

Policy Gradient methods

Sébastien Gros

Cybernetic, NTNU
Elec. Eng., Chalmers

TUM lectures on RL

Outline

- 
- 1 Introduction
 - 2 Deterministic policy gradient
 - 3 Stochastic policy gradient

Outline

1 Introduction

2 Deterministic policy gradient

3 Stochastic policy gradient

Why another technique?

The approximation problem:

We have seen how to do:

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) \approx Q_{*}(\mathbf{x}, \mathbf{u})$$

and defined:

$$\pi_{\theta}(\mathbf{x}) := \arg \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}, \mathbf{u})$$

For a “rich” parametrization of Q_{θ}

$$\pi_{\theta}(\mathbf{x}) = \pi_{*}(\mathbf{x})$$

holds. Otherwise all bets are off

Why another technique?

The approximation problem:

We have seen how to do:

$$Q_\theta(\mathbf{x}, \mathbf{u}) \approx Q_*(\mathbf{x}, \mathbf{u})$$

and defined:

$$\pi_\theta(\mathbf{x}) := \arg \min_{\mathbf{u}} Q_\theta(\mathbf{x}, \mathbf{u})$$

For a “rich” parametrization of Q_θ

$$\pi_\theta(\mathbf{x}) = \pi_*(\mathbf{x})$$

holds. Otherwise all bets are off

The minimization problem:

Minimization step:

$$\pi_\theta(\mathbf{x}) := \arg \min_{\mathbf{u}} Q_\theta(\mathbf{x}, \mathbf{u})$$

is “easy” only if $Q_\theta(\mathbf{x}, \mathbf{u})$ is convex in \mathbf{u} .

However, generic choices of Q_θ are typically not

Why another technique?

The approximation problem:

We have seen how to do:

$$Q_\theta(\mathbf{x}, \mathbf{u}) \approx Q_*(\mathbf{x}, \mathbf{u})$$

and defined:

$$\pi_\theta(\mathbf{x}) := \arg \min_{\mathbf{u}} Q_\theta(\mathbf{x}, \mathbf{u})$$

For a “rich” parametrization of Q_θ

$$\pi_\theta(\mathbf{x}) = \pi_*(\mathbf{x})$$

holds. Otherwise all bets are off

The minimization problem:

Minimization step:

$$\pi_\theta(\mathbf{x}) := \arg \min_{\mathbf{u}} Q_\theta(\mathbf{x}, \mathbf{u})$$

is “easy” only if $Q_\theta(\mathbf{x}, \mathbf{u})$ is convex in \mathbf{u} .

However, generic choices of Q_θ are typically not

Policy gradient methods build directly π_θ and find the θ such that the policy yields the best closed-loop possible performances on the real system

Policy performance

Recall:

$$J(\pi) = \mathbb{E}_{\tau_\pi} [L(\mathbf{x}, \mathbf{u})] = \sum_{k=0}^{\infty} \int L(\mathbf{x}_k, \pi(\mathbf{x}_k)) \mathbb{P}[\mathbf{x}_k | \mathbf{x}_0] \mathbb{P}[\mathbf{x}_0] d\mathbf{x}_0 d\mathbf{x}_k$$

where

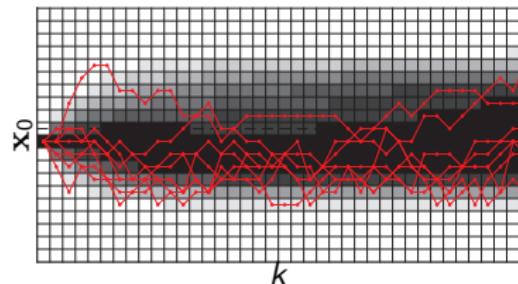
$$\mathbb{P}[\mathbf{x}_k | \mathbf{x}_0] = \int \prod_{i=0}^{k-1} \mathbb{P}[\mathbf{x}_{i+1} | \mathbf{x}_i, \pi(\mathbf{x}_i)] d\mathbf{x}_1 \dots d\mathbf{x}_{k-1}$$

Hence

$$J(\pi_\theta) = \mathbb{E}_{\tau_{\pi_\theta}} [L(\mathbf{x}, \pi_\theta)]$$

Note that

$$J(\pi_\theta) = \mathbb{E}[V_{\pi_\theta}(\mathbf{x}_0)] = \int V_{\pi_\theta}(\mathbf{x}_0) \mathbb{P}[\mathbf{x}_0] d\mathbf{x}_0$$



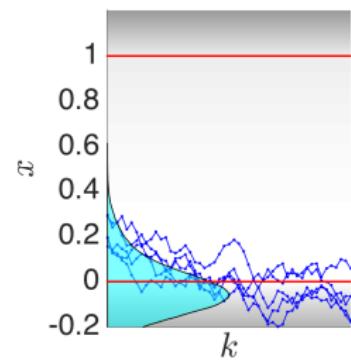
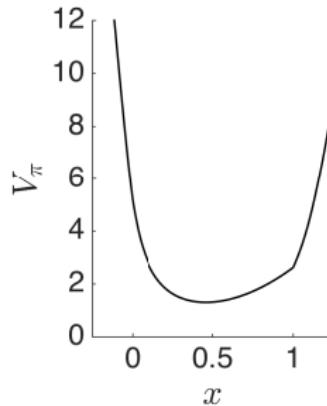
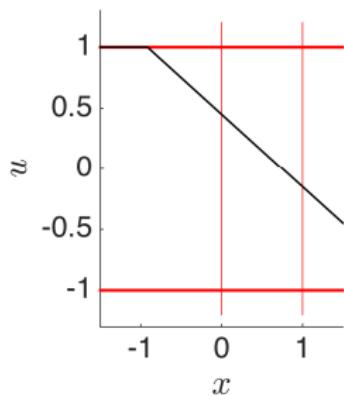
i.e. expectation is taken over initial conditions.

Policy performance - MDP distribution τ_π

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π

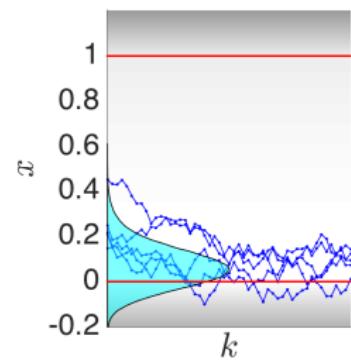
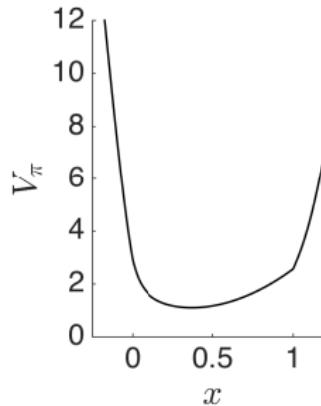
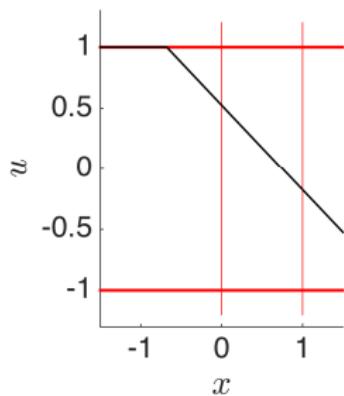
Policy performance - MDP distribution τ_π

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π



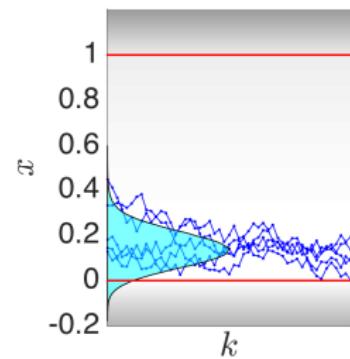
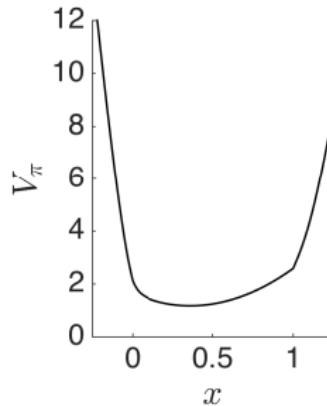
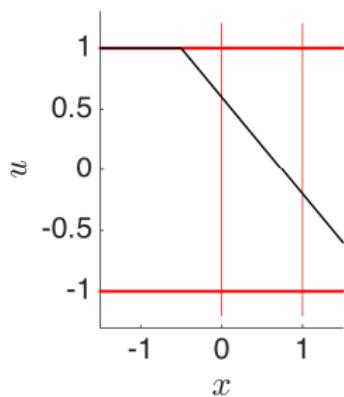
Policy performance - MDP distribution τ_π

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π



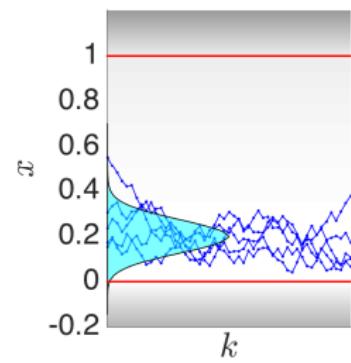
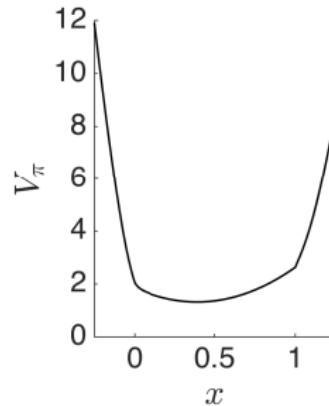
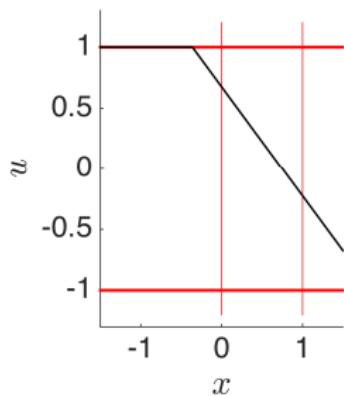
Policy performance - MDP distribution τ_π

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π



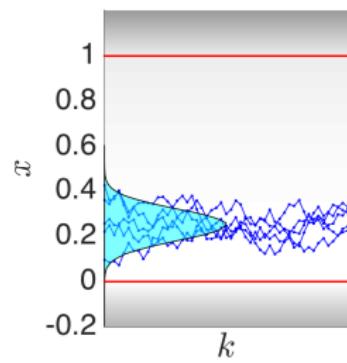
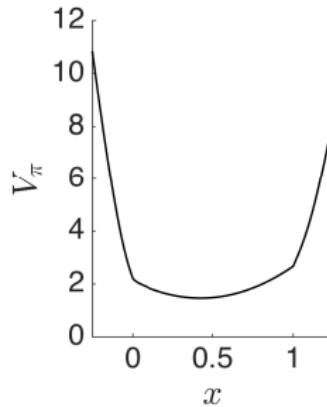
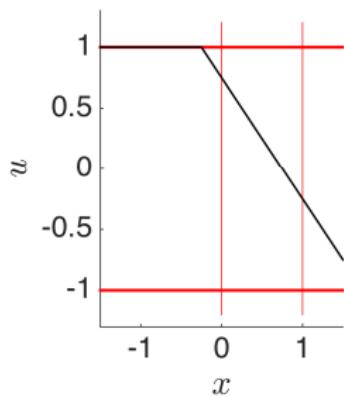
Policy performance - MDP distribution τ_π

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π



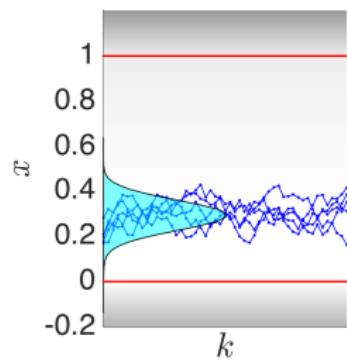
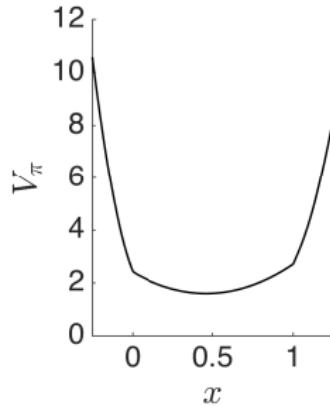
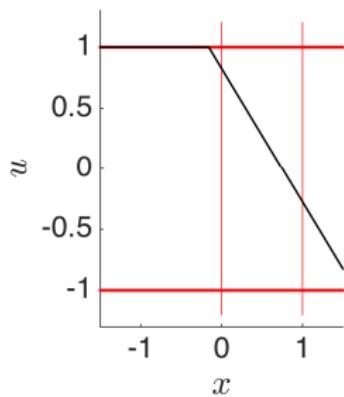
Policy performance - MDP distribution τ_π

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π



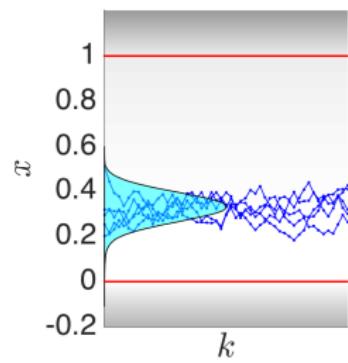
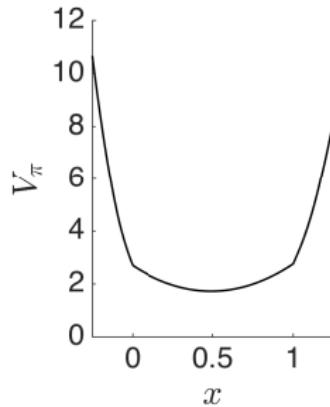
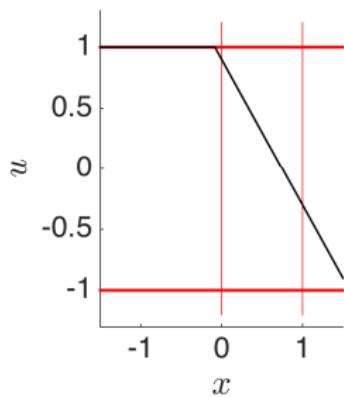
Policy performance - MDP distribution τ_π

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π



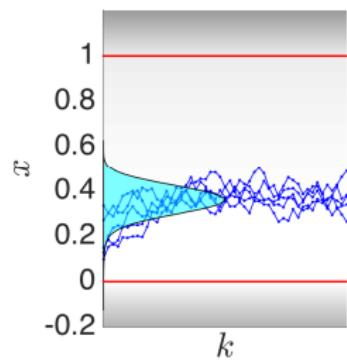
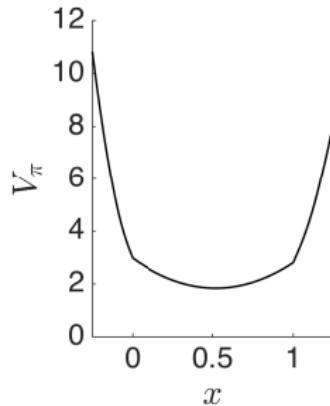
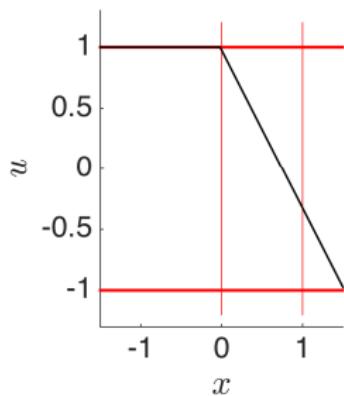
Policy performance - MDP distribution τ_π

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π



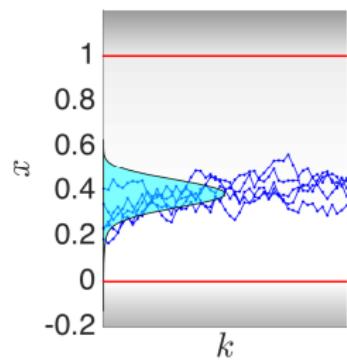
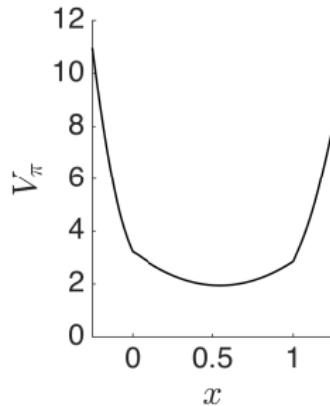
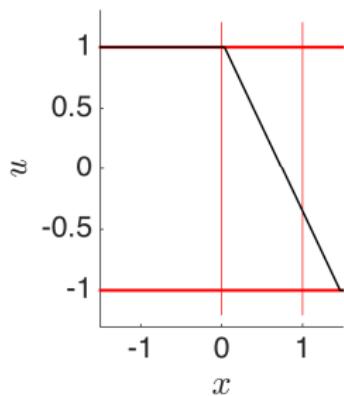
Policy performance - MDP distribution τ_π

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π



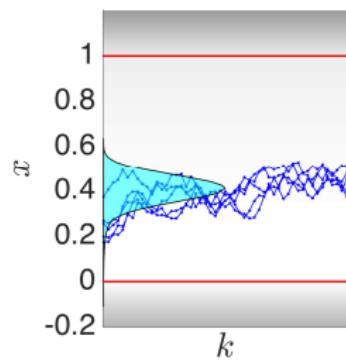
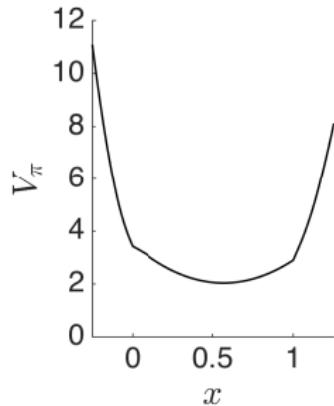
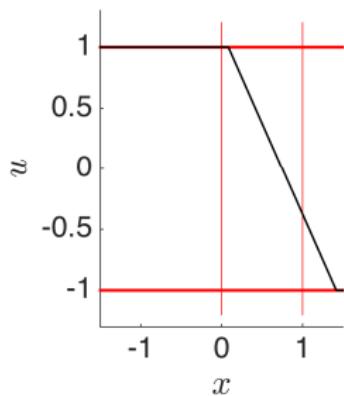
Policy performance - MDP distribution τ_π

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π



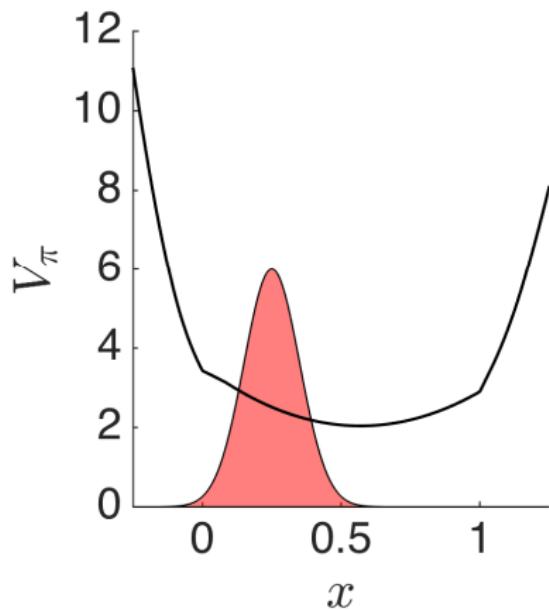
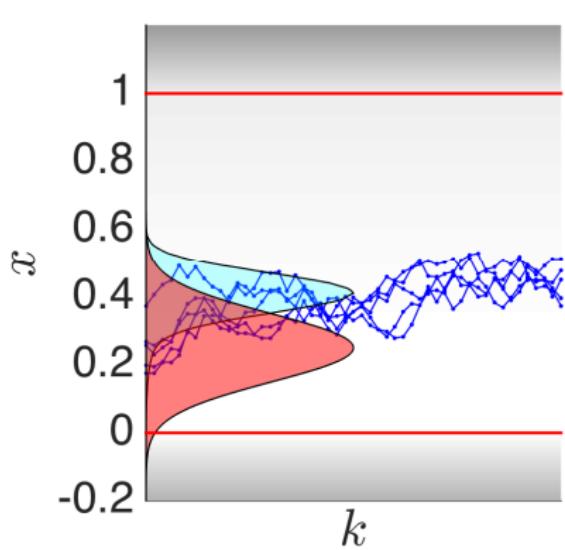
Policy performance - MDP distribution τ_π

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π



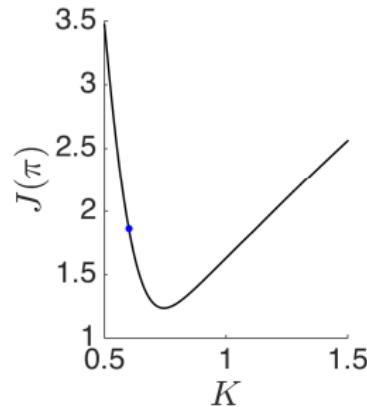
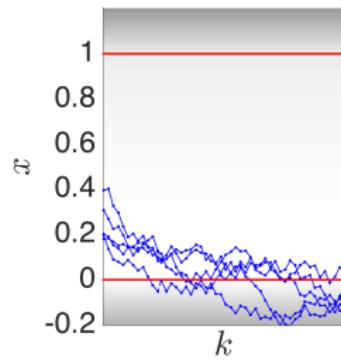
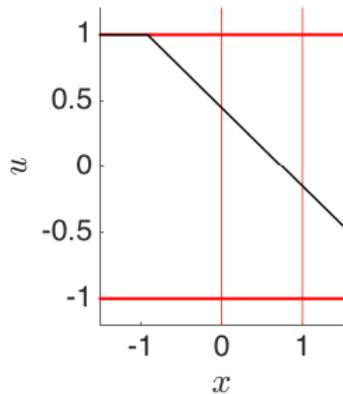
Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



