) + \gamma Q_\theta(\mathbf{x}_+, \boldsymbol{\pi}(\mathbf{x})) - Q_\theta(\mathbf{x}, \mathbf{u})$$

**LSTD and LSTDQ do not necessarily deliver**

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi^d} [(V_\pi - V_\theta)^2]$$

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \mathbb{E}_{\tau_\pi^d} [(Q_\pi - Q_\theta)^2]$$

if there is no  $\theta$  such that  $V_\pi = V_\theta$  holds

**Bottom line: a “rich” parametrization of the value functions / policy is important!**

*MC-like targets (deeper backup) alleviate this problem...*