

Policy Search Methods



Pieter Abbeel
UC Berkeley



Jan Peters
TU Darmstadt



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Motivation



2



Motivation

- Learning for high-dimensional robots is difficult:
 - Limit of Value Functions: fill-up state-space
 - Limit of Model Learning: accurate model!



2



Motivation

- Learning for high-dimensional robots is difficult:
 - Limit of Value Functions: fill-up state-space
 - Limit of Model Learning: accurate model!
- Starting with expert's knowledge helps!
 - Improving upon Demonstrations
 - Using Task-Appropriate Policies is possible



Motivation

- Learning for high-dimensional robots is difficult:
 - Limit of Value Functions: fill-up state-space
 - Limit of Model Learning: accurate model!
- Starting with expert's knowledge helps!
 - Improving upon Demonstrations
 - Using Task-Appropriate Policies is possible
- Exploring on the real system?



Motivation

- Learning for high-dimensional robots is difficult:
 - Limit of Value Functions: fill-up state-space
 - Limit of Model Learning: accurate model!
 - Starting with expert's knowledge helps!
 - Improving upon Demonstrations
 - Using Task-Appropriate Policies is possible
 - Exploring on the real system?
- ➔ Parametric Policy Search methods can do all that!



Bigger Picture



Reinforcement Learning Data

$$\mathcal{D} = \{(\mathbf{x}_i, \mathbf{u}_i, \mathbf{x}_{i+1}, r_i)\}$$

$$\mathbf{x}' \sim \mathcal{P}_{\mathbf{x}\mathbf{x}'}^{\mathbf{u}}$$
$$r \approx \mathcal{R}_{\mathbf{x}\mathbf{x}'}^{\mathbf{u}}$$

$$V^*(\mathbf{x})$$

$$\pi^*(\mathbf{u}|\mathbf{x})$$

Policy Search

$$V^*(\mathbf{x})$$

$$\pi^*(\mathbf{u}|\mathbf{x})$$

Value Function Methods

Now

$$\pi^*(\mathbf{u}|\mathbf{x})$$

Optimal Control with
Model Learning

Outline of the Lecture



4

Outline of the Lecture



1. Introduction with Policy Gradients

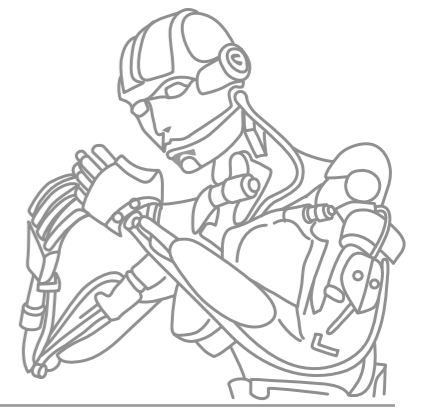


Outline of the Lecture



1. Introduction with Policy Gradients
2. Recent Advances in Policy Gradients





Outline of the Lecture

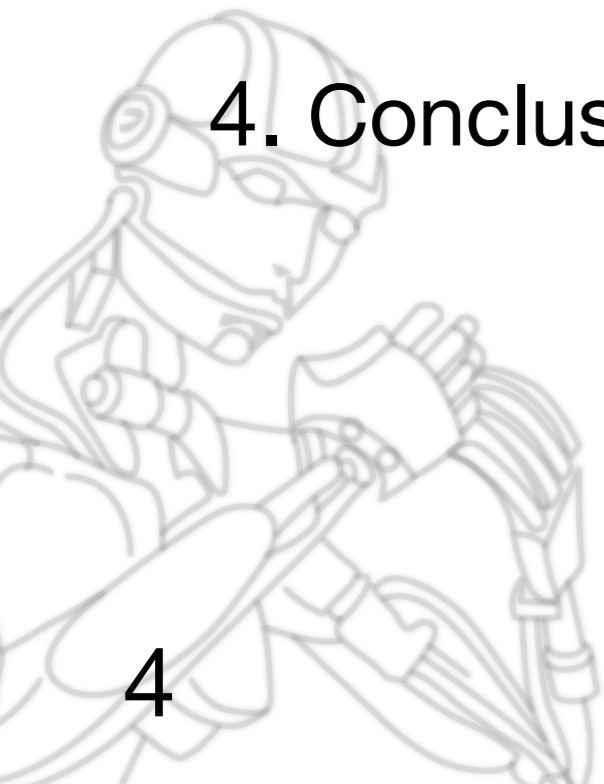
1. Introduction with Policy Gradients
2. Recent Advances in Policy Gradients
3. Probabilistic Policy Search with EM-like Approaches