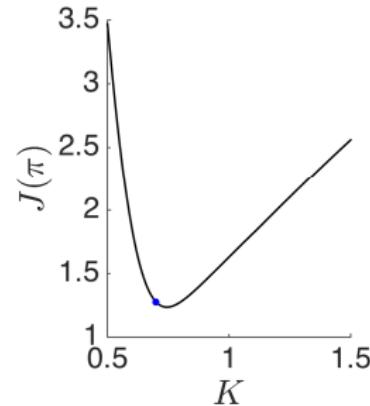
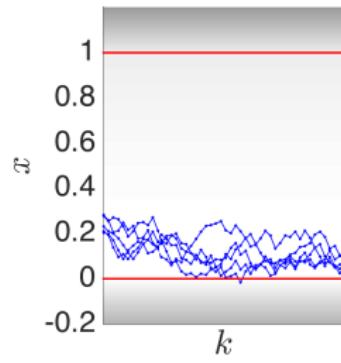
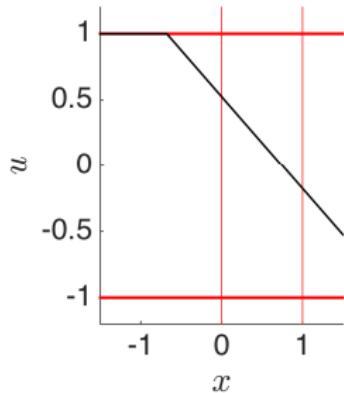
Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



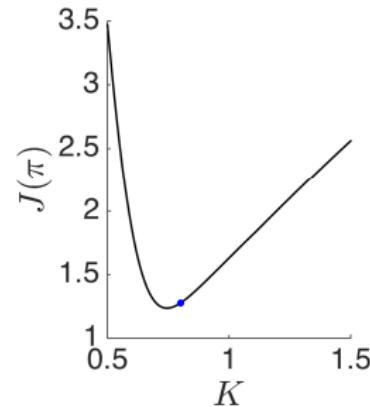
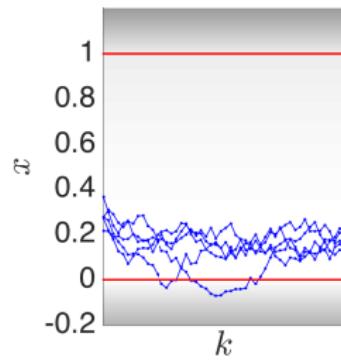
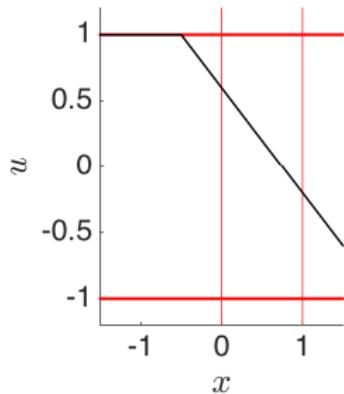
Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



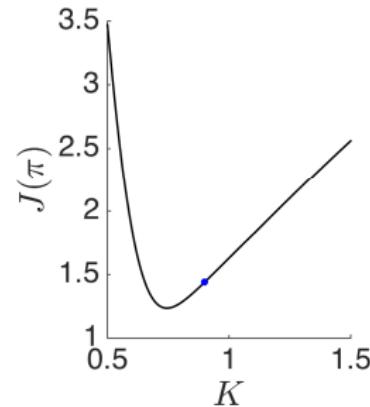
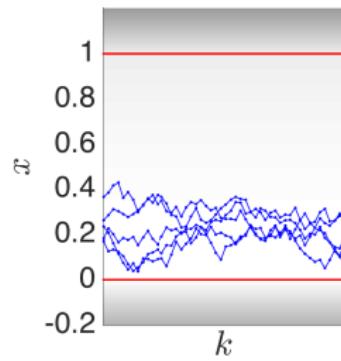
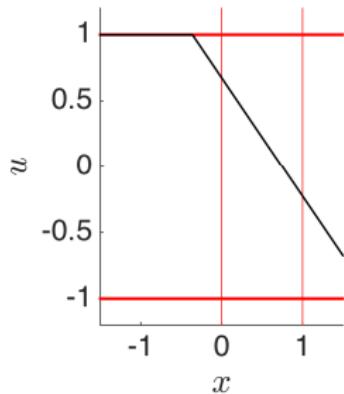
Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



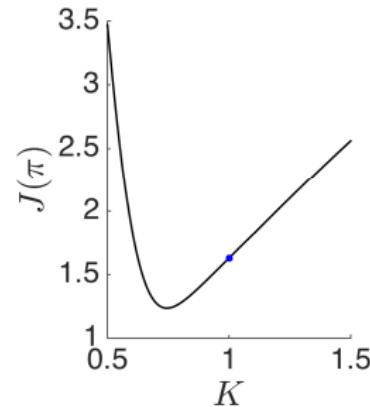
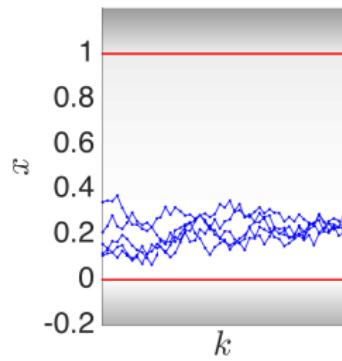
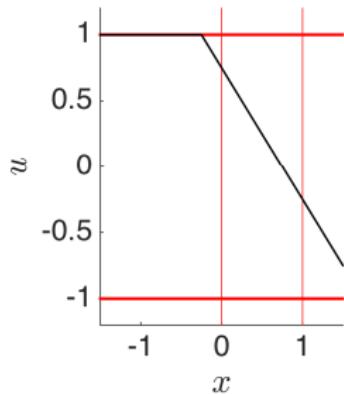
Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



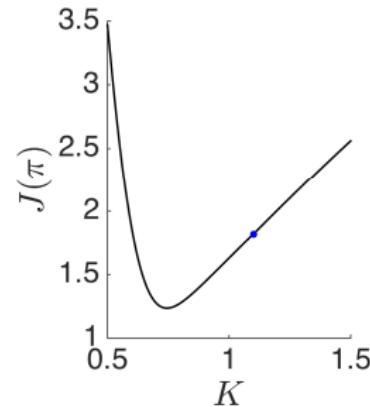
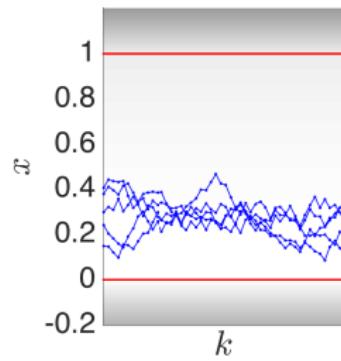
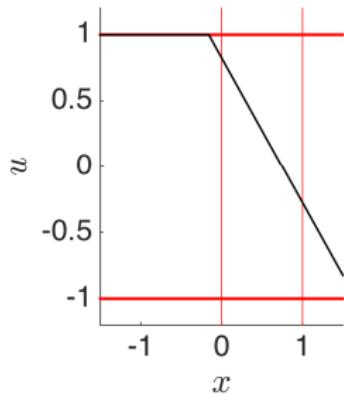
Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



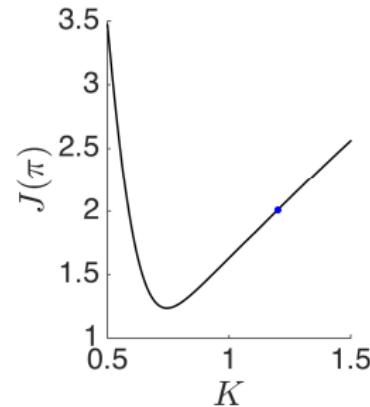
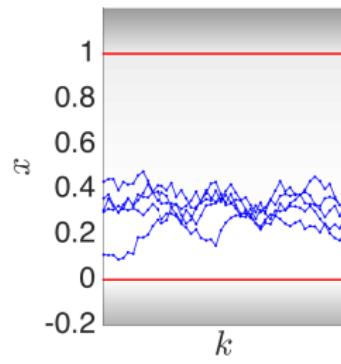
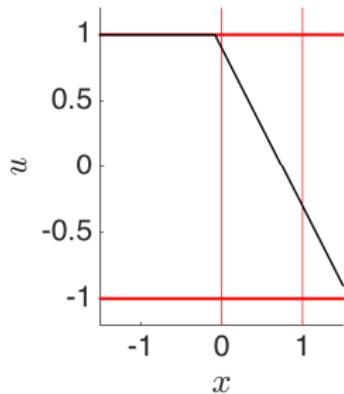
Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



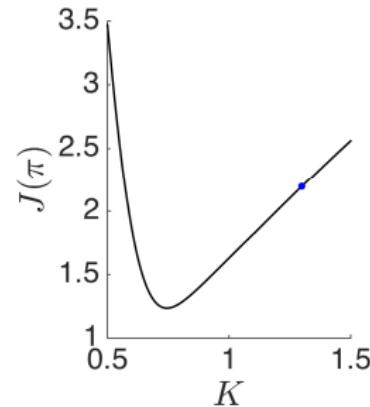
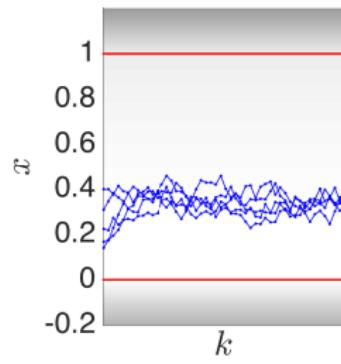
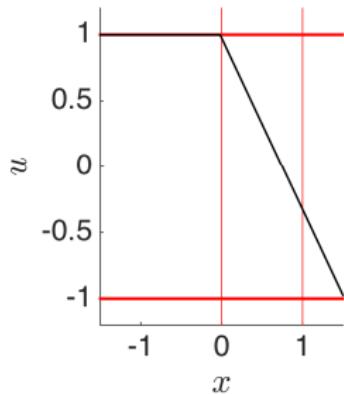
Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



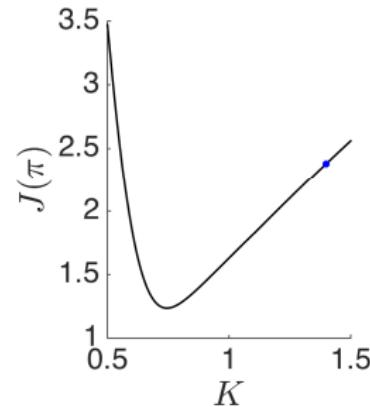
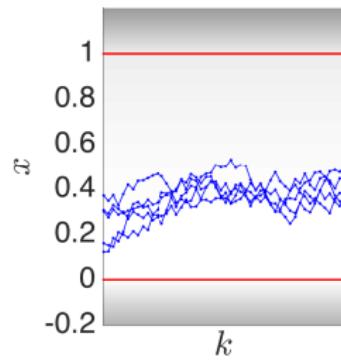
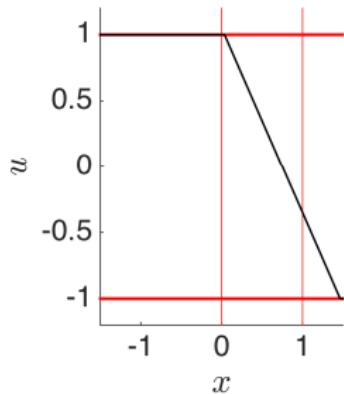
Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



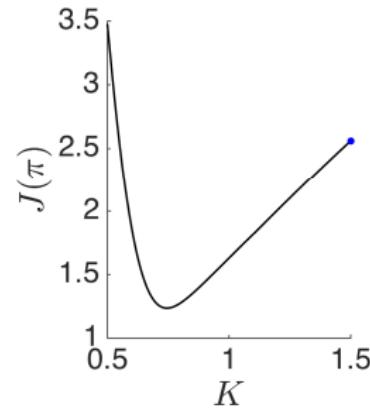
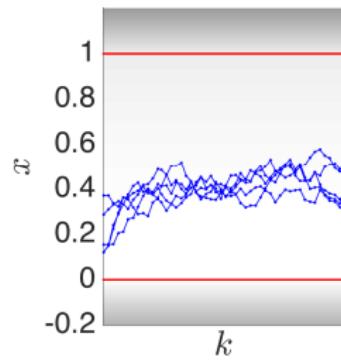
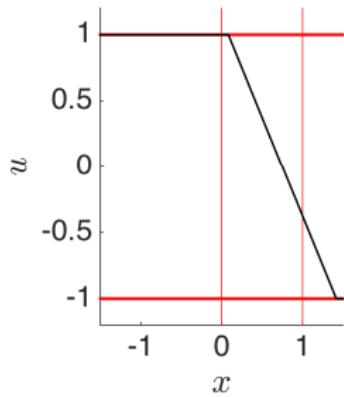
Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



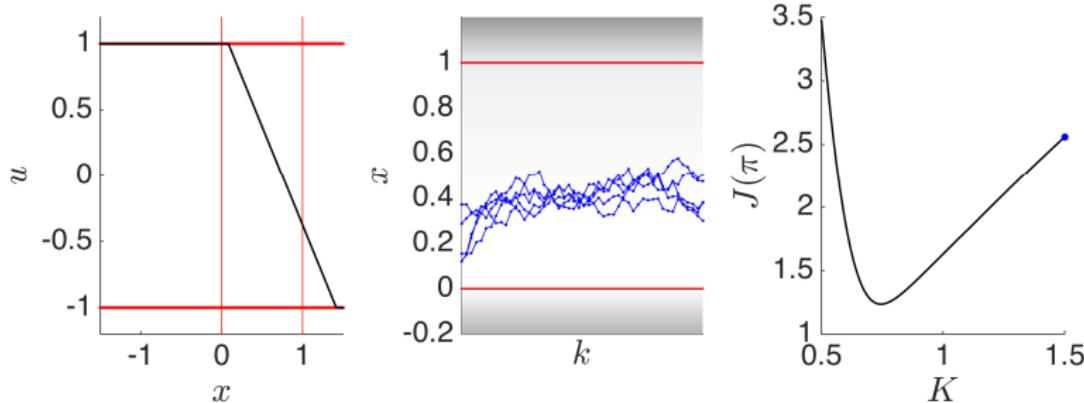
Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



Policy performance - $J(\pi)$

- Policy: $\pi(x) = \text{mid}(\underline{u}, -K(x - x_{\text{ref}}), \bar{u})$, with $x_{\text{ref}} = 0.75$
- Initial condition $x_0 \sim \mathcal{N}(0.25, 0.1)$
- Yields MDP distribution τ_π , value function V_π
- Yields policy performance $J(\pi) = \mathbb{E}[V_\pi(x_0)] = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))]$



Optimal θ must satisfy: $\nabla_\theta J(\pi_\theta) = 0$

Outline

1 Introduction

2 Deterministic policy gradient

3 Stochastic policy gradient



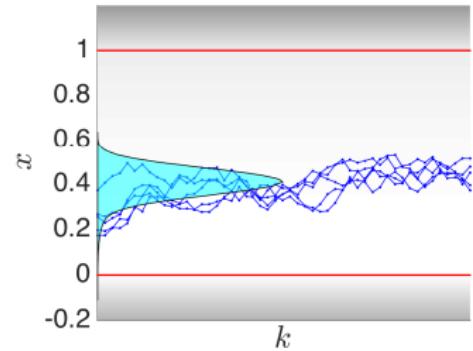
Deterministic Policy gradient theorem

How to compute $\nabla_{\theta} J(\pi_{\theta})$?

- changing θ changes both π_{θ} and $\tau_{\pi_{\theta}}$
- policy gradient $\nabla_{\theta} J(\pi_{\theta})$ is non-trivial...

Reminder:

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



Deterministic Policy gradient theorem

How to compute $\nabla_{\theta} J(\pi_{\theta})$?

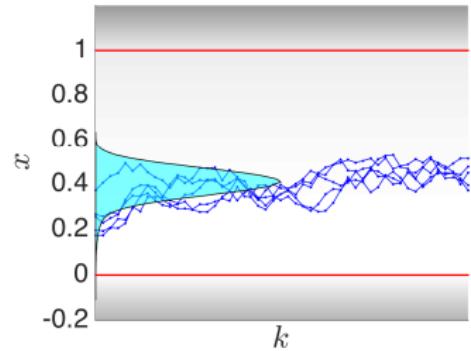
- changing θ changes both π_{θ} and $\tau_{\pi_{\theta}}$
- policy gradient $\nabla_{\theta} J(\pi_{\theta})$ is non-trivial...

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}} (\mathbf{x}, \pi_{\theta})]$$

Reminder:

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



Deterministic Policy gradient theorem

How to compute $\nabla_{\theta} J(\pi_{\theta})$?

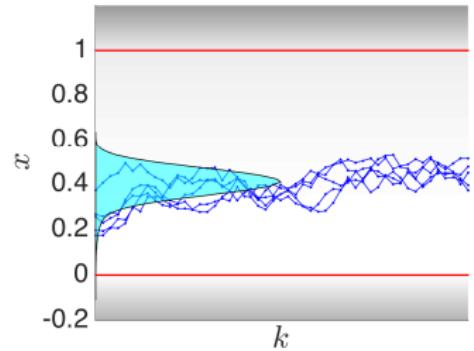
- changing θ changes both π_{θ} and $\tau_{\pi_{\theta}}$
- policy gradient $\nabla_{\theta} J(\pi_{\theta})$ is non-trivial...

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$

Reminder:

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



where $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ is the action-value function associated to policy π_{θ} i.e.

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

$$V_{\pi_{\theta}}(\mathbf{x}) = L(\mathbf{x}, \pi_{\theta}) + \gamma \mathbb{E}[V_{\pi_{\theta}}(\mathbf{x}_+) | \mathbf{x}, \pi_{\theta}]$$

Deterministic Policy gradient theorem

How to compute $\nabla_{\theta} J(\pi_{\theta})$?

- changing θ changes both π_{θ} and $\tau_{\pi_{\theta}}$
- policy gradient $\nabla_{\theta} J(\pi_{\theta})$ is non-trivial...

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$

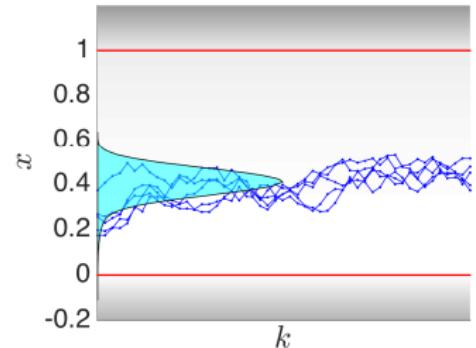
(i.e. we can do a chain-rule-like operation:

$$\nabla_{\theta} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta}) = \nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})$$

and ignore $\tau_{\pi_{\theta}}$)

Reminder:

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



where $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ is the action-value function associated to policy π_{θ} i.e.

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

$$V_{\pi_{\theta}}(\mathbf{x}) = L(\mathbf{x}, \pi_{\theta}) + \gamma \mathbb{E}[V_{\pi_{\theta}}(\mathbf{x}_+) | \mathbf{x}, \pi_{\theta}]$$

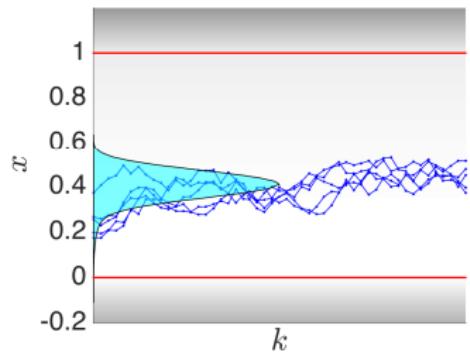
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



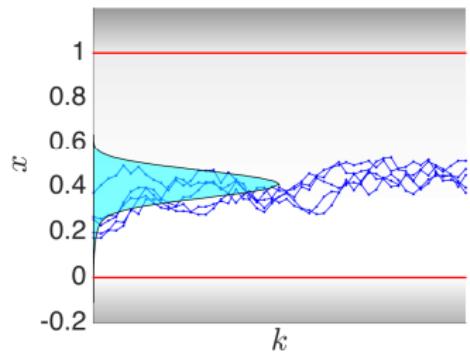
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



Remarks

- Need π_{θ} to be (almost everywhere) differentiable w.r.t. θ (why would it not be?)
- Need $Q_{\pi_{\theta}}$ to be (almost everywhere) differentiable w.r.t. θ (why would it not be?)
- Evaluating $\nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})$ from data requires departing from policy π_{θ}
- Requires some regularity assumptions (non-trivial, not discussed here)
- Measure theory may kick in (because of $\mathbb{E}_{\tau_{\pi_{\theta}}}$)

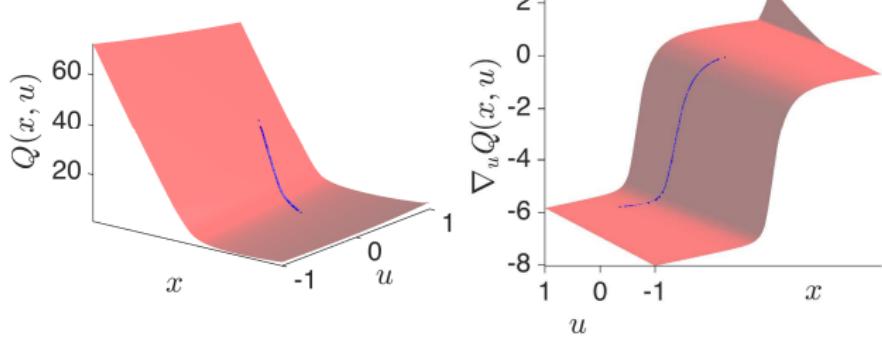
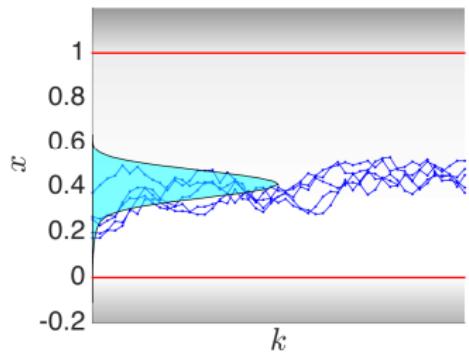
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_u Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



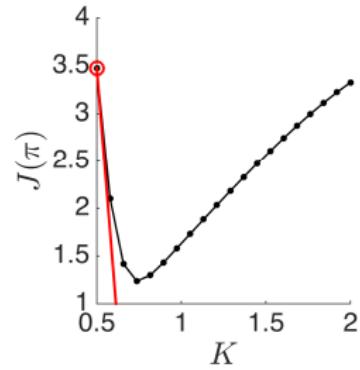
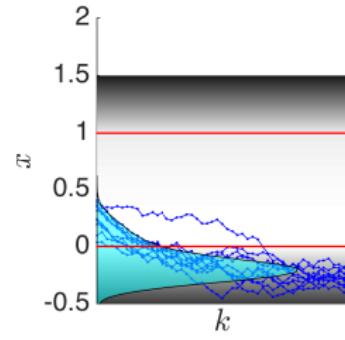
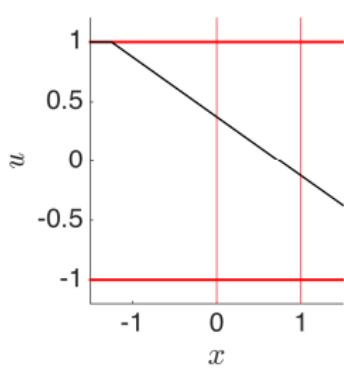
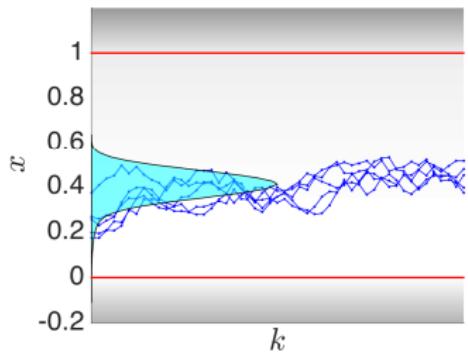
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



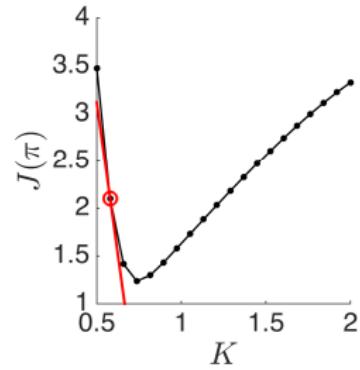
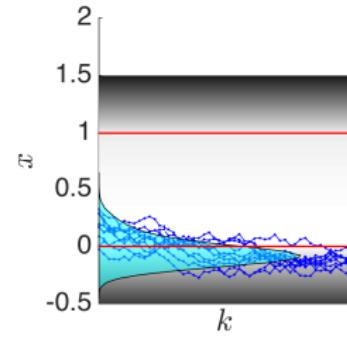
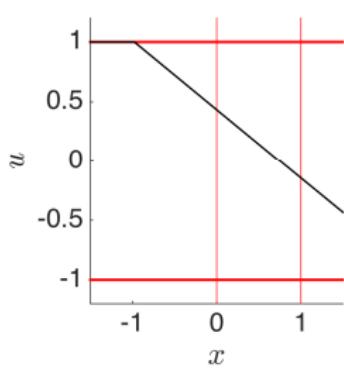
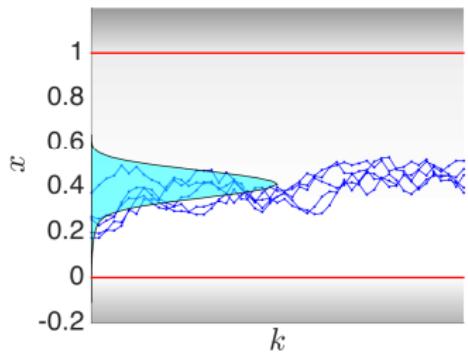
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



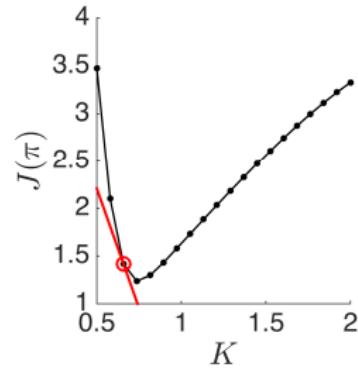
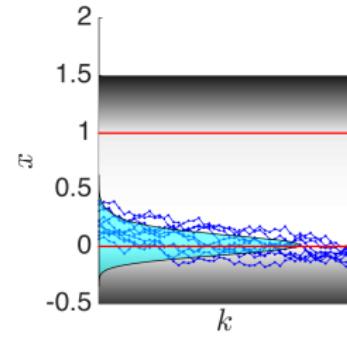
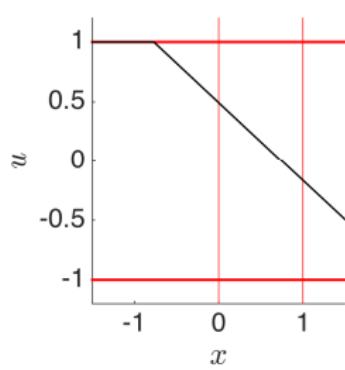
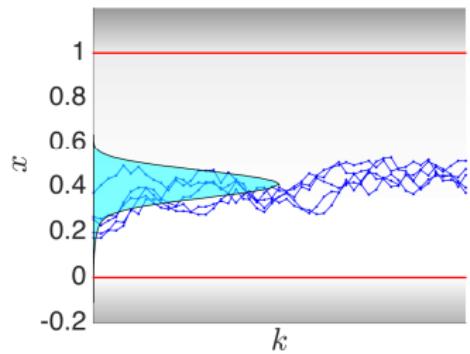
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



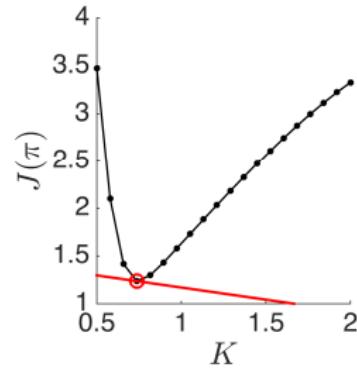
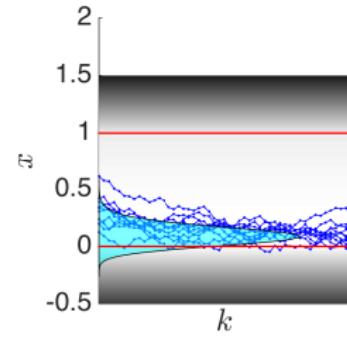
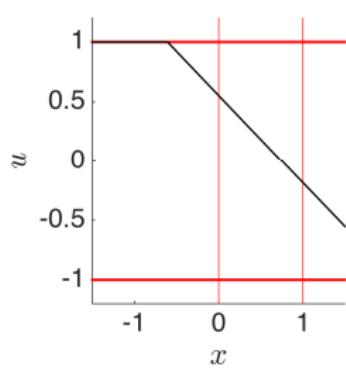
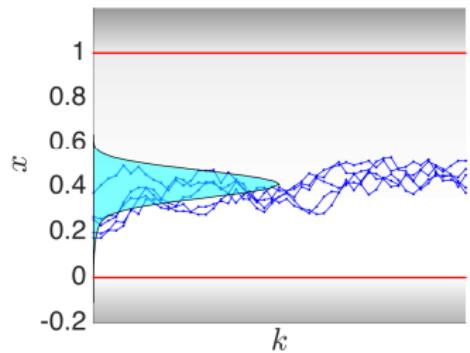
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



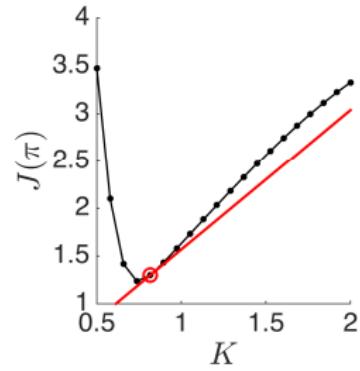
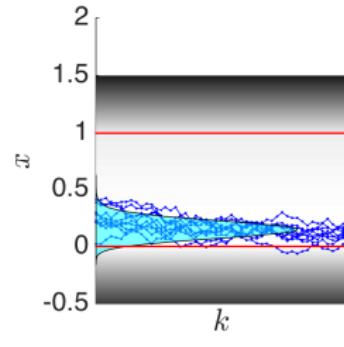
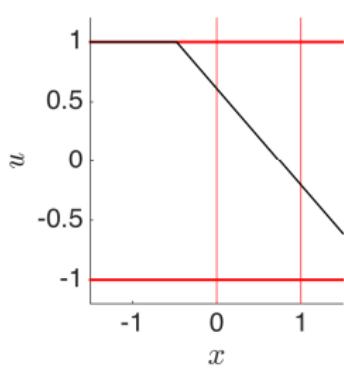
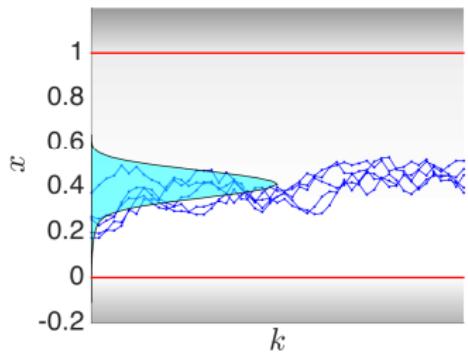
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



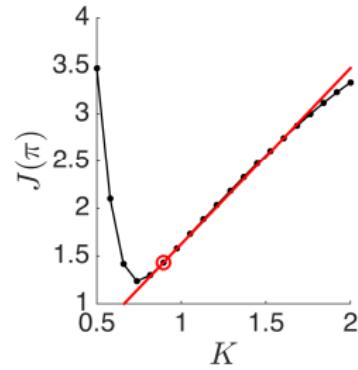
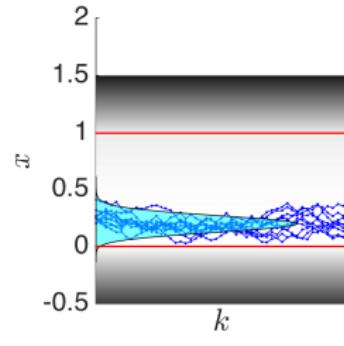
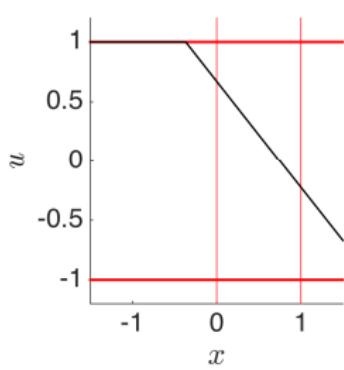
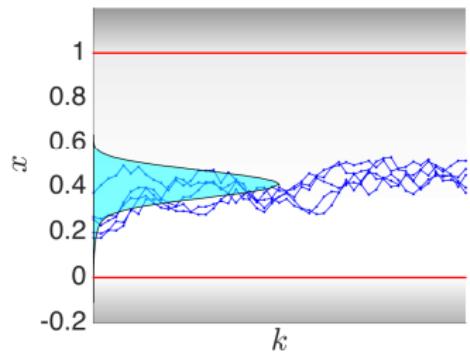
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



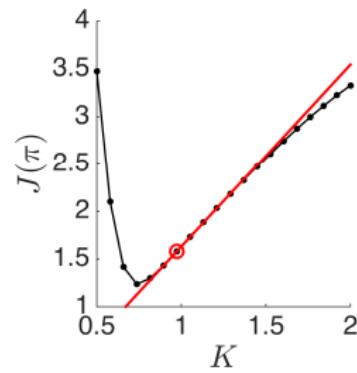
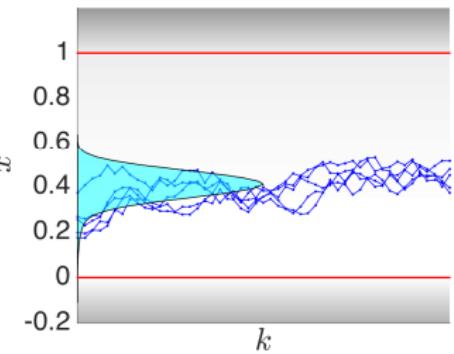
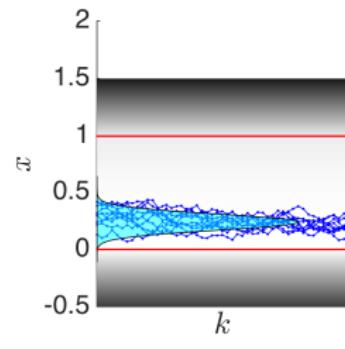
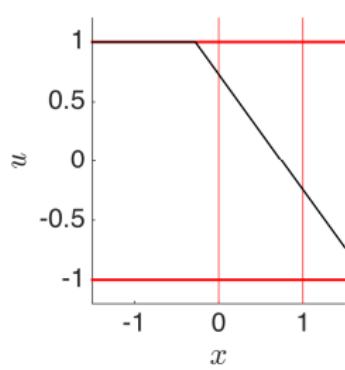
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



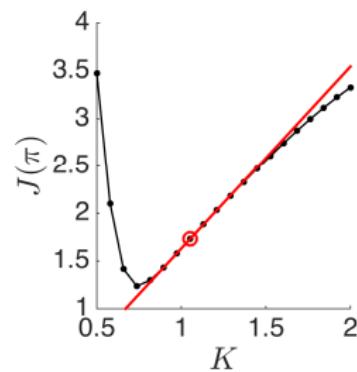
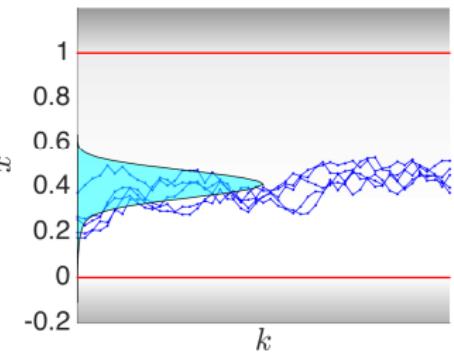
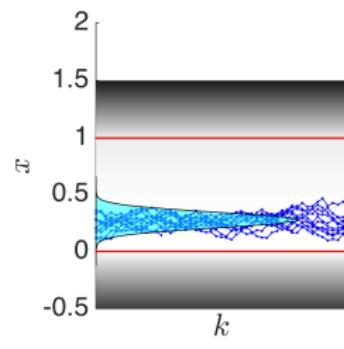
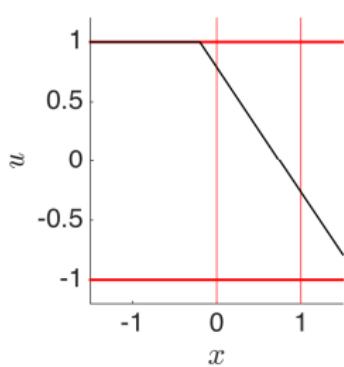
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



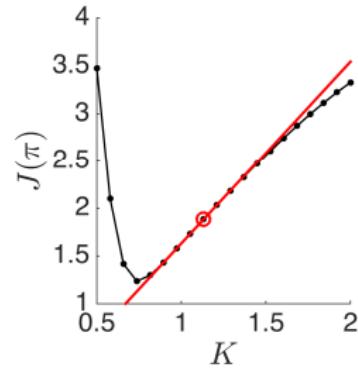
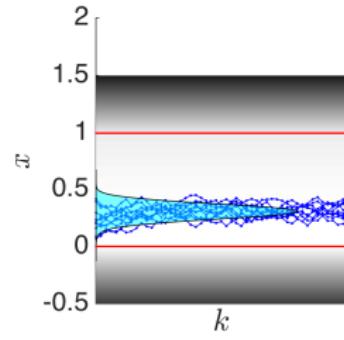
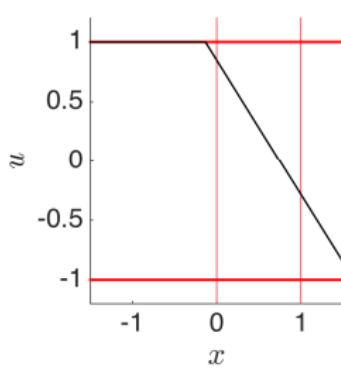
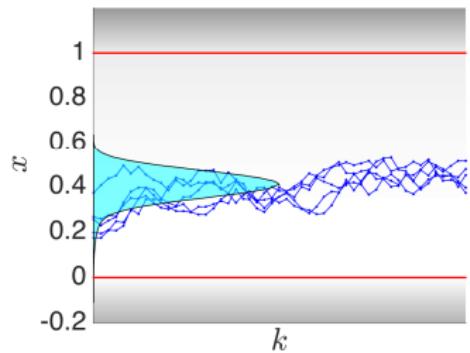
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



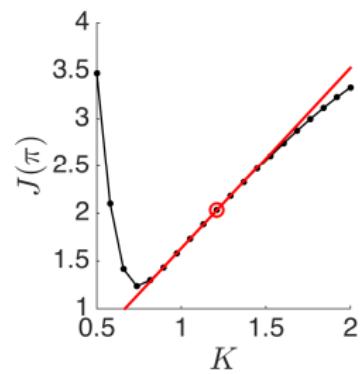
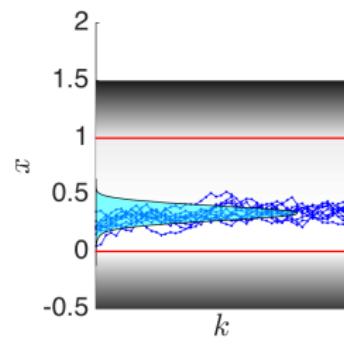
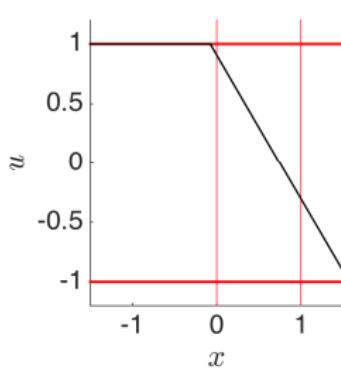
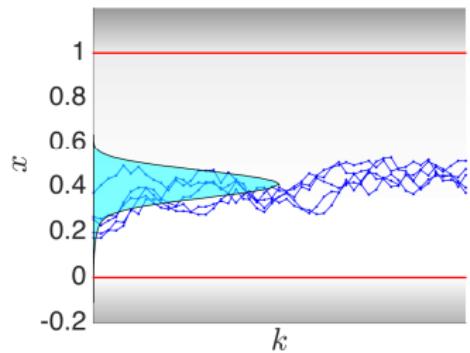
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



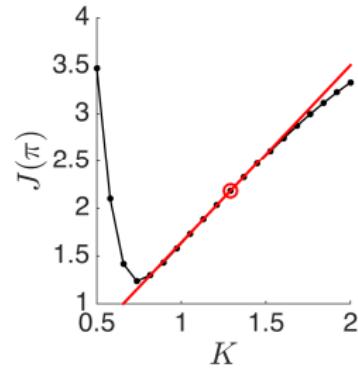
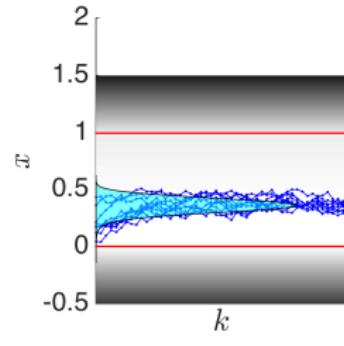
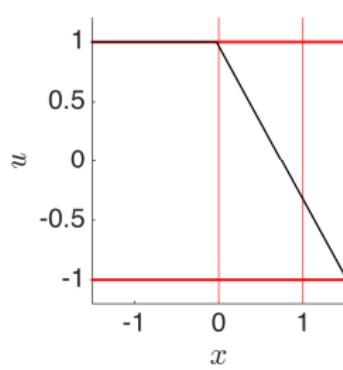
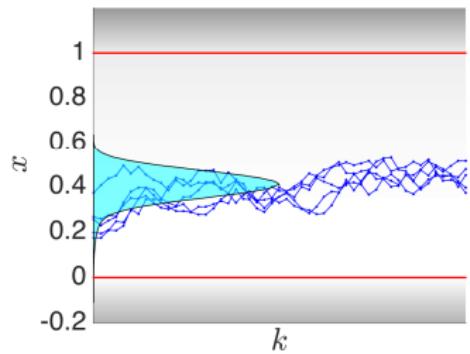
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



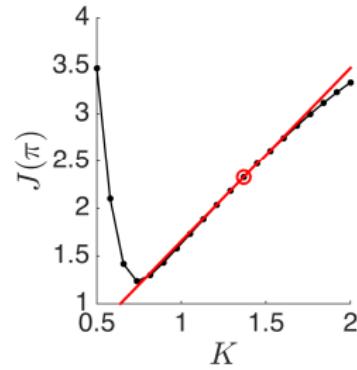
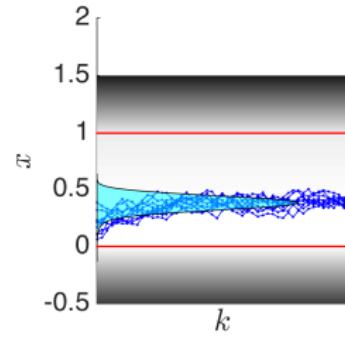
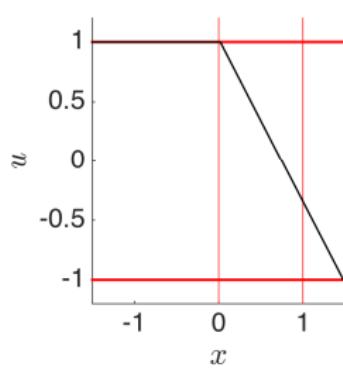
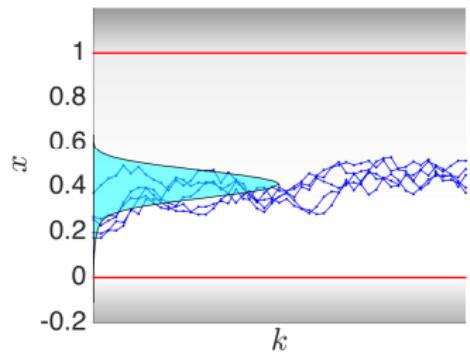
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