Outline of the Lecture

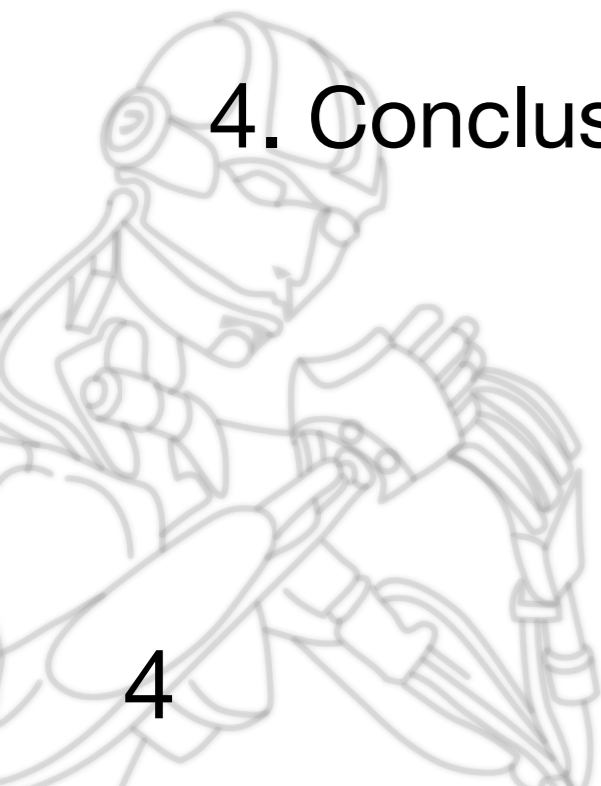
1. Introduction with Policy Gradients
2. Recent Advances in Policy Gradients
3. Probabilistic Policy Search with EM-like Approaches
4. Conclusion