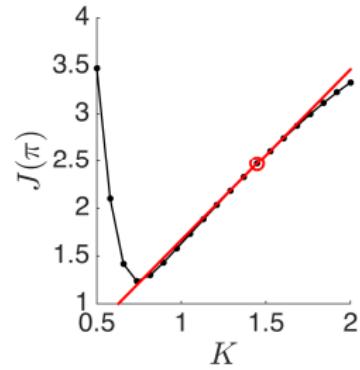
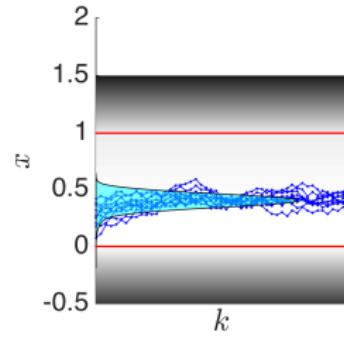
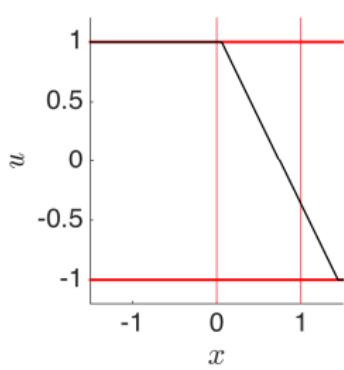
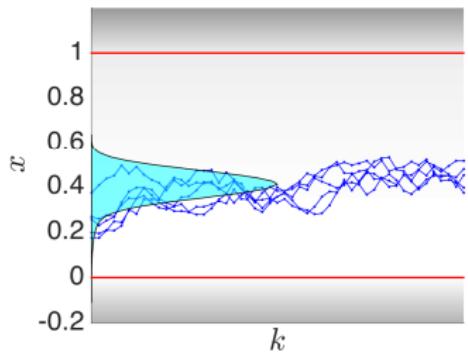
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



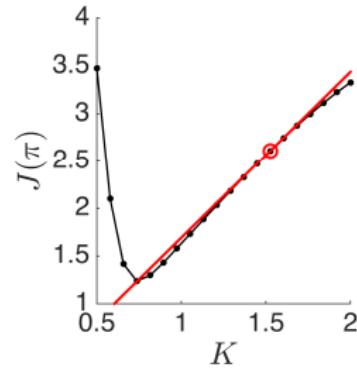
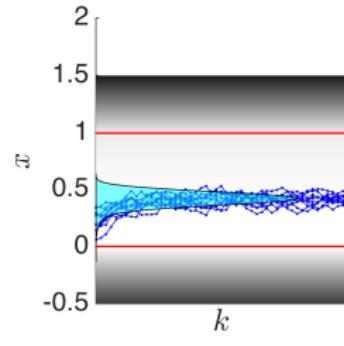
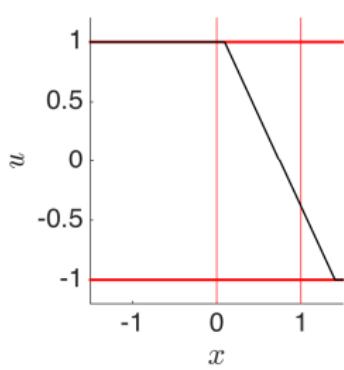
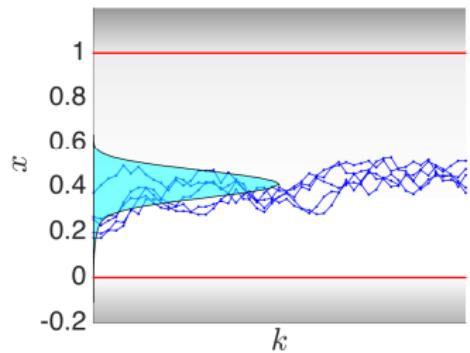
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



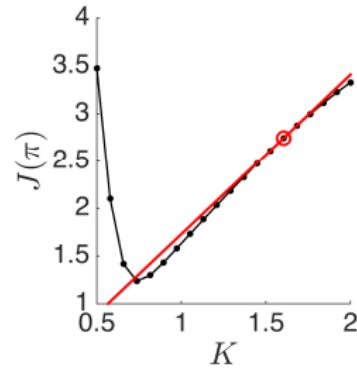
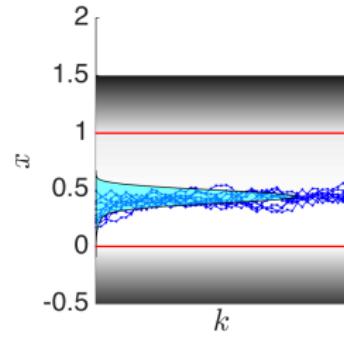
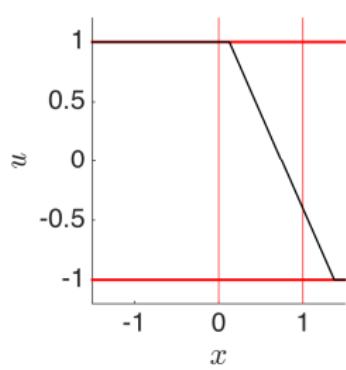
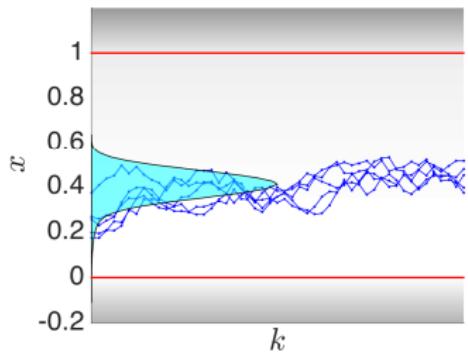
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



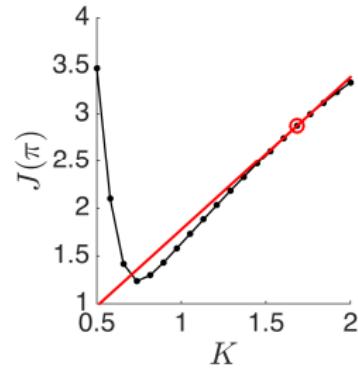
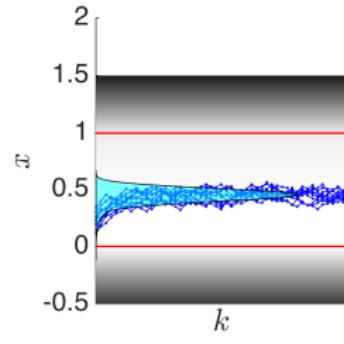
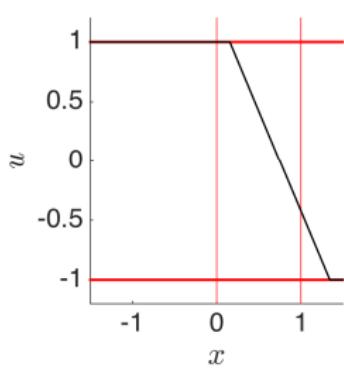
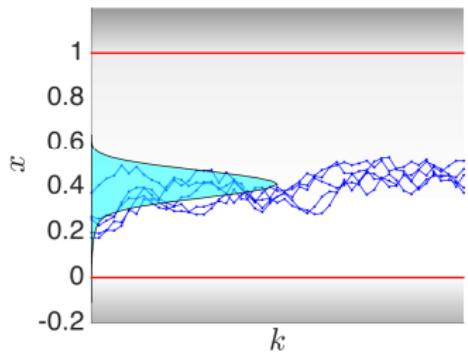
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



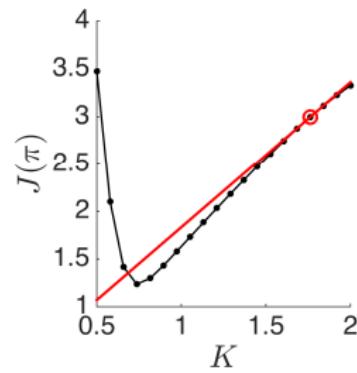
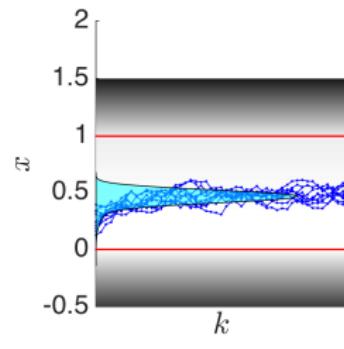
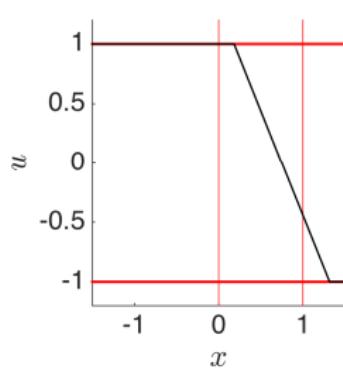
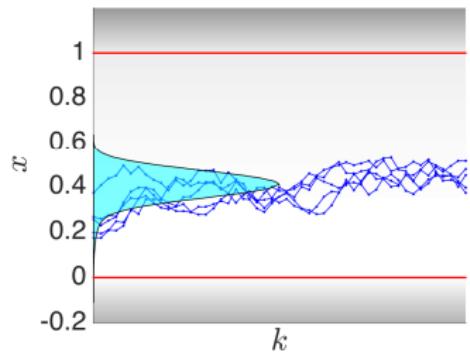
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



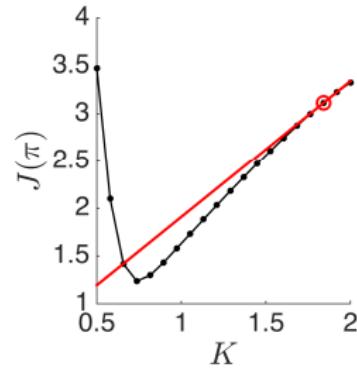
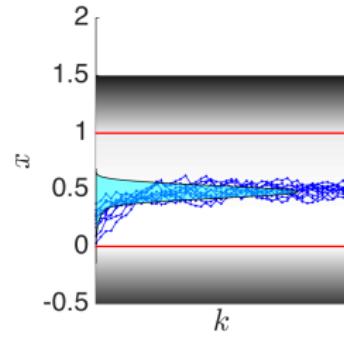
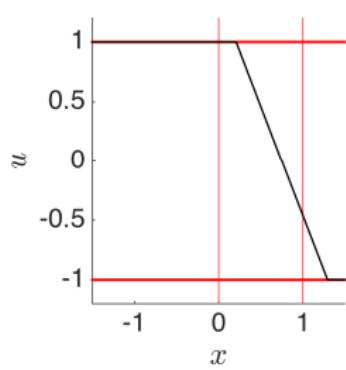
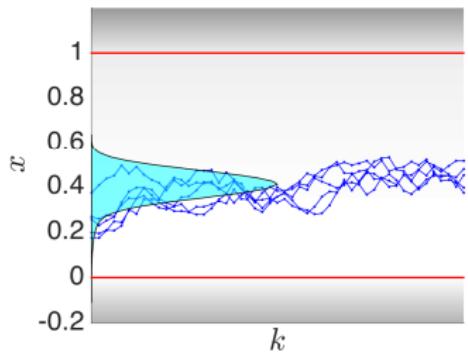
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



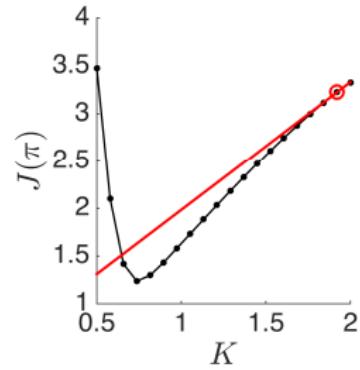
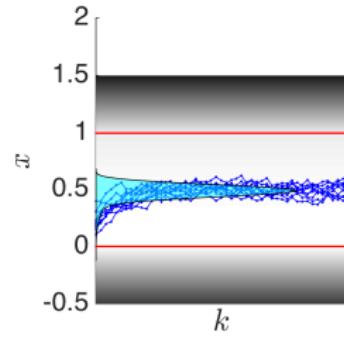
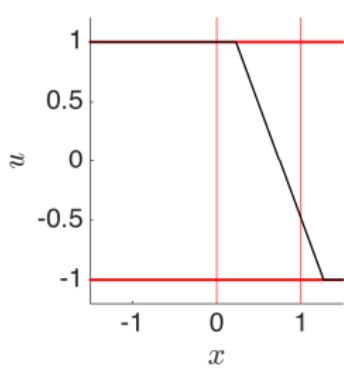
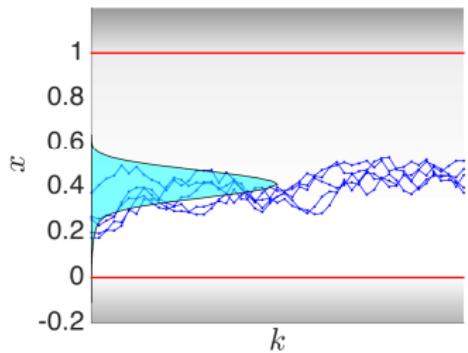
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



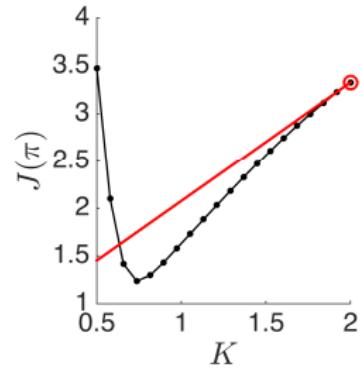
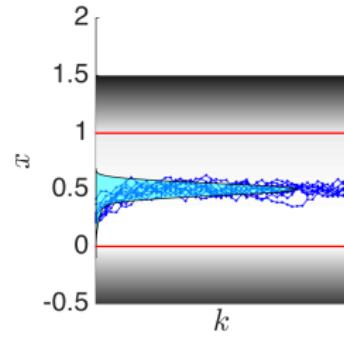
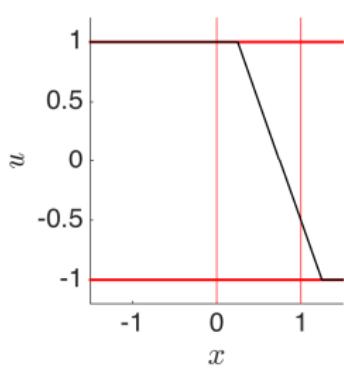
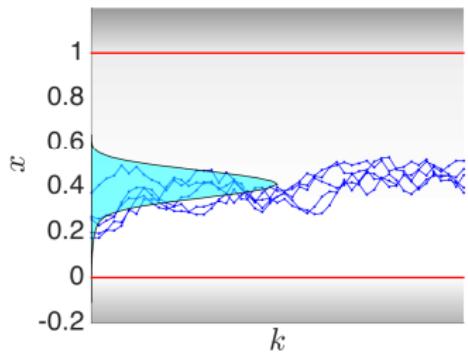
Deterministic Policy gradient theorem

Reminder:

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

Policy gradient theorem

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$

Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

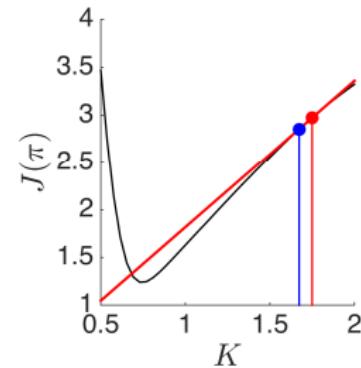
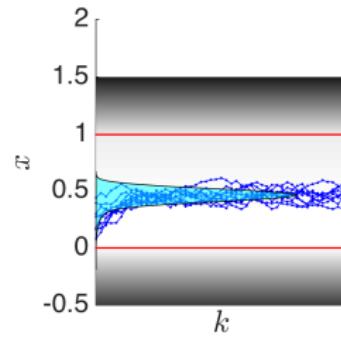
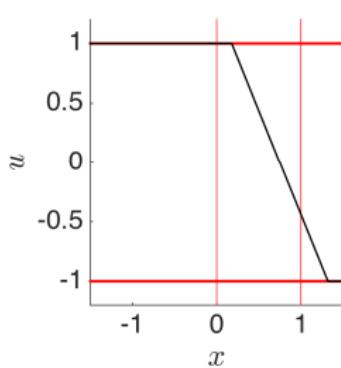
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

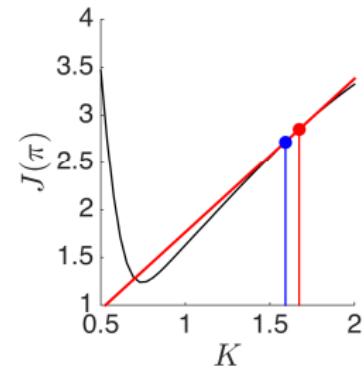
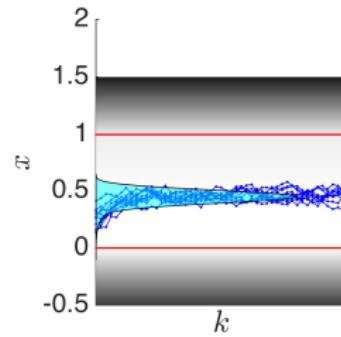
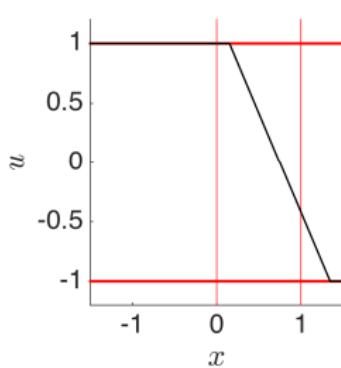
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

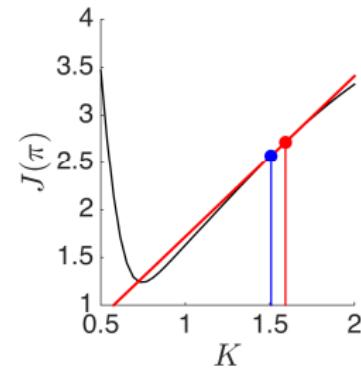
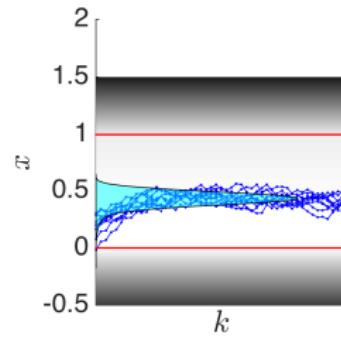
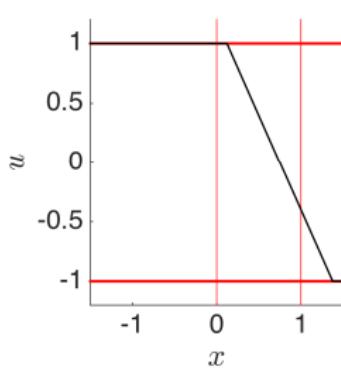
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

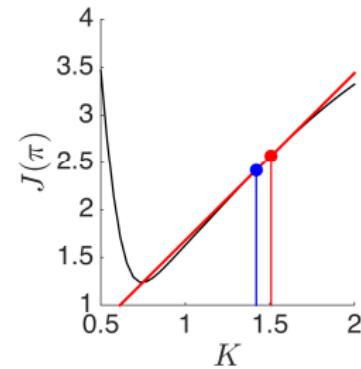
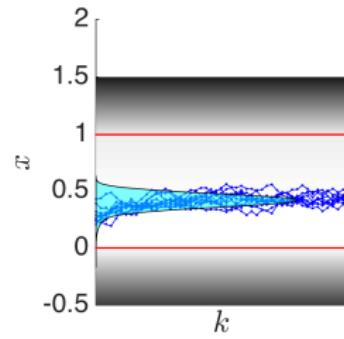
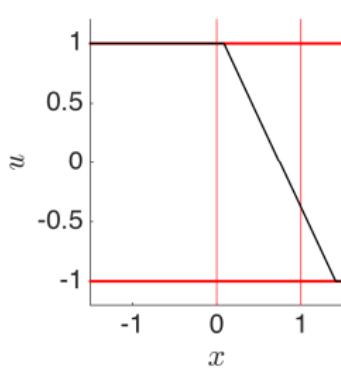
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

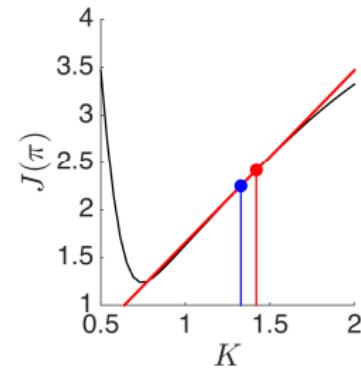
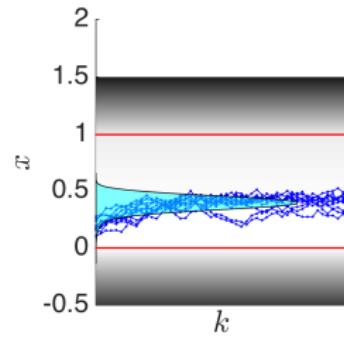
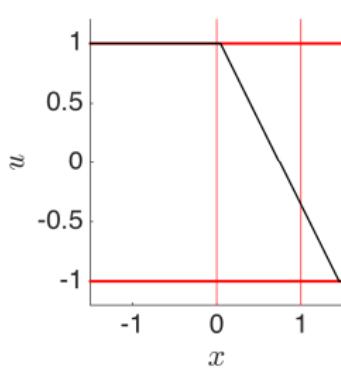
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

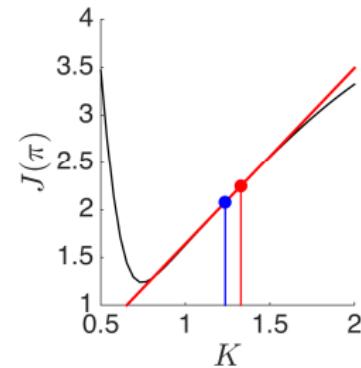
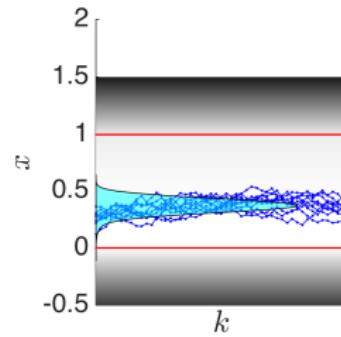
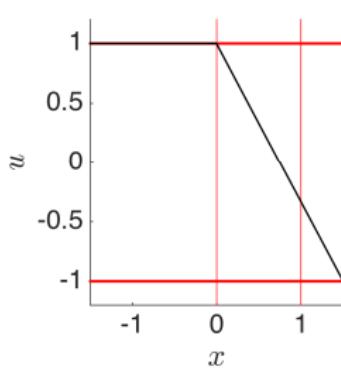
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