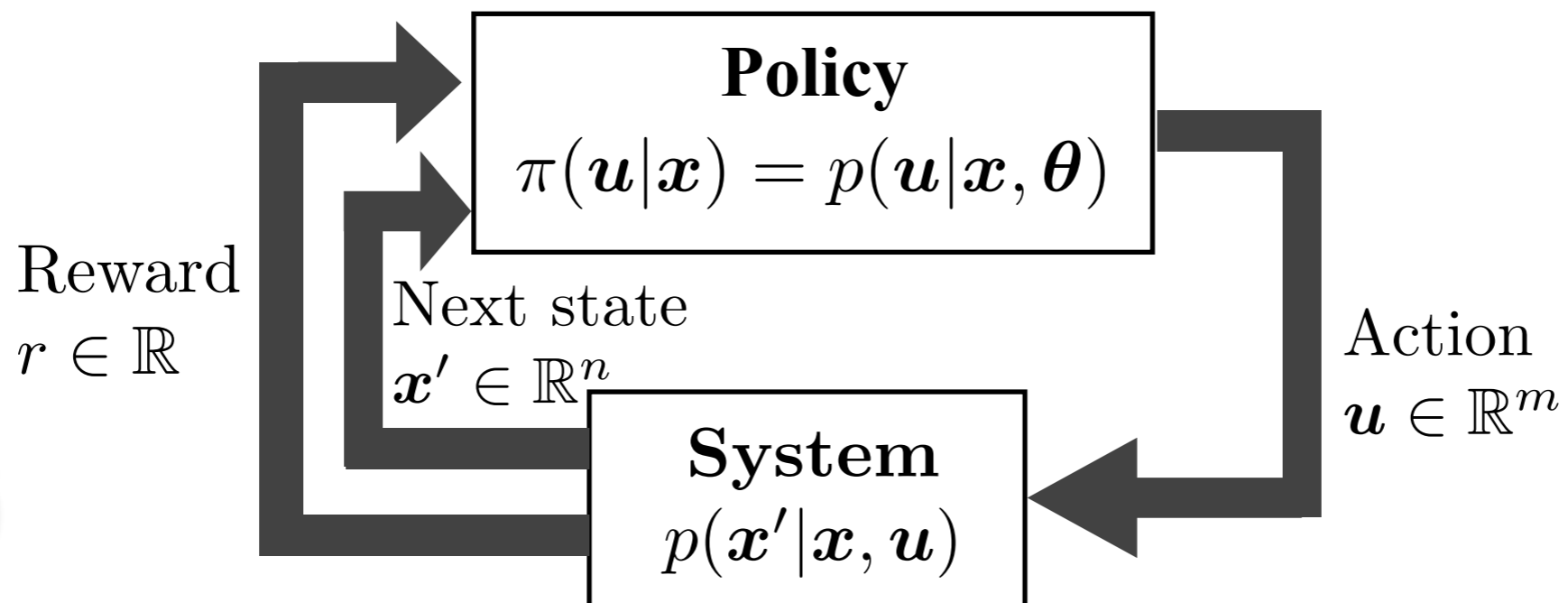
Outline of the Lecture



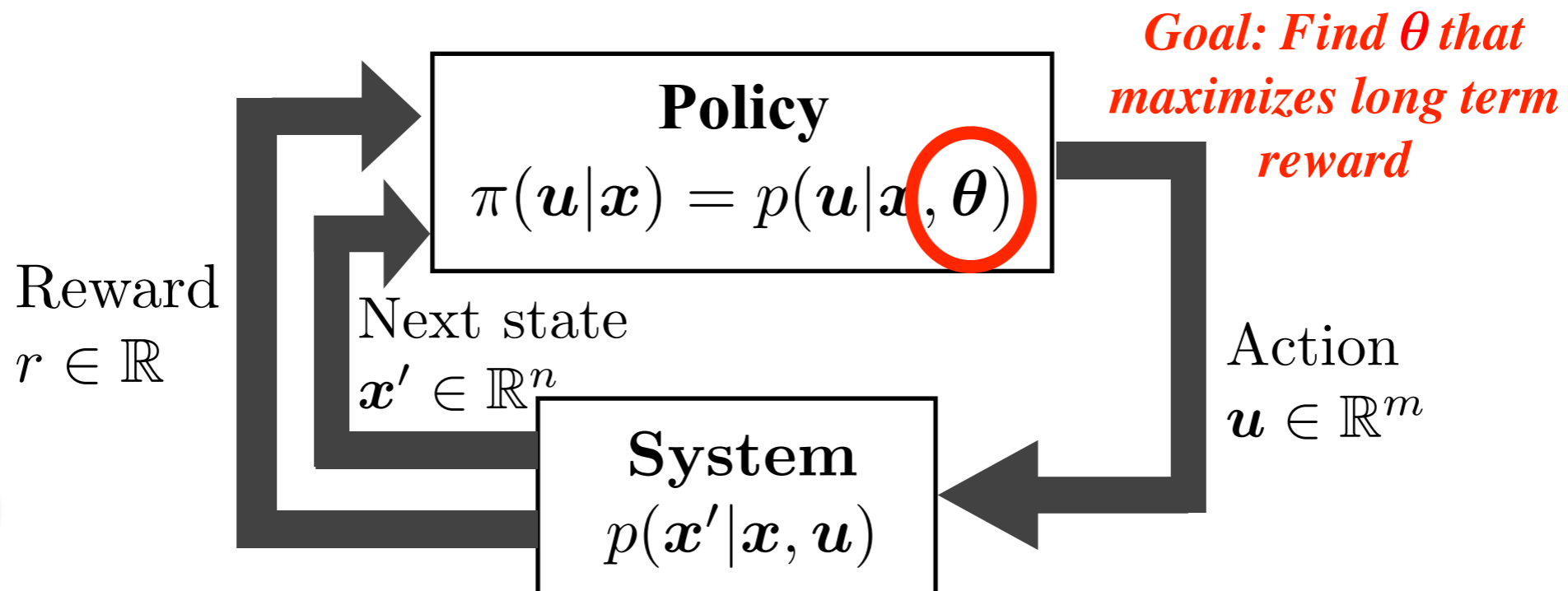
- ▶ 1. Introduction with Policy Gradients
- 2. Recent Advances in Policy Gradients
- 3. Probabilistic Policy Search with EM-like Approaches
- 4. Conclusion