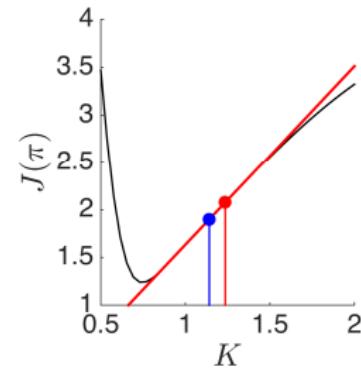
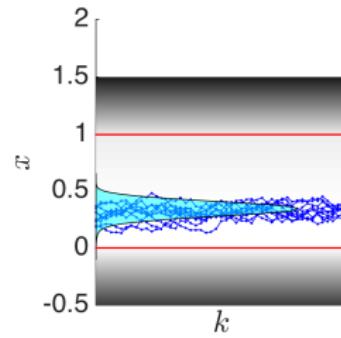
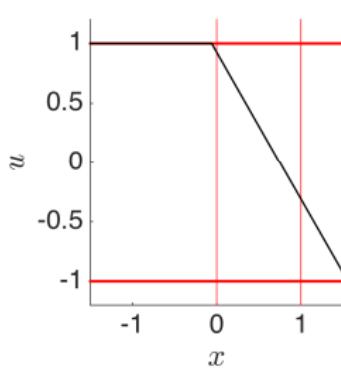
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

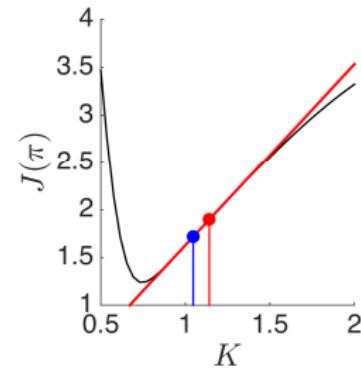
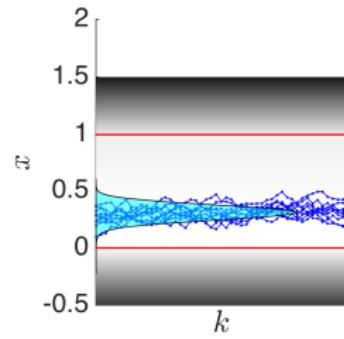
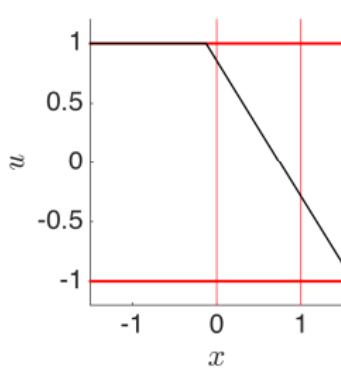
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

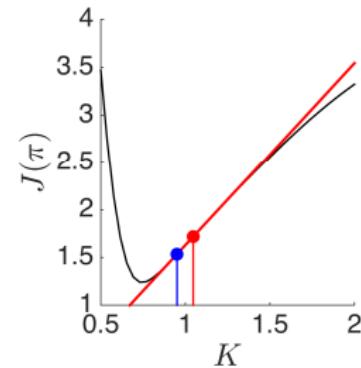
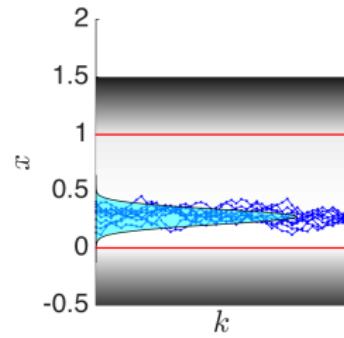
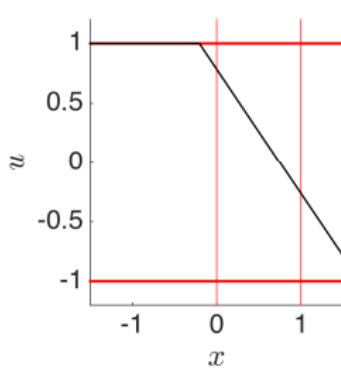
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

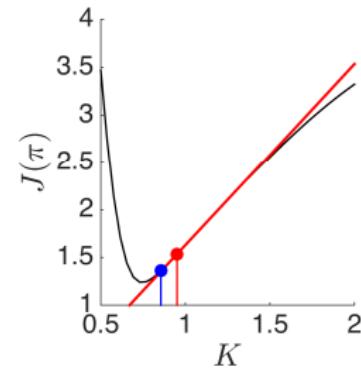
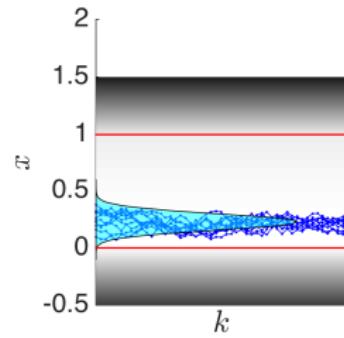
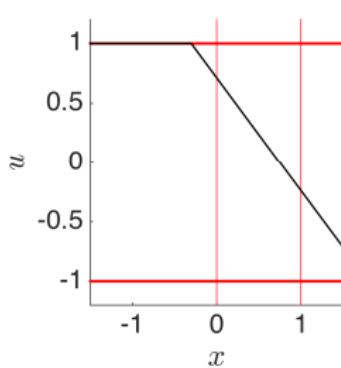
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

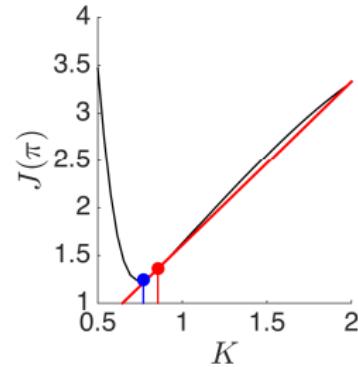
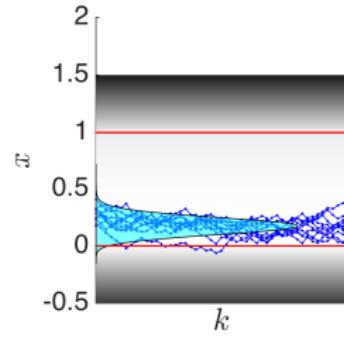
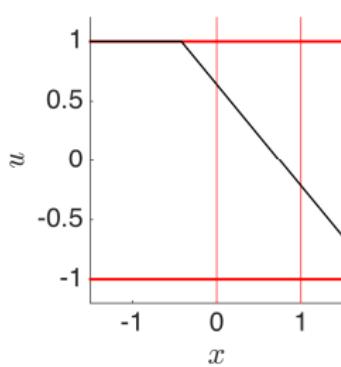
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

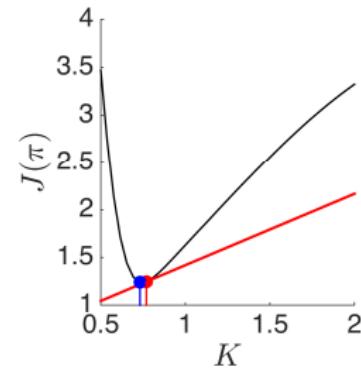
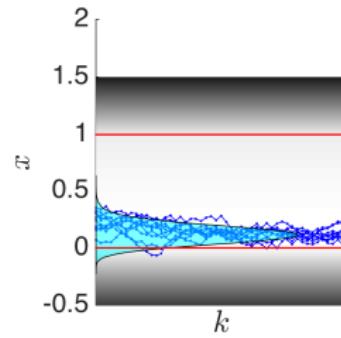
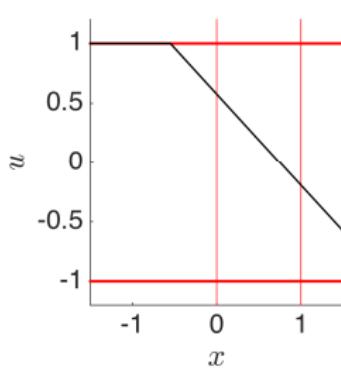
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

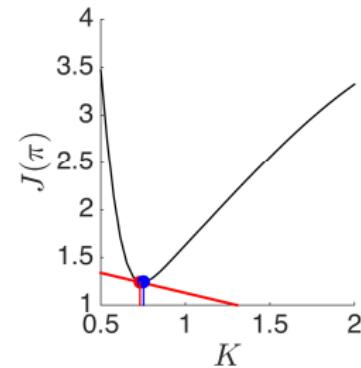
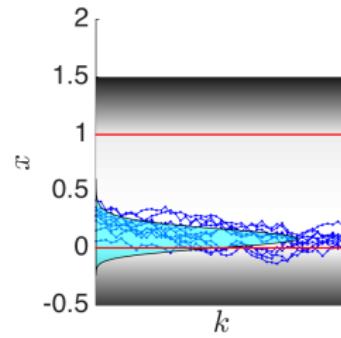
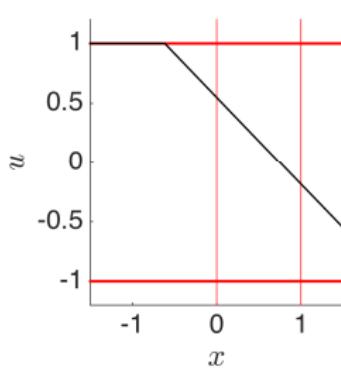
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

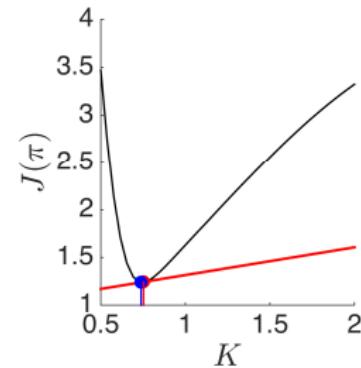
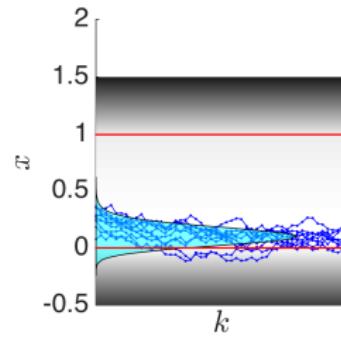
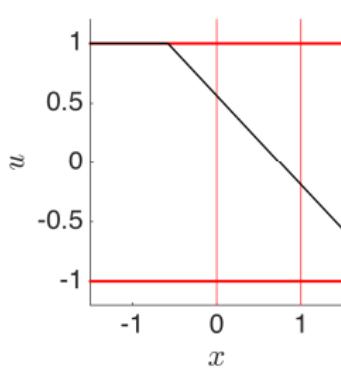
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

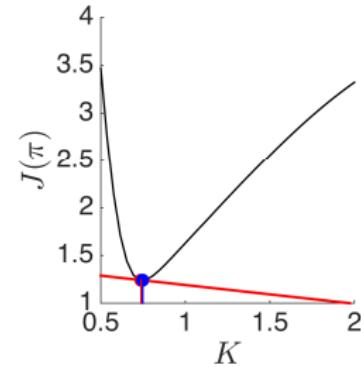
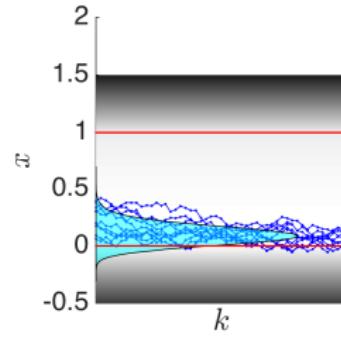
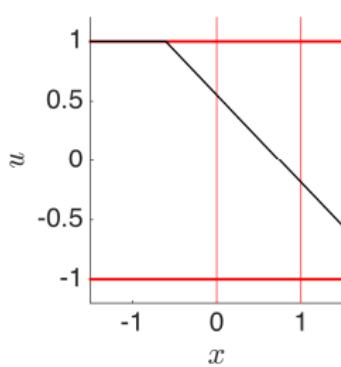
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

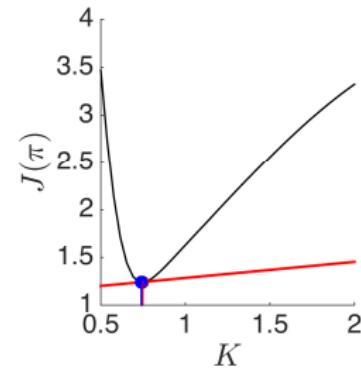
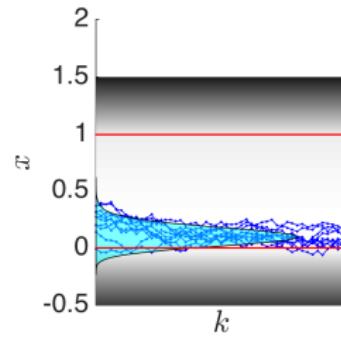
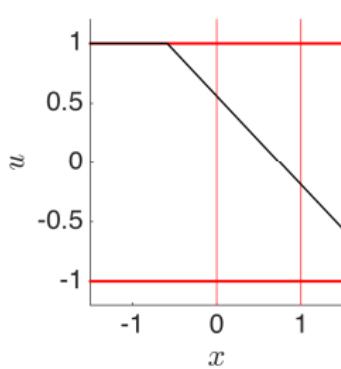
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

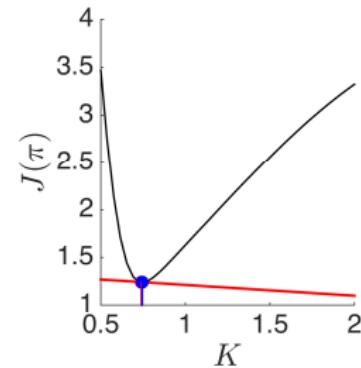
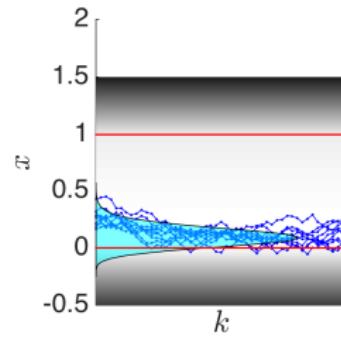
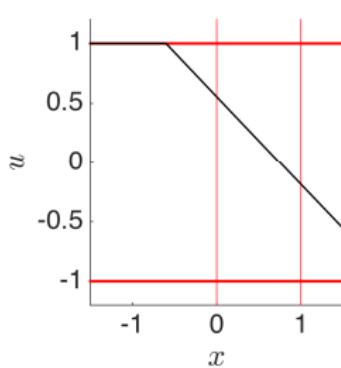
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

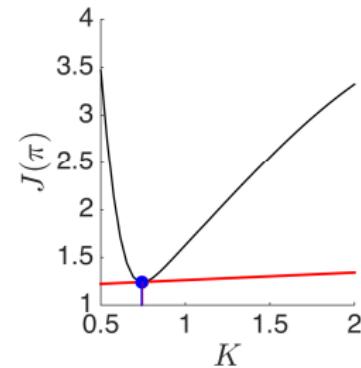
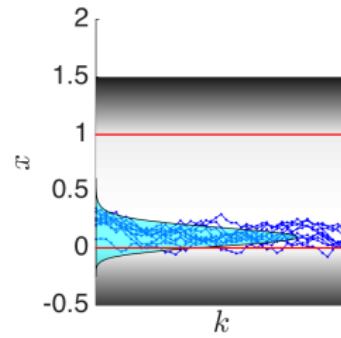
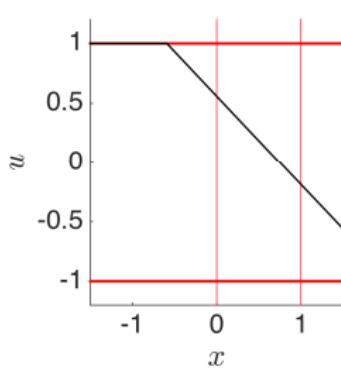
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

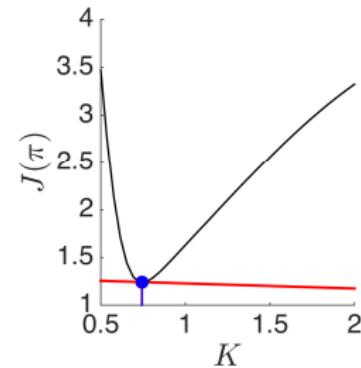
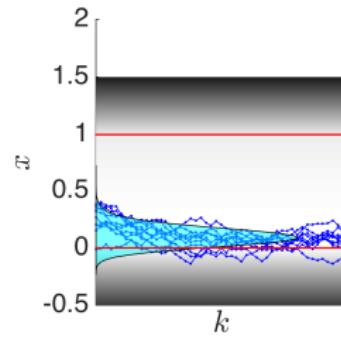
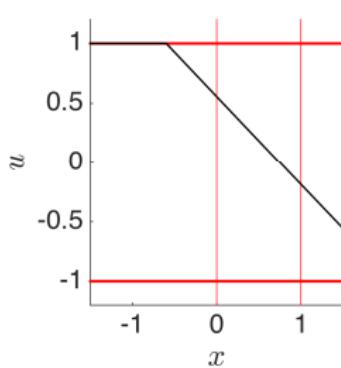
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

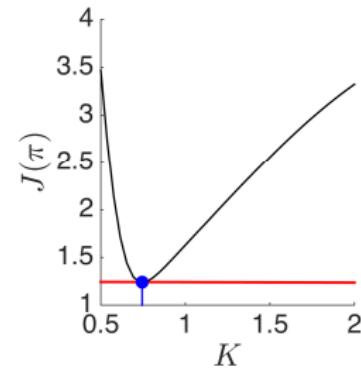
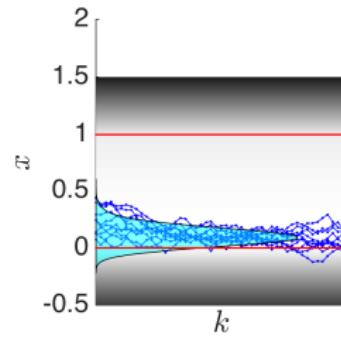
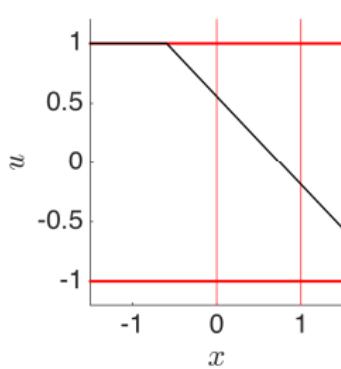
while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$



Deterministic policy gradient

Algorithm: Policy gradient (prototype)

Input: Initial policy parameters θ , step-size $\alpha > 0$

while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

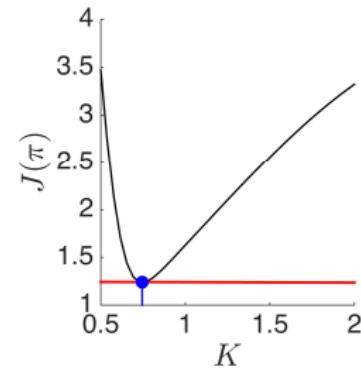
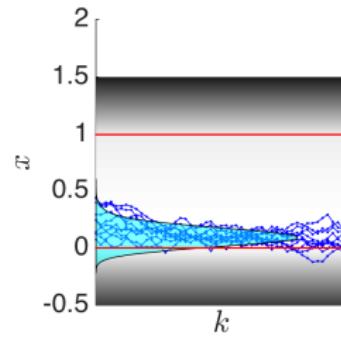
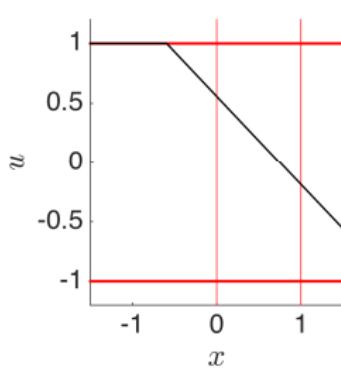
 Run policy $\pi_{\theta} + \mathbf{d}$ for $N \rightarrow \infty$ samples (yields $\tilde{\tau}_{\pi_{\theta}}$)

 Compute $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using e.g. TD-learning, MC, etc.

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau_{\pi_{\theta}}^{\mathbf{d}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta})]$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$

On-the-fly versions
apply

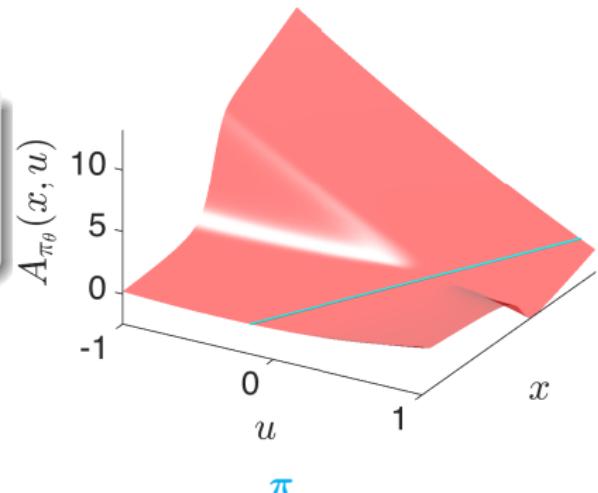


Approximation of Q_{π_θ}

Recall that:

$$Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = A_{\pi_\theta}(\mathbf{x}, \mathbf{u}) + V_{\pi_\theta}(\mathbf{x})$$

$$\text{with } A_{\pi_\theta}(\mathbf{x}, \pi_\theta(\mathbf{x})) = 0 \text{ for all } \mathbf{x}$$



Approximation of Q_{π_θ}

Recall that:

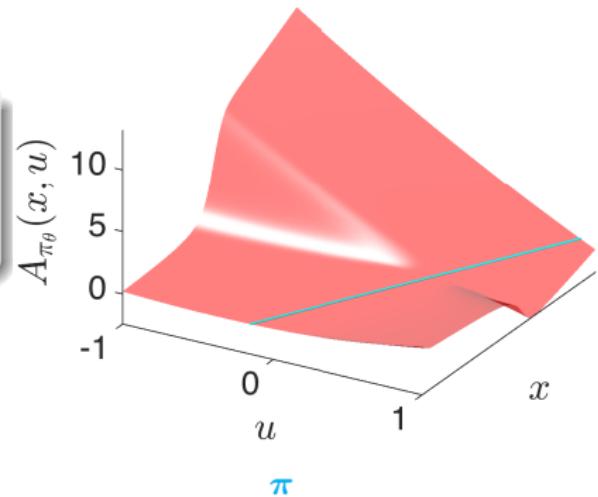
$$Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = A_{\pi_\theta}(\mathbf{x}, \mathbf{u}) + V_{\pi_\theta}(\mathbf{x})$$

with $A_{\pi_\theta}(\mathbf{x}, \pi_\theta(\mathbf{x})) = 0$ for all \mathbf{x}

For $\mathbf{u} \approx \pi_\theta(\mathbf{x})$, approximate A_{π_θ} with:

$$\hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \mathbf{a}(\mathbf{x})^\top (\mathbf{u} - \pi_\theta(\mathbf{x}))$$

for some state-dependent $\mathbf{a}(\mathbf{x})$



Approximation of Q_{π_θ}

Recall that:

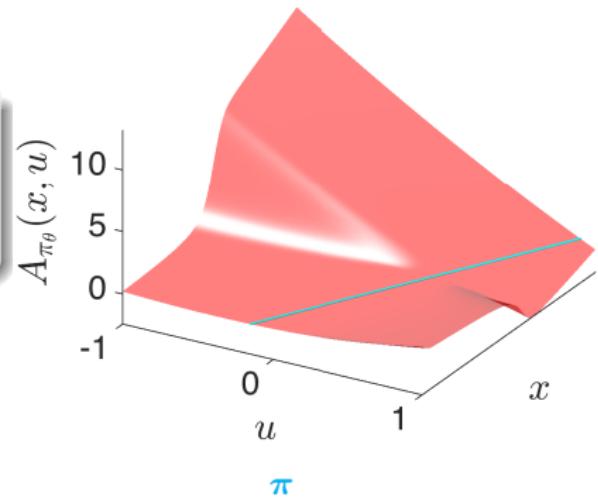
$$Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = A_{\pi_\theta}(\mathbf{x}, \mathbf{u}) + V_{\pi_\theta}(\mathbf{x})$$

with $A_{\pi_\theta}(\mathbf{x}, \pi_\theta(\mathbf{x})) = 0$ for all \mathbf{x}

For $\mathbf{u} \approx \pi_\theta(\mathbf{x})$, approximate A_{π_θ} with:

$$\hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \mathbf{a}(\mathbf{x})^\top (\mathbf{u} - \pi_\theta(\mathbf{x}))$$

for some state-dependent $\mathbf{a}(\mathbf{x})$



Approximation of Q_{π_θ}

Recall that:

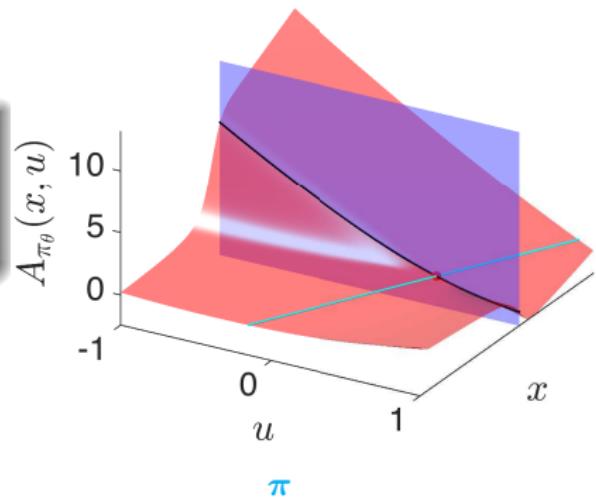
$$Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = A_{\pi_\theta}(\mathbf{x}, \mathbf{u}) + V_{\pi_\theta}(\mathbf{x})$$

$$\text{with } A_{\pi_\theta}(\mathbf{x}, \pi_\theta(\mathbf{x})) = 0 \text{ for all } \mathbf{x}$$

For $\mathbf{u} \approx \pi_\theta(\mathbf{x})$, approximate A_{π_θ} with:

$$\hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \mathbf{a}(\mathbf{x})^\top (\mathbf{u} - \pi_\theta(\mathbf{x}))$$

for some state-dependent $\mathbf{a}(\mathbf{x})$



Approximation of Q_{π_θ}

Recall that:

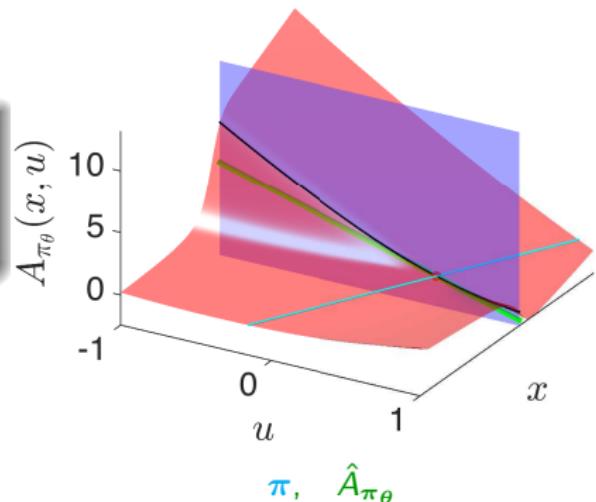
$$Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = A_{\pi_\theta}(\mathbf{x}, \mathbf{u}) + V_{\pi_\theta}(\mathbf{x})$$

$$\text{with } A_{\pi_\theta}(\mathbf{x}, \pi_\theta(\mathbf{x})) = 0 \text{ for all } \mathbf{x}$$

For $\mathbf{u} \approx \pi_\theta(\mathbf{x})$, approximate A_{π_θ} with:

$$\hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \mathbf{a}(\mathbf{x})^\top (\mathbf{u} - \pi_\theta(\mathbf{x}))$$

for some state-dependent $\mathbf{a}(\mathbf{x})$



Approximation of Q_{π_θ}

Recall that:

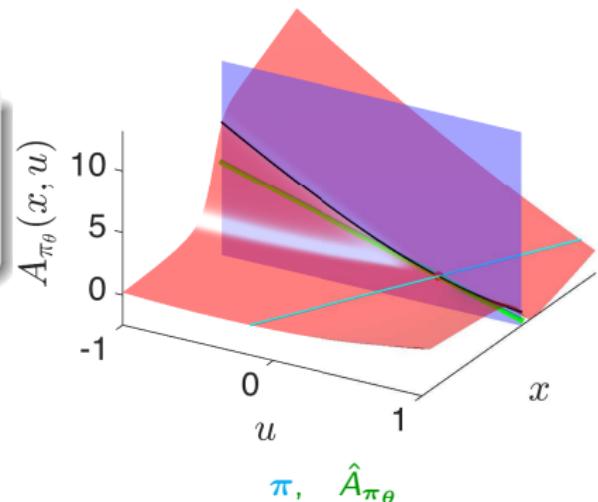
$$Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = A_{\pi_\theta}(\mathbf{x}, \mathbf{u}) + V_{\pi_\theta}(\mathbf{x})$$

$$\text{with } A_{\pi_\theta}(\mathbf{x}, \pi_\theta(\mathbf{x})) = 0 \text{ for all } \mathbf{x}$$

For $\mathbf{u} \approx \pi_\theta(\mathbf{x})$, approximate A_{π_θ} with:

$$\hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \mathbf{a}(\mathbf{x})^\top (\mathbf{u} - \pi_\theta(\mathbf{x}))$$

for some state-dependent $\mathbf{a}(\mathbf{x})$



E.g. **compatible function approximation**

$$\hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \mathbf{w}^\top \nabla_\theta \pi_\theta(\mathbf{x}) (\mathbf{u} - \pi_\theta(\mathbf{x}))$$

for some parameters \mathbf{w} of the size of θ

Approximation of Q_{π_θ}

Recall that:

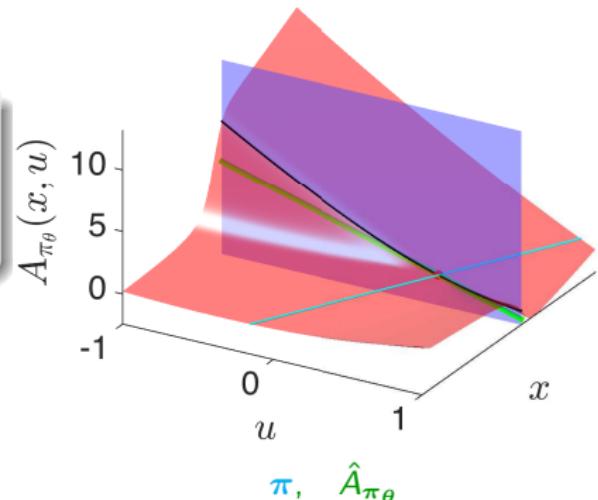
$$Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = A_{\pi_\theta}(\mathbf{x}, \mathbf{u}) + V_{\pi_\theta}(\mathbf{x})$$

$$\text{with } A_{\pi_\theta}(\mathbf{x}, \pi_\theta(\mathbf{x})) = 0 \text{ for all } \mathbf{x}$$

For $\mathbf{u} \approx \pi_\theta(\mathbf{x})$, approximate A_{π_θ} with:

$$\hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \mathbf{a}(\mathbf{x})^\top (\mathbf{u} - \pi_\theta(\mathbf{x}))$$

for some state-dependent $\mathbf{a}(\mathbf{x})$



E.g. **compatible function approximation**

$$\hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \mathbf{w}^\top \nabla_\theta \pi_\theta(\mathbf{x}) (\mathbf{u} - \pi_\theta(\mathbf{x}))$$

for some parameters \mathbf{w} of the size of θ

We then use

$$\hat{Q}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) + \hat{V}_{\pi_\theta}(\mathbf{x})$$

with $\hat{V}_{\pi_\theta}(\mathbf{x}) \approx V_{\pi_\theta}(\mathbf{x})$

Approximation of Q_{π_θ}

Recall that:

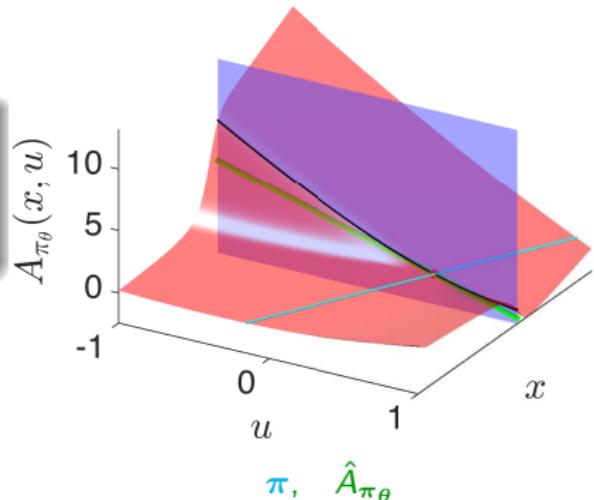
$$Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = A_{\pi_\theta}(\mathbf{x}, \mathbf{u}) + V_{\pi_\theta}(\mathbf{x})$$

with $A_{\pi_\theta}(\mathbf{x}, \pi_\theta(\mathbf{x})) = 0$ for all \mathbf{x}

For $\mathbf{u} \approx \pi_\theta(\mathbf{x})$, approximate A_{π_θ} with:

$$\hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \mathbf{a}(\mathbf{x})^\top (\mathbf{u} - \pi_\theta(\mathbf{x}))$$

for some state-dependent $\mathbf{a}(\mathbf{x})$



E.g. **compatible function approximation**

$$\hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \mathbf{w}^\top \nabla_\theta \pi_\theta(\mathbf{x}) (\mathbf{u} - \pi_\theta(\mathbf{x}))$$

for some parameters \mathbf{w} of the size of θ

We then use

$$\hat{Q}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) = \hat{A}_{\pi_\theta}(\mathbf{x}, \mathbf{u}) + \hat{V}_{\pi_\theta}(\mathbf{x})$$

with $\hat{V}_{\pi_\theta}(\mathbf{x}) \approx V_{\pi_\theta}(\mathbf{x})$

Then if gradients fit, we get exact policy gradients!

$$\mathbf{w} = \arg \min_{\mathbf{w}} \mathbb{E}_{\tau_{\pi_\theta}} \left[\left\| \nabla_{\mathbf{u}} Q_{\pi_\theta} - \nabla_{\mathbf{u}} \hat{Q}_{\pi_\theta} \right\|^2 \right] \Rightarrow \nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau_{\pi_\theta}} \left[\nabla_\theta \pi_\theta \nabla_{\mathbf{u}} \hat{Q}_{\pi_\theta} \mathbf{w} \right]$$

Gradient fitting?

Fitting the gradients

$$\nabla_u Q_{\pi_\theta} \approx \nabla_u \hat{Q}_{\pi_\theta}$$

using data is very difficult. What about fitting $Q_{\pi_\theta} \approx \hat{Q}_{\pi_\theta}$ (classic π evaluation)?

$$\min_w \mathbb{E}_{\tau_{\pi_\theta}} \left[\left\| \nabla_u Q_{\pi_\theta} - \nabla_u \hat{Q}_{\pi_\theta} \right\|^2 \right]$$

vs.

$$\min_w \mathbb{E}_{\tau_{\pi_\theta}} \left[\left\| Q_{\pi_\theta} - \hat{Q}_{\pi_\theta} \right\|^2 \right]$$

- Compatible approximation
- Learning rule?...

- Can use TD or MC-learning
- \neq fitting gradients...

Gradient fitting?

Fitting the gradients

$$\nabla_{\mathbf{u}} Q_{\pi_{\theta}} \approx \nabla_{\mathbf{u}} \hat{Q}_{\pi_{\theta}}$$

using data is very difficult. What about fitting $Q_{\pi_{\theta}} \approx \hat{Q}_{\pi_{\theta}}$ (classic π evaluation)?

$$\min_{\mathbf{w}} \mathbb{E}_{\tau_{\pi_{\theta}}} \left[\left\| \nabla_{\mathbf{u}} Q_{\pi_{\theta}} - \nabla_{\mathbf{u}} \hat{Q}_{\pi_{\theta}} \right\|^2 \right]$$

vs.

$$\min_{\mathbf{w}} \mathbb{E}_{\tau_{\pi_{\theta}}} \left[\left\| Q_{\pi_{\theta}} - \hat{Q}_{\pi_{\theta}} \right\|^2 \right]$$

- Compatible approximation
- Learning rule?...

- Can use TD or MC-learning
- \neq fitting gradients...

Fitting $\hat{Q}_{\pi_{\theta}}$ instead of $\nabla_{\mathbf{u}} \hat{Q}_{\pi_{\theta}}$ yields correct policy gradient if

- Input sequence $\mathbf{u}_k = \pi(\mathbf{x}_k) + \mathbf{d}_k$ where $\mathbf{d}_k \rightarrow 0$ and $\mathbb{E}[\mathbf{d}_k] = 0$

Gradient fitting?

Fitting the gradients

$$\nabla_{\mathbf{u}} Q_{\pi_{\theta}} \approx \nabla_{\mathbf{u}} \hat{Q}_{\pi_{\theta}}$$

using data is very difficult. What about fitting $Q_{\pi_{\theta}} \approx \hat{Q}_{\pi_{\theta}}$ (classic π evaluation)?

$$\min_{\mathbf{w}} \mathbb{E}_{\tau_{\pi_{\theta}}} \left[\left\| \nabla_{\mathbf{u}} Q_{\pi_{\theta}} - \nabla_{\mathbf{u}} \hat{Q}_{\pi_{\theta}} \right\|^2 \right]$$

vs.

$$\min_{\mathbf{w}} \mathbb{E}_{\tau_{\pi_{\theta}}} \left[\left\| Q_{\pi_{\theta}} - \hat{Q}_{\pi_{\theta}} \right\|^2 \right]$$

- Compatible approximation
- Learning rule?...

- Can use TD or MC-learning
- \neq fitting gradients...

Fitting $\hat{Q}_{\pi_{\theta}}$ instead of $\nabla_{\mathbf{u}} \hat{Q}_{\pi_{\theta}}$ yields correct policy gradient if

- Input sequence $\mathbf{u}_k = \pi(\mathbf{x}_k) + \mathbf{d}_k$ where $\mathbf{d}_k \rightarrow 0$ and $\mathbb{E}[\mathbf{d}_k] = 0$
- What about input constraints?!?

Gradient fitting?

Fitting the gradients

$$\nabla_{\mathbf{u}} Q_{\pi_{\theta}} \approx \nabla_{\mathbf{u}} \hat{Q}_{\pi_{\theta}}$$

using data is very difficult. What about fitting $Q_{\pi_{\theta}} \approx \hat{Q}_{\pi_{\theta}}$ (classic π evaluation)?

$$\min_{\mathbf{w}} \mathbb{E}_{\tau_{\pi_{\theta}}} \left[\left\| \nabla_{\mathbf{u}} Q_{\pi_{\theta}} - \nabla_{\mathbf{u}} \hat{Q}_{\pi_{\theta}} \right\|^2 \right]$$

vs.

$$\min_{\mathbf{w}} \mathbb{E}_{\tau_{\pi_{\theta}}} \left[\left\| Q_{\pi_{\theta}} - \hat{Q}_{\pi_{\theta}} \right\|^2 \right]$$

- Compatible approximation
- Learning rule?...

- Can use TD or MC-learning
- \neq fitting gradients...

Fitting $\hat{Q}_{\pi_{\theta}}$ instead of $\nabla_{\mathbf{u}} \hat{Q}_{\pi_{\theta}}$ yields correct policy gradient if

- Input sequence $\mathbf{u}_k = \pi(\mathbf{x}_k) + \mathbf{d}_k$ where $\mathbf{d}_k \rightarrow 0$ and $\mathbb{E}[\mathbf{d}_k] = 0$
- What about input constraints?!? Fine if \mathbf{d}_k is in the “border” of the feasible set and centred in the null space of the constraints

Outline

1 Introduction

2 Deterministic policy gradient

3 Stochastic policy gradient



Deterministic policy

$$\mathbf{u}_k = \pi(\mathbf{x}_k)$$

Policy gradient

$$\nabla_{\theta} J = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta} \nabla_{\mathbf{u}} Q_{\pi_{\theta}}]$$

what if inputs are discrete?

Stochastic policy

Deterministic policy

$$\mathbf{u}_k = \pi(\mathbf{x}_k)$$

Policy gradient

$$\nabla_{\theta} J = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta} \nabla_{\mathbf{u}} Q_{\pi_{\theta}}]$$

what if inputs are discrete?

Stochastic policy

$$\pi[\mathbf{u}_k | \mathbf{x}_k]$$

gives input distribution for a given state \mathbf{x}_k

Policy gradient?

Stochastic policy

Deterministic policy

$$\mathbf{u}_k = \pi(\mathbf{x}_k)$$

Policy gradient

$$\nabla_{\theta} J = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta} \nabla_{\mathbf{u}} Q_{\pi_{\theta}}]$$

what if inputs are discrete?

Stochastic policy

$$\pi[\mathbf{u}_k | \mathbf{x}_k]$$

gives input distribution for a given state \mathbf{x}_k

Policy gradient?

E.g. Gaussian stochastic policy:

$$\pi_{\theta}[\mathbf{u}_k | \mathbf{x}_k] = \mathcal{N}(\bar{\pi}_{\theta}(\mathbf{x}_k), \Sigma)$$

restricted to feasible \mathbf{u}

Covariance Σ can be part of policy parameters θ .

Stochastic policy

Deterministic policy

$$\mathbf{u}_k = \pi(\mathbf{x}_k)$$

Policy gradient

$$\nabla_{\theta} J = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta} \nabla_{\mathbf{u}} Q_{\pi_{\theta}}]$$

what if inputs are discrete?

Stochastic policy

$$\pi[\mathbf{u}_k | \mathbf{x}_k]$$

gives input distribution for a given state \mathbf{x}_k

Policy gradient?

E.g. Gaussian stochastic policy:

$$\pi_{\theta}[\mathbf{u}_k | \mathbf{x}_k] = \mathcal{N}(\bar{\pi}_{\theta}(\mathbf{x}_k), \Sigma)$$

restricted to feasible \mathbf{u}

Covariance Σ can be part of policy parameters θ .

$$\bar{\pi}_{\theta}(\mathbf{x}_k), \pi_{\theta}[\mathbf{u}_k | \mathbf{x}_k]$$

Stochastic policy

Deterministic policy

$$\mathbf{u}_k = \pi(\mathbf{x}_k)$$

Policy gradient

$$\nabla_{\theta} J = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta} \nabla_{\mathbf{u}} Q_{\pi_{\theta}}]$$

what if inputs are discrete?

Stochastic policy

$$\pi[\mathbf{u}_k | \mathbf{x}_k]$$

gives input distribution for a given state \mathbf{x}_k

Policy gradient?

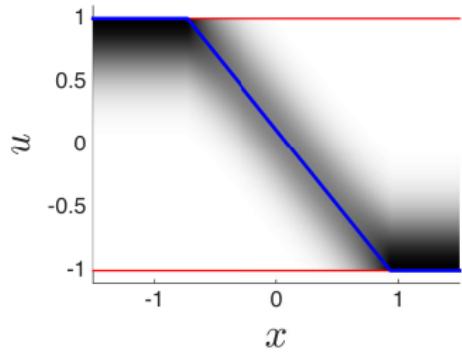
E.g. Gaussian stochastic policy:

$$\pi_{\theta}[\mathbf{u}_k | \mathbf{x}_k] = \mathcal{N}(\bar{\pi}_{\theta}(\mathbf{x}_k), \Sigma)$$

restricted to feasible \mathbf{u}

Covariance Σ can be part of policy parameters θ .

$$\bar{\pi}_{\theta}(\mathbf{x}_k), \pi_{\theta}[\mathbf{u}_k | \mathbf{x}_k]$$



Bellman equation with stochastic policy?

Bellman equations

$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

$$V_{\pi}(x) = Q_{\pi}(x, \pi(x))$$

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

become...

Bellman equation with stochastic policy?

Bellman equations

$$V_\pi(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_\pi(x_+) | x, \pi]$$

$$V_\pi(x) = Q_\pi(x, \pi(x))$$

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

become...

$$V_\pi(x) = \mathbb{E}[L(x, u) + \gamma V_\pi(x_+) | x, u \sim \pi]$$

Bellman equation with stochastic policy?