Basics & Notation



Basics & Notation



Generic Reinforcement Learning Loop



- Learning requires an iteration through Policy Evaluation and Policy Improvement.





Generic Reinforcement Learning Loop

- Learning requires an iteration through Policy Evaluation and Policy Improvement.

Critic: Policy Evaluation

$$Q^\pi(\mathbf{x}, \mathbf{u}) = E \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \mid \mathbf{x}, \mathbf{u} \right\}$$

$$V^\pi(\mathbf{x}) = E \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \mid \mathbf{x} \right\}$$



Generic Reinforcement Learning Loop

- Learning requires an iteration through Policy Evaluation and Policy Improvement.

Critic: Policy Evaluation

$$Q^\pi(\mathbf{x}, \mathbf{u}) = E \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \mid \mathbf{x}, \mathbf{u} \right\}$$

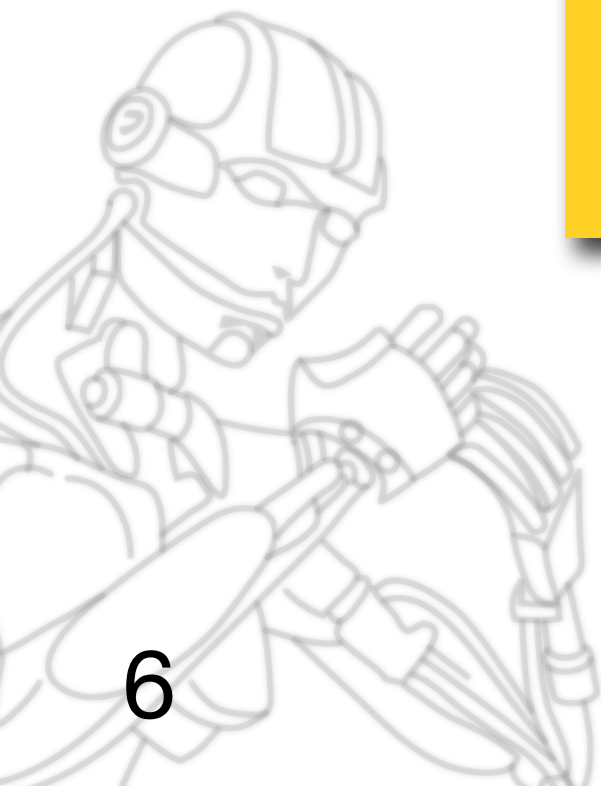
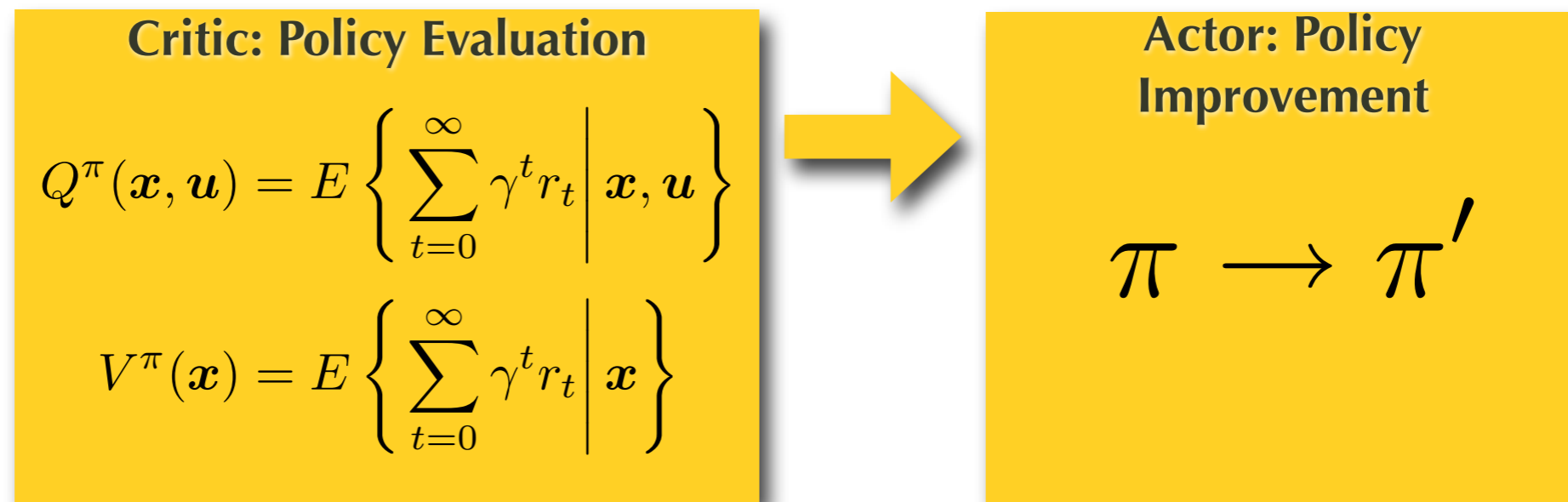
$$V^\pi(\mathbf{x}) = E \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \mid \mathbf{x} \right\}$$