Bellman equations

$$V_\pi(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_\pi(x_+) | x, \pi]$$

$$V_\pi(x) = Q_\pi(x, \pi(x))$$

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

become...

$$V_\pi(x) = \mathbb{E}[L(x, u) + \gamma V_\pi(x_+) | x, u \sim \pi]$$

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

Bellman equation with stochastic policy?

Bellman equations

$$V_\pi(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_\pi(x_+) | x, \pi]$$

$$V_\pi(x) = Q_\pi(x, \pi(x))$$

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

become...

$$V_\pi(x) = \mathbb{E}[L(x, u) + \gamma V_\pi(x_+) | x, u \sim \pi]$$

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

$$V_\pi(x) = \mathbb{E}[Q_\pi(x, u) | u \sim \pi]$$

Bellman equation with stochastic policy?

Bellman equations

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

become...

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

- ... but are otherwise still valid for $\pi(u|x)$

Bellman equation with stochastic policy?

Bellman equations

$$V_\pi(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_\pi(x_+) | x, \pi]$$

$$V_\pi(x) = Q_\pi(x, \pi(x))$$

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

become...

$$V_\pi(x) = \mathbb{E}[L(x, u) + \gamma V_\pi(x_+) | x, u \sim \pi]$$

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

$$V_\pi(x) = \mathbb{E}[Q_\pi(x, u) | u \sim \pi]$$

- ... but are otherwise still valid for $\pi[u | x]$
- Expected values $\mathbb{E}[\cdot]$ over state transitions are taken as:

$$\mathbb{E}[\zeta | x, \pi] = \int \zeta \cdot \mathbb{P}[x_+ | x, u] \color{red}{\pi[u | x]} dx_+ du$$

Score function

$$\begin{aligned}\nabla_{\theta} \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] &= \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] \frac{\nabla_{\theta} \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]}{\pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]} \\ &= \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] \underbrace{\nabla_{\theta} \log \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]}_{\text{Score function}}\end{aligned}$$

Score function

$$\begin{aligned}\nabla_{\theta} \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] &= \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] \frac{\nabla_{\theta} \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]}{\pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]} \\ &= \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] \underbrace{\nabla_{\theta} \log \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]}_{\text{Score function}}\end{aligned}$$

E.g. Gaussian stochastic policy:

$$\pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] = \frac{1}{\sqrt{2\pi \det(\Sigma)}} e^{-\frac{1}{2}(\mathbf{u}_k - \bar{\pi}_{\theta}(\mathbf{x}_k))^{\top} \Sigma^{-1} (\mathbf{u}_k - \bar{\pi}_{\theta}(\mathbf{x}_k))}$$

has score function: $\nabla_{\theta} \log \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] = (\mathbf{u}_k - \bar{\pi}_{\theta}(\mathbf{x}_k))^{\top} \Sigma^{-1} \nabla_{\theta} \bar{\pi}_{\theta}(\mathbf{x}_k)$

Score function

$$\begin{aligned}\nabla_{\theta} \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] &= \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] \frac{\nabla_{\theta} \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]}{\pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]} \\ &= \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] \underbrace{\nabla_{\theta} \log \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]}_{\text{Score function}}\end{aligned}$$

E.g. Gaussian stochastic policy:

$$\pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] = \frac{1}{\sqrt{2\pi \det(\Sigma)}} e^{-\frac{1}{2}(\mathbf{u}_k - \bar{\pi}_{\theta}(\mathbf{x}_k))^{\top} \Sigma^{-1} (\mathbf{u}_k - \bar{\pi}_{\theta}(\mathbf{x}_k))}$$

has score function: $\nabla_{\theta} \log \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] = (\mathbf{u}_k - \bar{\pi}_{\theta}(\mathbf{x}_k))^{\top} \Sigma^{-1} \nabla_{\theta} \bar{\pi}_{\theta}(\mathbf{x}_k)$

Why is the score function useful?

$$\nabla_{\theta} \mathbb{E}_{\pi_{\theta}} [\varphi(\mathbf{x}, \mathbf{u})] = \nabla_{\theta} \int \pi_{\theta} [\mathbf{u} | \mathbf{x}] \varphi(\mathbf{x}, \mathbf{u}) d\mathbf{u} = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [\mathbf{u} | \mathbf{x}] \varphi(\mathbf{x}, \mathbf{u})]$$

conversion $\nabla_{\theta} \mathbb{E}_{\pi_{\theta}} [.] \rightarrow \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} .]$

Score function

$$\begin{aligned}\nabla_{\theta} \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] &= \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] \frac{\nabla_{\theta} \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]}{\pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]} \\ &= \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] \underbrace{\nabla_{\theta} \log \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]}_{\text{Score function}}\end{aligned}$$

Score \equiv relative change:

$$\nabla_{\theta} \log \pi_{\theta} = \frac{\nabla_{\theta} \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]}{\pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k]}$$

E.g. Gaussian stochastic policy:

$$\pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] = \frac{1}{\sqrt{2\pi \det(\Sigma)}} e^{-\frac{1}{2}(\mathbf{u}_k - \bar{\pi}_{\theta}(\mathbf{x}_k))^{\top} \Sigma^{-1} (\mathbf{u}_k - \bar{\pi}_{\theta}(\mathbf{x}_k))}$$

has score function: $\nabla_{\theta} \log \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] = (\mathbf{u}_k - \bar{\pi}_{\theta}(\mathbf{x}_k))^{\top} \Sigma^{-1} \nabla_{\theta} \bar{\pi}_{\theta}(\mathbf{x}_k)$

Why is the score function useful?

$$\nabla_{\theta} \mathbb{E}_{\pi_{\theta}} [\varphi(\mathbf{x}, \mathbf{u})] = \nabla_{\theta} \int \pi_{\theta} [\mathbf{u} | \mathbf{x}] \varphi(\mathbf{x}, \mathbf{u}) d\mathbf{u} = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [\mathbf{u} | \mathbf{x}] \varphi(\mathbf{x}, \mathbf{u})]$$

conversion $\nabla_{\theta} \mathbb{E}_{\pi_{\theta}} [.] \rightarrow \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} .]$

Score function & policy gradient

Deterministic policy $\pi_\theta(\mathbf{x}_k)$ performance

$$J(\pi_\theta) = \sum_{k=0}^{\infty} \int \gamma^k L(\mathbf{x}_k, \pi_\theta(\mathbf{x}_k)) \prod_{i=0}^{k-1} \mathbb{P}[\mathbf{x}_{i+1} | \mathbf{x}_i, \pi_\theta(\mathbf{x}_i)] \mathbb{P}[\mathbf{x}_0] \cdot d\mathbf{x}_0 \dots d\mathbf{x}_{k-1}$$

Stochastic policy $\pi_\theta[\mathbf{u}_k | \mathbf{x}_k]$ performance

$$J(\pi_\theta) =$$

$$= \sum_{k=0}^{\infty} \int \gamma^k \pi_\theta[\mathbf{u}_k | \mathbf{x}_k] L(\mathbf{x}_k, \mathbf{u}_k) \prod_{i=0}^{k-1} \pi_\theta[\mathbf{u}_i | \mathbf{x}_i] \mathbb{P}[\mathbf{x}_{i+1} | \mathbf{x}_i, \mathbf{u}_i] \mathbb{P}[\mathbf{x}_0] d\mathbf{x}_{0,\dots,k-1} d\mathbf{u}_{0,\dots,k-1}$$

Score function & policy gradient

Deterministic policy $\pi_\theta(\mathbf{x}_k)$ performance

$$J(\pi_\theta) = \sum_{k=0}^{\infty} \int \gamma^k L(\mathbf{x}_k, \pi_\theta(\mathbf{x}_k)) \prod_{i=0}^{k-1} \mathbb{P}[\mathbf{x}_{i+1} | \mathbf{x}_i, \pi_\theta(\mathbf{x}_i)] \mathbb{P}[\mathbf{x}_0] \cdot d\mathbf{x}_0 \dots d\mathbf{x}_{k-1}$$

Stochastic policy $\pi_\theta [\mathbf{u}_k | \mathbf{x}_k]$ performance

$$\begin{aligned} J(\pi_\theta) &= \\ &= \sum_{k=0}^{\infty} \int \gamma^k \pi_\theta [\mathbf{u}_k | \mathbf{x}_k] L(\mathbf{x}_k, \mathbf{u}_k) \prod_{i=0}^{k-1} \pi_\theta [\mathbf{u}_i | \mathbf{x}_i] \mathbb{P}[\mathbf{x}_{i+1} | \mathbf{x}_i, \mathbf{u}_i] \mathbb{P}[\mathbf{x}_0] d\mathbf{x}_{0,\dots,k-1} d\mathbf{u}_{0,\dots,k-1} \end{aligned}$$

Gradient of stochastic policy performance

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] Q_{\pi_\theta} (\mathbf{x}, \mathbf{u})]$$

i.e. correlation between score function and action-value function

Policy gradient deterministic vs. stochastic

Gradient of stochastic policy

$$\begin{aligned}\nabla_{\theta} J(\pi_{\theta}) \\ = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} (\mathbf{u} | \mathbf{x}) Q_{\pi_{\theta}} (\mathbf{x}, \mathbf{u})]\end{aligned}$$

Gradient of deterministic policy

$$\begin{aligned}\nabla_{\theta} J(\bar{\pi}_{\theta}) \\ = \mathbb{E}_{\bar{\pi}_{\theta}} [\nabla_{\theta} \bar{\pi}_{\theta} (\mathbf{x}) \nabla_{\mathbf{u}} Q_{\bar{\pi}_{\theta}} (\mathbf{x}, \mathbf{u})]\end{aligned}$$

Policy gradient deterministic vs. stochastic

Gradient of stochastic policy

$$\begin{aligned}\nabla_{\theta} J(\pi_{\theta}) \\ = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} (\mathbf{u} | \mathbf{x}) Q_{\pi_{\theta}} (\mathbf{x}, \mathbf{u})]\end{aligned}$$

Gradient of deterministic policy

$$\begin{aligned}\nabla_{\theta} J(\bar{\pi}_{\theta}) \\ = \mathbb{E}_{\bar{\pi}_{\theta}} [\nabla_{\theta} \bar{\pi}_{\theta} (\mathbf{x}) \nabla_{\mathbf{u}} Q_{\bar{\pi}_{\theta}} (\mathbf{x}, \mathbf{u})]\end{aligned}$$

Let Σ be the covariance of $\pi_{\theta} [\mathbf{u} | \mathbf{x}]$ then as stochastic \rightarrow deterministic...

$$\lim_{\Sigma \rightarrow 0} \pi_{\theta} [\mathbf{u} | \mathbf{x}] = \delta(\mathbf{u} - \bar{\pi}_{\theta} (\mathbf{x}))$$

Policy gradient deterministic vs. stochastic

Gradient of stochastic policy

$$\begin{aligned}\nabla_{\theta} J(\pi_{\theta}) \\ = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} (\mathbf{u} | \mathbf{x}) Q_{\pi_{\theta}} (\mathbf{x}, \mathbf{u})]\end{aligned}$$

Gradient of deterministic policy

$$\begin{aligned}\nabla_{\theta} J(\bar{\pi}_{\theta}) \\ = \mathbb{E}_{\bar{\pi}_{\theta}} [\nabla_{\theta} \bar{\pi}_{\theta} (\mathbf{x}) \nabla_{\mathbf{u}} Q_{\bar{\pi}_{\theta}} (\mathbf{x}, \mathbf{u})]\end{aligned}$$

Let Σ be the covariance of $\pi_{\theta} [\mathbf{u} | \mathbf{x}]$ then as stochastic \rightarrow deterministic...

$$\lim_{\Sigma \rightarrow 0} \pi_{\theta} [\mathbf{u} | \mathbf{x}] = \delta(\mathbf{u} - \bar{\pi}_{\theta} (\mathbf{x}))$$

...policy gradients match:

$$\lim_{\Sigma \rightarrow 0} \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} (\mathbf{u} | \mathbf{x}) Q_{\pi_{\theta}} (\mathbf{x}, \mathbf{u})] = \mathbb{E}_{\bar{\pi}_{\theta}} [\nabla_{\theta} \bar{\pi}_{\theta} (\mathbf{x}) \nabla_{\mathbf{u}} Q_{\bar{\pi}_{\theta}} (\mathbf{x}, \mathbf{u})]$$

...under some regularity assumptions

Interpretation of stochastic policy gradient?

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)]$$

What are the terms $\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)$?

Interpretation of stochastic policy gradient?

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)]$$

What are the terms $\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)$?

- $Q_{\pi_{\theta}}$ is scalar, acts like a “weight”

Interpretation of stochastic policy gradient?

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)]$$

What are the terms $\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)$?

- $Q_{\pi_{\theta}}$ is scalar, acts like a “weight”
- Change θ in “direction” $-\nabla_{\theta} \log \pi_{\theta} [u | x]$...
...decreases the probability of drawing u in state x

Interpretation of stochastic policy gradient?

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [\mathbf{u} | \mathbf{x}] Q_{\pi_{\theta}} (\mathbf{x}, \mathbf{u})]$$

What are the terms $\nabla_{\theta} \log \pi_{\theta} [\mathbf{u} | \mathbf{x}] Q_{\pi_{\theta}} (\mathbf{x}, \mathbf{u})$?

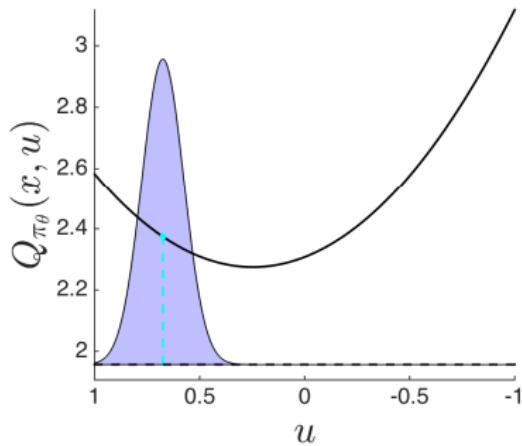
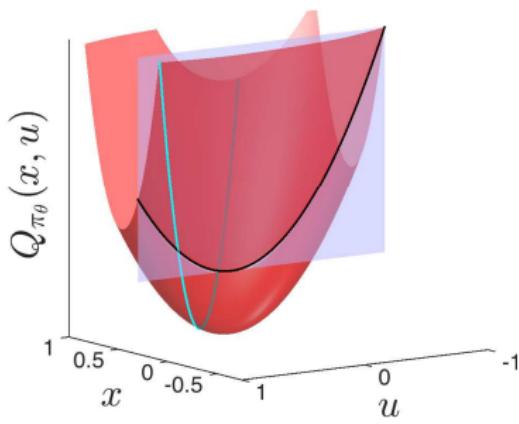
- $Q_{\pi_{\theta}}$ is scalar, acts like a “weight”
- Change θ in “direction” $-\nabla_{\theta} \log \pi_{\theta} [\mathbf{u} | \mathbf{x}] \dots$
...decreases the probability of drawing \mathbf{u} in state \mathbf{x}
- Consider $\pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] = \mathcal{N}(\bar{\pi}_{\theta} (\mathbf{x}_k), \Sigma)$

Interpretation of stochastic policy gradient?

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)]$$

What are the terms $\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)$?

- $Q_{\pi_{\theta}}$ is scalar, acts like a “weight”
- Change θ in “direction” $-\nabla_{\theta} \log \pi_{\theta} [u | x]$...
...decreases the probability of drawing u in state x

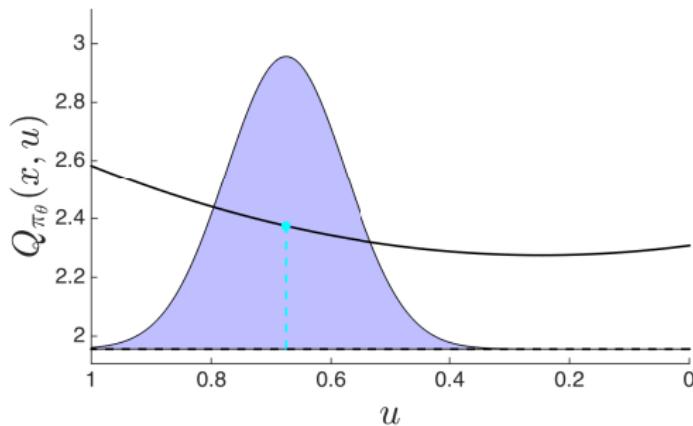


Interpretation of stochastic policy gradient?

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)]$$

What are the terms $\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)$?

- $Q_{\pi_{\theta}}$ is scalar, acts like a “weight”
- Change θ in “direction” $-\nabla_{\theta} \log \pi_{\theta} [u | x]$...
...decreases the probability of drawing u in state x
- Consider $\pi_{\theta} [u_k | x_k] = \mathcal{N}(\bar{\pi}_{\theta}(x_k), \Sigma)$

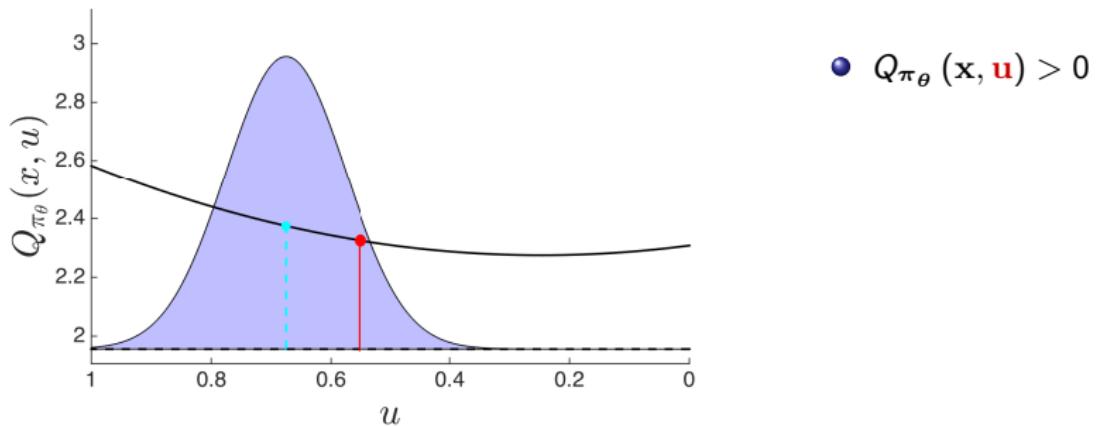


Interpretation of stochastic policy gradient?

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)]$$

What are the terms $\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)$?

- $Q_{\pi_{\theta}}$ is scalar, acts like a “weight”
- Change θ in “direction” $-\nabla_{\theta} \log \pi_{\theta} [u | x]$...
...decreases the probability of drawing u in state x
- Consider $\pi_{\theta} [u_k | x_k] = \mathcal{N}(\bar{\pi}_{\theta}(x_k), \Sigma)$

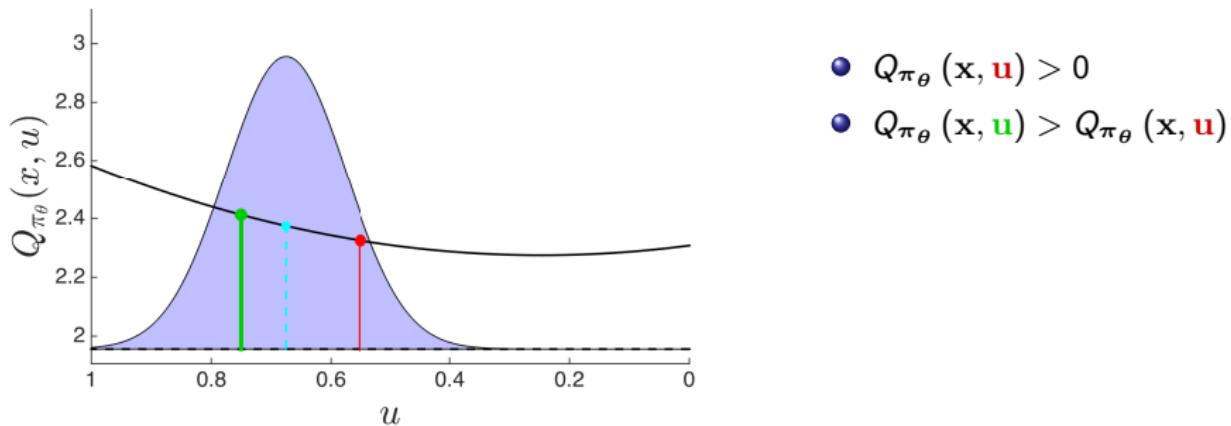


Interpretation of stochastic policy gradient?

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)]$$

What are the terms $\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)$?

- $Q_{\pi_{\theta}}$ is scalar, acts like a “weight”
- Change θ in “direction” $-\nabla_{\theta} \log \pi_{\theta} [u | x]$...
...decreases the probability of drawing u in state x
- Consider $\pi_{\theta} [u_k | x_k] = \mathcal{N}(\bar{\pi}_{\theta}(x_k), \Sigma)$

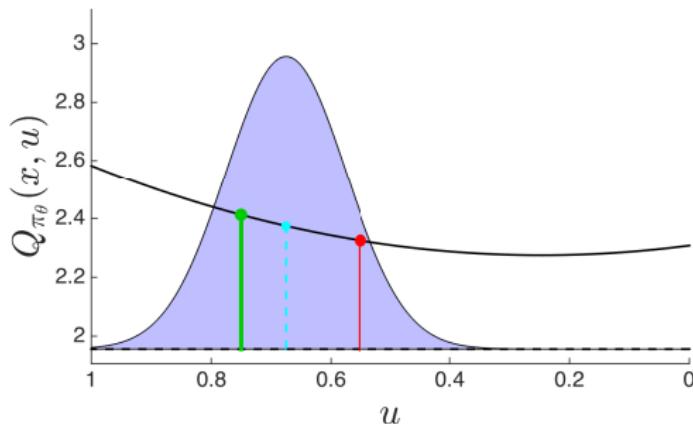


Interpretation of stochastic policy gradient?

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)]$$

What are the terms $\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)$?