Generic Reinforcement Learning Loop

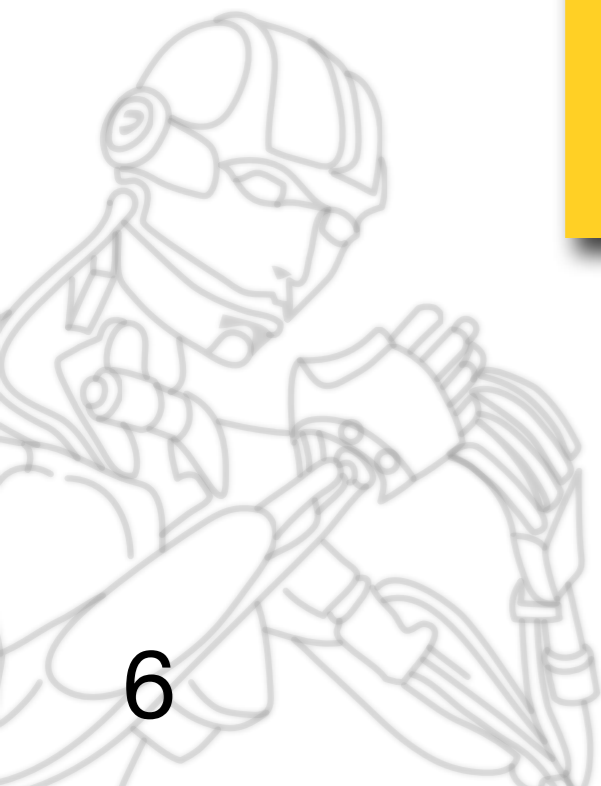
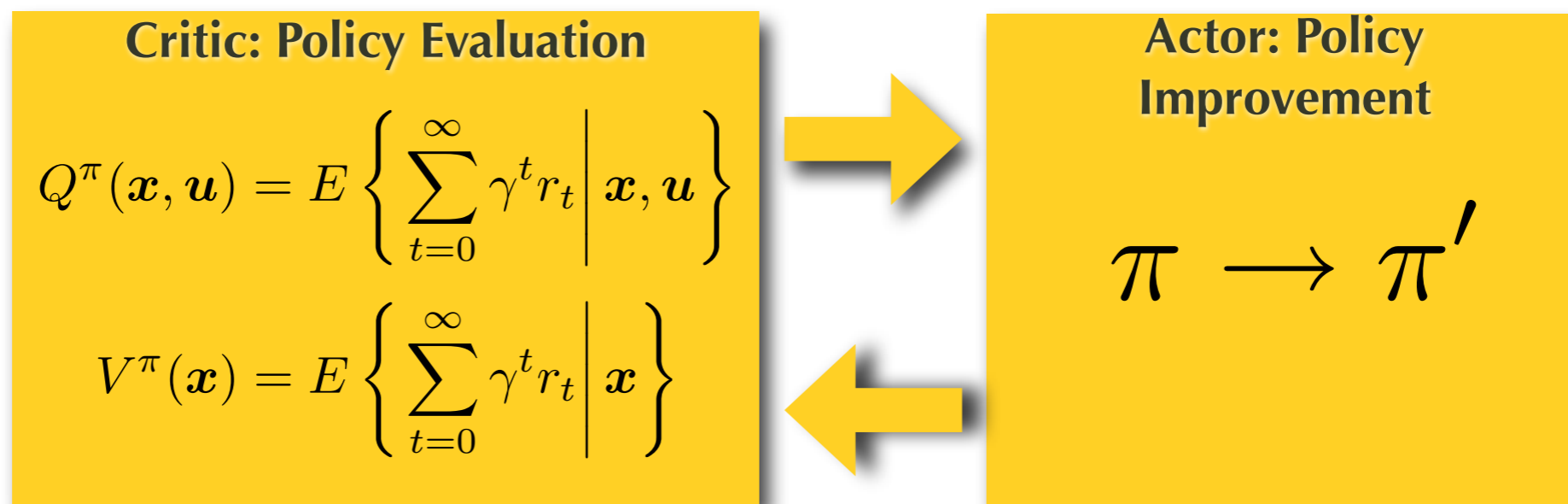
- Learning requires an iteration through Policy Evaluation and Policy Improvement.





Generic Reinforcement Learning Loop

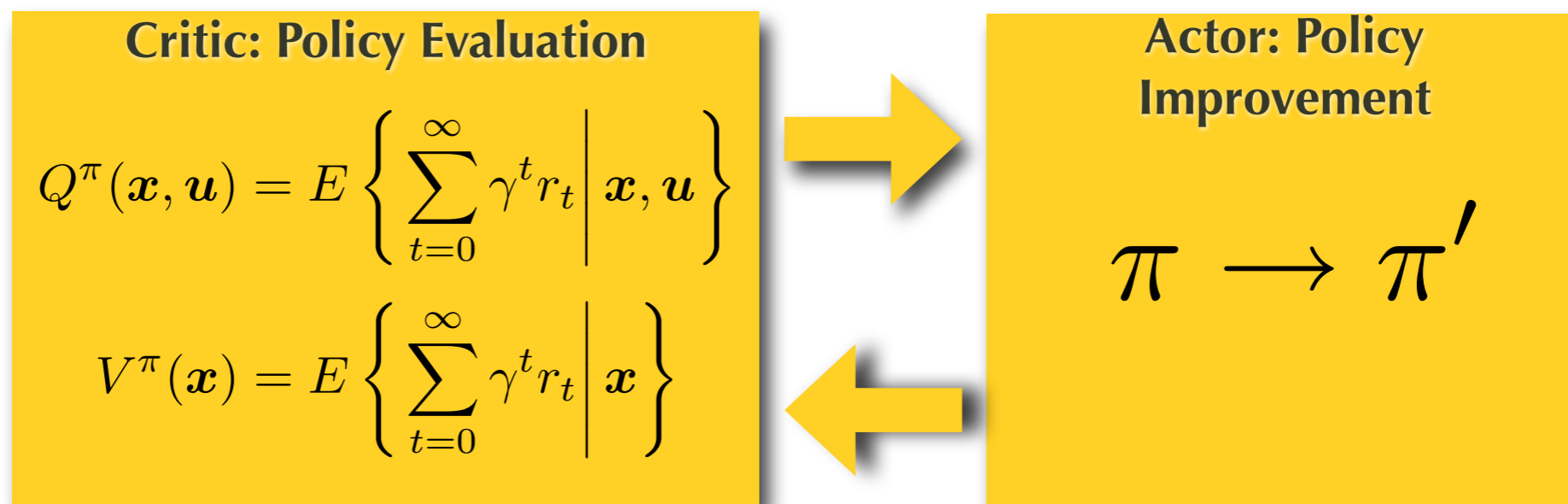
- Learning requires an iteration through Policy Evaluation and Policy Improvement.