- $Q_{\pi_{\theta}}$ is scalar, acts like a “weight”
- Change θ in “direction” $-\nabla_{\theta} \log \pi_{\theta} [u | x]$...
...decreases the probability of drawing u in state x
- Consider $\pi_{\theta} [u_k | x_k] = \mathcal{N}(\bar{\pi}_{\theta}(x_k), \Sigma)$



- $Q_{\pi_{\theta}} (x, u) > 0$
- $Q_{\pi_{\theta}} (x, u) > Q_{\pi_{\theta}} (x, \bar{u})$

Step in direction

$$-\nabla_{\theta} \log \pi_{\theta} Q_{\pi_{\theta}}$$

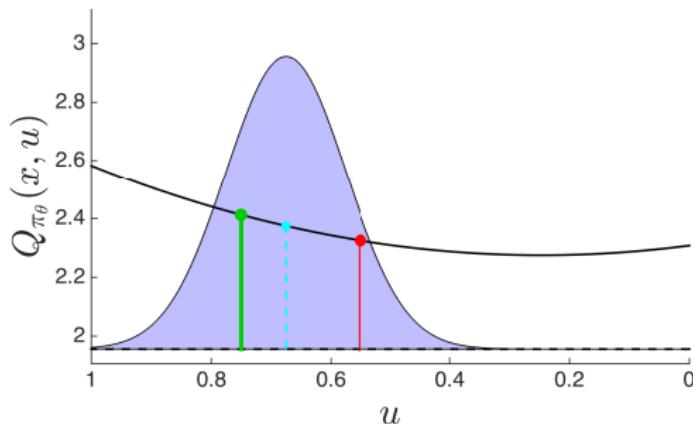
“pushes” $\bar{\pi}_{\theta}$ away

Interpretation of stochastic policy gradient?

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)]$$

What are the terms $\nabla_{\theta} \log \pi_{\theta} [u | x] Q_{\pi_{\theta}} (x, u)$?

- $Q_{\pi_{\theta}}$ is scalar, acts like a “weight”
- Change θ in “direction” $-\nabla_{\theta} \log \pi_{\theta} [u | x]$...
...decreases the probability of drawing u in state x
- Consider $\pi_{\theta} [u_k | x_k] = \mathcal{N}(\bar{\pi}_{\theta}(x_k), \Sigma)$



- $Q_{\pi_{\theta}} (x, u) > 0$
- $Q_{\pi_{\theta}} (x, u) > Q_{\pi_{\theta}} (x, \bar{u})$

Step in direction

$$-\nabla_{\theta} \log \pi_{\theta} Q_{\pi_{\theta}}$$

“pushes” $\bar{\pi}_{\theta}$ away

Both u and \bar{u} “push” $\bar{\pi}_{\theta}$ away from them, u “pushes” harder, $\bar{\pi}_{\theta}$ in average moves to the right

Stochastic policy gradient - Advantage critic

Baseline:

$$\mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] B(\mathbf{x})] = 0$$

holds for any baseline function B

Key idea: use baseline to prevent “everyone from pushing”

Stochastic policy gradient - Advantage critic

Baseline:

$$\mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] B(\mathbf{x})] = 0$$

holds for any baseline function B

Key idea: use baseline to prevent “everyone from pushing”

$$\begin{aligned}\nabla_\theta J(\pi_\theta) &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] Q_{\pi_\theta} (\mathbf{x}, \mathbf{u})] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] (Q_{\pi_\theta} (\mathbf{x}, \mathbf{u}) - V_{\pi_\theta} (\mathbf{x}))] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] A_{\pi_\theta} (\mathbf{x}, \mathbf{u})]\end{aligned}$$

A_{π_θ} is the Advantage function

Stochastic policy gradient - Advantage critic

Baseline:

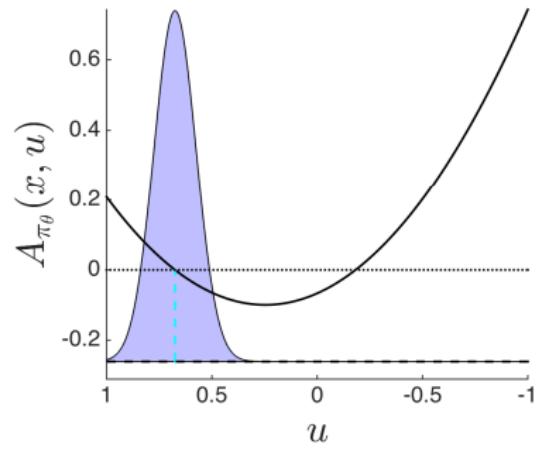
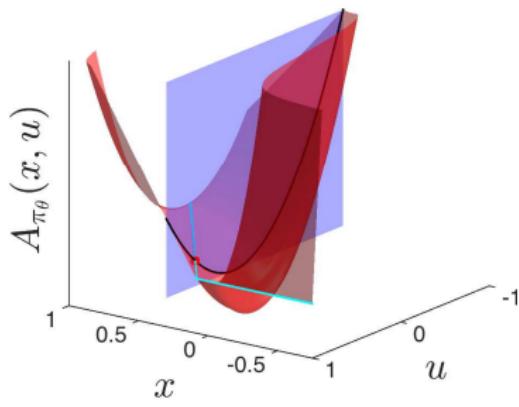
$$\mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] B(\mathbf{x})] = 0$$

holds for any baseline function B

Key idea: use baseline to prevent “everyone from pushing”

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] A_{\pi_\theta} (\mathbf{x}, \mathbf{u})]$$

A_{π_θ} is the Advantage function



Stochastic policy gradient - Advantage critic

Baseline:

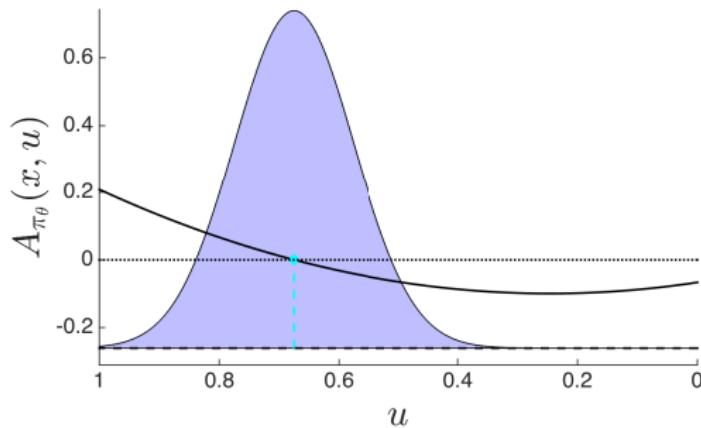
$$\mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] B(\mathbf{x})] = 0$$

holds for any baseline function B

Key idea: use baseline to prevent “everyone from pushing”

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] A_{\pi_\theta} (\mathbf{x}, \mathbf{u})]$$

A_{π_θ} is the Advantage function



Step towards $-\nabla_\theta \log \pi_\theta A_{\pi_\theta}$

Stochastic policy gradient - Advantage critic

Baseline:

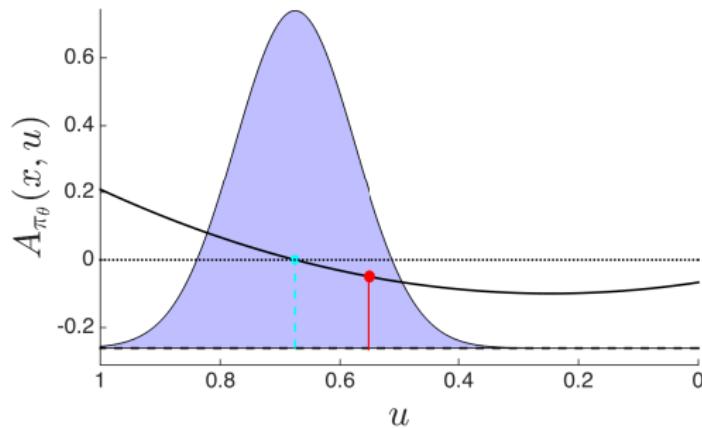
$$\mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] B(\mathbf{x})] = 0$$

holds for any baseline function B

Key idea: use baseline to prevent “everyone from pushing”

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] A_{\pi_\theta} (\mathbf{x}, \mathbf{u})]$$

A_{π_θ} is the Advantage function



Step towards $-\nabla_\theta \log \pi_\theta A_{\pi_\theta}$

- $A_{\pi_\theta} (\mathbf{x}, \mathbf{u}) < 0$ pulls

Stochastic policy gradient - Advantage critic

Baseline:

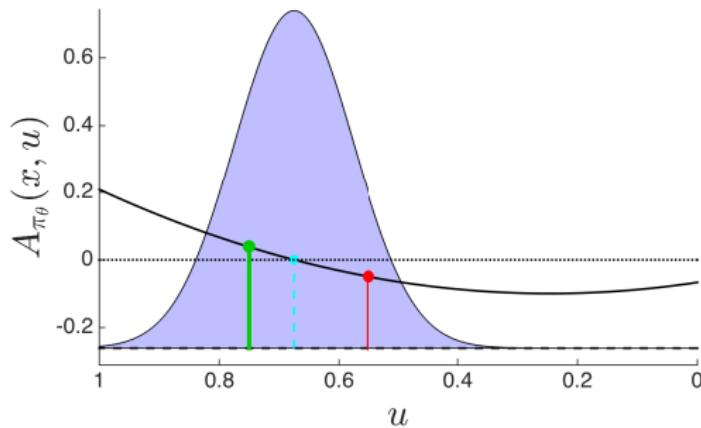
$$\mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] B(\mathbf{x})] = 0$$

holds for any baseline function B

Key idea: use baseline to prevent “everyone from pushing”

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] A_{\pi_\theta} (\mathbf{x}, \mathbf{u})]$$

A_{π_θ} is the Advantage function



Step towards $-\nabla_\theta \log \pi_\theta A_{\pi_\theta}$

- $A_{\pi_\theta} (\mathbf{x}, \mathbf{u}) < 0$ pulls
- $A_{\pi_\theta} (\mathbf{x}, \mathbf{u}) > 0$ pushes

Stochastic policy gradient - Advantage critic

Baseline:

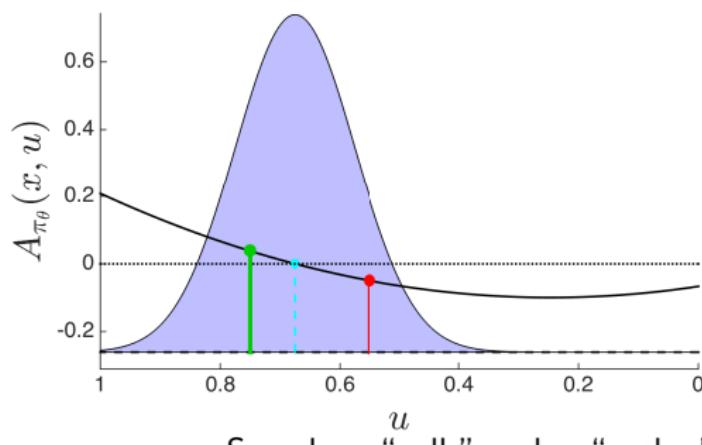
$$\mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] B(\mathbf{x})] = 0$$

holds for any baseline function B

Key idea: use baseline to prevent “everyone from pushing”

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [\mathbf{u} | \mathbf{x}] A_{\pi_\theta} (\mathbf{x}, \mathbf{u})]$$

A_{π_θ} is the Advantage function



Sample \mathbf{u} “pulls” and \mathbf{u} “pushes” π_θ , drive it to $A_{\pi_\theta} = 0$

Step towards $-\nabla_\theta \log \pi_\theta A_{\pi_\theta}$

- $A_{\pi_\theta} (\mathbf{x}, \mathbf{u}) < 0$ pulls
- $A_{\pi_\theta} (\mathbf{x}, \mathbf{u}) > 0$ pushes

Stochastic policy gradient

Algorithm: RL using stochastic π -gradient

Input: Initial policy parameters θ , step-size $\alpha > 0$

while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

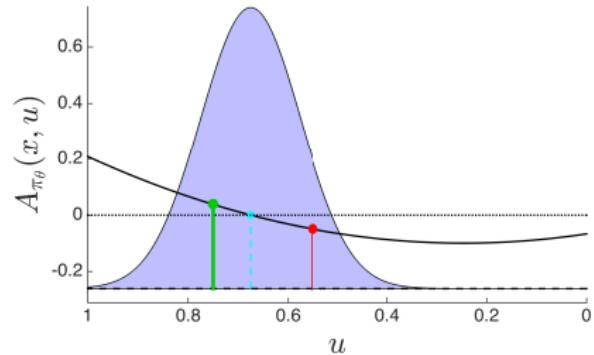
 Run policy π_{θ} for $N \rightarrow \infty$ samples (yields $\mathbf{x}_0, \mathbf{u}_0, \dots, \mathbf{x}_{N-1}, \mathbf{u}_{N-1}$)

 Evaluate $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using (LS)TD-learning, MC-like, ...

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \frac{1}{N} \sum_{k=0}^{N-1} \nabla_{\theta} \log \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] Q_{\pi_{\theta}} (\mathbf{x}_k, \mathbf{u}_k)$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$

return Optimized policy parameters θ



Stochastic policy gradient

Algorithm: RL using stochastic π -gradient

Input: Initial policy parameters θ , step-size $\alpha > 0$

while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy π_{θ} for $N \rightarrow \infty$ samples (yields $\mathbf{x}_0, \mathbf{u}_0, \dots, \mathbf{x}_{N-1}, \mathbf{u}_{N-1}$)

 Evaluate $Q_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using (LS)TD-learning, MC-like, ...

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \frac{1}{N} \sum_{k=0}^{N-1} \nabla_{\theta} \log \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] Q_{\pi_{\theta}} (\mathbf{x}_k, \mathbf{u}_k)$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$

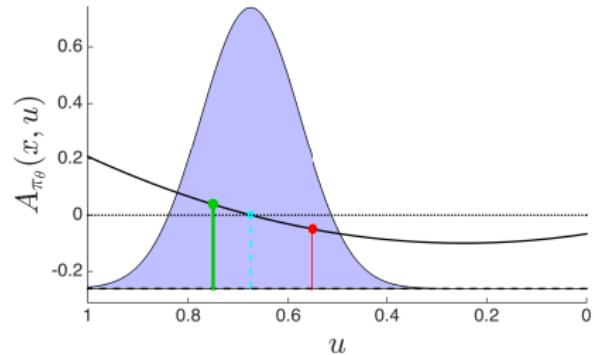
return Optimized policy parameters θ

In practice, $Q_{\pi_{\theta}}$ approximated:

$$Q_w \approx Q_{\pi_{\theta}}$$

Labelled **Actor-Critic**:

- π_{θ} is the “actor”
- Q_w is the “critic”



Stochastic policy gradient

Algorithm: RL using stochastic π -gradient

Input: Initial policy parameters θ , step-size $\alpha > 0$

while $\|\nabla_{\theta} J(\pi_{\theta})\| > \text{Tol}$ **do**

 Run policy π_{θ} for $N \rightarrow \infty$ samples (yields $\mathbf{x}_0, \mathbf{u}_0, \dots, \mathbf{x}_{N-1}, \mathbf{u}_{N-1}$)

 Evaluate $A_{\pi_{\theta}}(\mathbf{x}, \mathbf{u})$ using (LS)TD-learning, MC-like, ...

 Compute $\nabla_{\theta} J(\pi_{\theta}) = \frac{1}{N} \sum_{k=0}^{N-1} \nabla_{\theta} \log \pi_{\theta} [\mathbf{u}_k | \mathbf{x}_k] A_{\pi_{\theta}} (\mathbf{x}_k, \mathbf{u}_k)$

 Gradient step: $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta})$

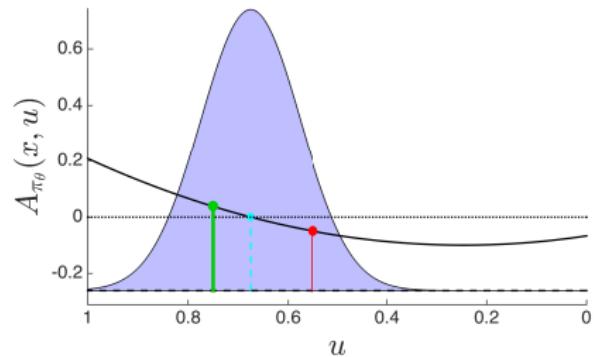
return Optimized policy parameters θ

In practice, $Q_{\pi_{\theta}}$ approximated:

$$Q_w \approx Q_{\pi_{\theta}}$$

Labelled **Actor-Critic**:

- π_{θ} is the “actor”
- Q_w is the “critic”



Stochastic policy gradient - Compatible function approximation

What is the impact of approximation $Q_w \approx Q_{\pi_\theta}$?

If the following conditions hold:

- $\nabla_w Q_w = \nabla_\theta \log \pi_\theta$
- $w = \min_w \mathbb{E}_{\pi_\theta} [(Q_w - Q_{\pi_\theta})^2]$

then

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta [u | x] Q_w(x, u)]$$

What about $A_w \approx Q_{\pi_\theta} - V_{\pi_\theta}$? Same story if

$$A_w = Q_w - \hat{V}$$

where Q_w is compatible and \hat{V} evaluated independently of w

Stochastic policy gradient - TD critic

Consider TD error:

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

We observe that:

$$\begin{aligned}\mathbb{E}[\delta | \mathbf{x}, \mathbf{u}] &= \mathbb{E}[L(\mathbf{x}, \mathbf{u}) + \gamma V_{\pi_\theta}(\mathbf{x}_+)] - V_{\pi_\theta}(\mathbf{x}) \\ &= Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) - V_{\pi_\theta}(\mathbf{x}) \\ &= A_{\pi_\theta}(\mathbf{x}, \mathbf{u})\end{aligned}$$

Stochastic policy gradient - TD critic

Consider TD error:

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

We observe that:

$$\begin{aligned}\mathbb{E}[\delta | \mathbf{x}, \mathbf{u}] &= \mathbb{E}[L(\mathbf{x}, \mathbf{u}) + \gamma V_{\pi_\theta}(\mathbf{x}_+) - V_{\pi_\theta}(\mathbf{x})] \\ &= Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) - V_{\pi_\theta}(\mathbf{x}) \\ &= A_{\pi_\theta}(\mathbf{x}, \mathbf{u})\end{aligned}$$

Hence we have three possible critics:

$$\begin{aligned}\nabla_\theta J(\pi_\theta) &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) A_{\pi_\theta}(\mathbf{x}, \mathbf{u})] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) \delta(\mathbf{x}, \mathbf{u}, \mathbf{x}_+)] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) Q_{\pi_\theta}(\mathbf{x}, \mathbf{u})]\end{aligned}$$

Stochastic policy gradient - TD critic

Consider TD error:

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

We observe that:

$$\begin{aligned}\mathbb{E}[\delta | \mathbf{x}, \mathbf{u}] &= \mathbb{E}[L(\mathbf{x}, \mathbf{u}) + \gamma V_{\pi_\theta}(\mathbf{x}_+) - V_{\pi_\theta}(\mathbf{x})] \\ &= Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) - V_{\pi_\theta}(\mathbf{x}) \\ &= A_{\pi_\theta}(\mathbf{x}, \mathbf{u})\end{aligned}$$

- $Q_{\pi_\theta} \rightarrow$ all samples \mathbf{x}, \mathbf{u} “push” the policy, bad samples push harder, policy moves to higher performance

Hence we have three possible critics:

$$\begin{aligned}\nabla_\theta J(\pi_\theta) &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) A_{\pi_\theta}(\mathbf{x}, \mathbf{u})] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) \delta(\mathbf{x}, \mathbf{u}, \mathbf{x}_+)] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) Q_{\pi_\theta}(\mathbf{x}, \mathbf{u})]\end{aligned}$$

Stochastic policy gradient - TD critic

Consider TD error:

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

We observe that:

$$\begin{aligned}\mathbb{E}[\delta | \mathbf{x}, \mathbf{u}] &= \mathbb{E}[L(\mathbf{x}, \mathbf{u}) + \gamma V_{\pi_\theta}(\mathbf{x}_+) - V_{\pi_\theta}(\mathbf{x})] \\ &= Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) - V_{\pi_\theta}(\mathbf{x}) \\ &= A_{\pi_\theta}(\mathbf{x}, \mathbf{u})\end{aligned}$$

- $Q_{\pi_\theta} \rightarrow$ all samples \mathbf{x}, \mathbf{u} “push” the policy, bad samples push harder, policy moves to higher performance
- $A_{\pi_\theta} \rightarrow$ good samples “pull” the policy towards them, bad samples “push” the policy away

Hence we have three possible critics:

$$\begin{aligned}\nabla_\theta J(\pi_\theta) &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) A_{\pi_\theta}(\mathbf{x}, \mathbf{u})] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) \delta(\mathbf{x}, \mathbf{u}, \mathbf{x}_+)] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) Q_{\pi_\theta}(\mathbf{x}, \mathbf{u})]\end{aligned}$$

Stochastic policy gradient - TD critic

Consider TD error:

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

We observe that:

$$\begin{aligned}\mathbb{E}[\delta | \mathbf{x}, \mathbf{u}] &= \mathbb{E}[L(\mathbf{x}, \mathbf{u}) + \gamma V_{\pi_\theta}(\mathbf{x}_+) - V_{\pi_\theta}(\mathbf{x})] \\ &= Q_{\pi_\theta}(\mathbf{x}, \mathbf{u}) - V_{\pi_\theta}(\mathbf{x}) \\ &= A_{\pi_\theta}(\mathbf{x}, \mathbf{u})\end{aligned}$$

Hence we have three possible critics:

$$\begin{aligned}\nabla_\theta J(\pi_\theta) &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) A_{\pi_\theta}(\mathbf{x}, \mathbf{u})] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) \delta(\mathbf{x}, \mathbf{u}, \mathbf{x}_+)] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(\mathbf{u} | \mathbf{x}) Q_{\pi_\theta}(\mathbf{x}, \mathbf{u})]\end{aligned}$$

- $Q_{\pi_\theta} \rightarrow$ all samples \mathbf{x}, \mathbf{u} “push” the policy, bad samples push harder, policy moves to higher performance
- $A_{\pi_\theta} \rightarrow$ good samples “pull” the policy towards them, bad samples “push” the policy away
- $\delta \rightarrow$ noisy version of A_{π_θ} , estimation of Q_{π_θ} not needed