Generic Reinforcement Learning Loop

- Learning requires an iteration through Policy Evaluation and Policy Improvement.



Requires Function Approximation



Greedy vs Incremental



7



Greedy vs Incremental



Greedy Updates:

$$\theta_{\pi'} = \operatorname{argmax}_{\tilde{\theta}} E_{\pi_{\tilde{\theta}}} \{Q^{\pi}(\mathbf{x}, \mathbf{u})\}$$



7





Greedy vs Incremental

Greedy Updates:

$$\theta_{\pi'} = \operatorname{argmax}_{\tilde{\theta}} E_{\pi_{\tilde{\theta}}} \{Q^{\pi}(\mathbf{x}, \mathbf{u})\}$$

Policy Gradient Updates:

$$\theta_{\pi'} = \theta_{\pi} + \alpha \left. \frac{dJ(\theta)}{d\theta} \right|_{\theta=\theta_{\pi}}$$



7

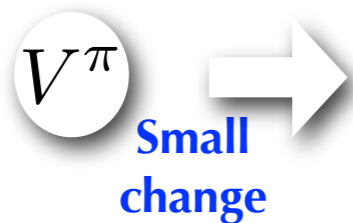


Greedy vs Incremental



Greedy Updates:

$$\theta_{\pi'} = \operatorname{argmax}_{\tilde{\theta}} E_{\pi_{\tilde{\theta}}} \{Q^{\pi}(\mathbf{x}, \mathbf{u})\}$$



Policy Gradient Updates:

$$\theta_{\pi'} = \theta_{\pi} + \alpha \left. \frac{dJ(\theta)}{d\theta} \right|_{\theta=\theta_{\pi}}$$

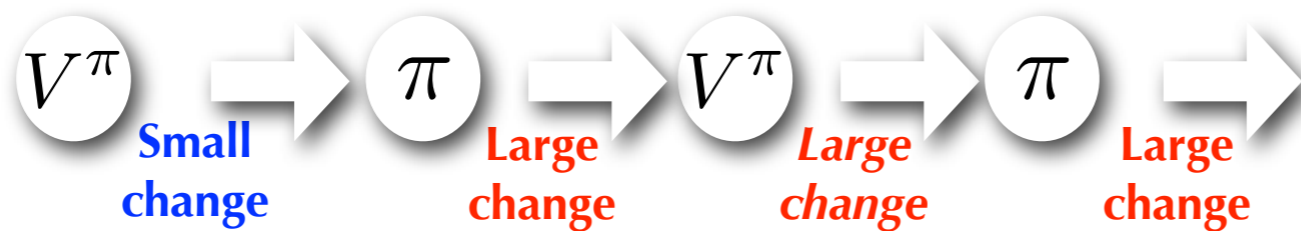


Greedy vs Incremental



Greedy Updates:

$$\theta_{\pi'} = \operatorname{argmax}_{\tilde{\theta}} E_{\pi_{\tilde{\theta}}} \{Q^{\pi}(x, u)\}$$



Policy Gradient Updates:

$$\theta_{\pi'} = \theta_{\pi} + \alpha \left. \frac{dJ(\theta)}{d\theta} \right|_{\theta=\theta_{\pi}}$$

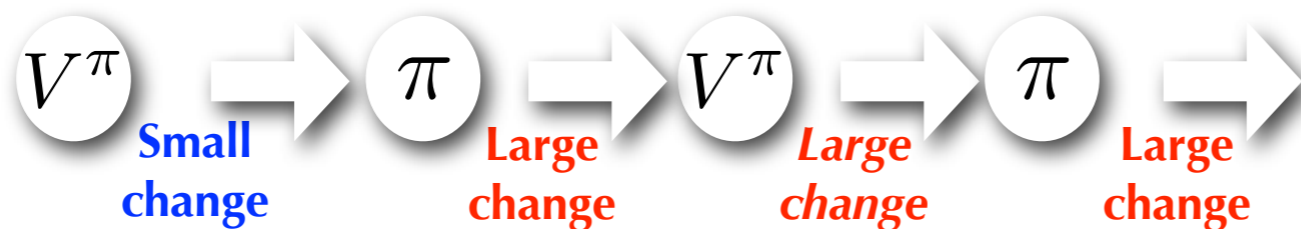


Greedy vs Incremental



Greedy Updates:

$$\theta_{\pi'} = \operatorname{argmax}_{\tilde{\theta}} E_{\pi_{\tilde{\theta}}} \{Q^{\pi}(x, u)\}$$



**potentially
unstable learning
process with large
policy jumps**

Policy Gradient Updates:

$$\theta_{\pi'} = \theta_{\pi} + \alpha \left. \frac{dJ(\theta)}{d\theta} \right|_{\theta=\theta_{\pi}}$$

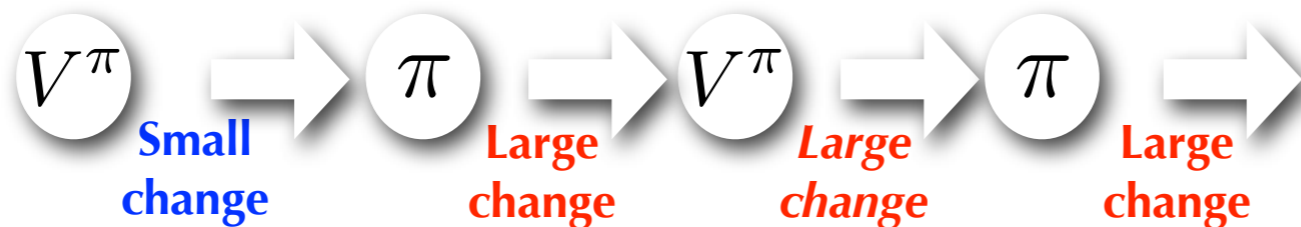


Greedy vs Incremental



Greedy Updates:

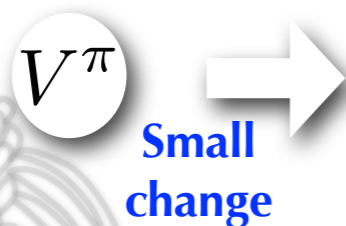
$$\theta_{\pi'} = \operatorname{argmax}_{\tilde{\theta}} E_{\pi_{\tilde{\theta}}} \{Q^{\pi}(x, u)\}$$



**potentially
unstable learning
process with large
policy jumps**

Policy Gradient Updates:

$$\theta_{\pi'} = \theta_{\pi} + \alpha \left. \frac{dJ(\theta)}{d\theta} \right|_{\theta=\theta_{\pi}}$$

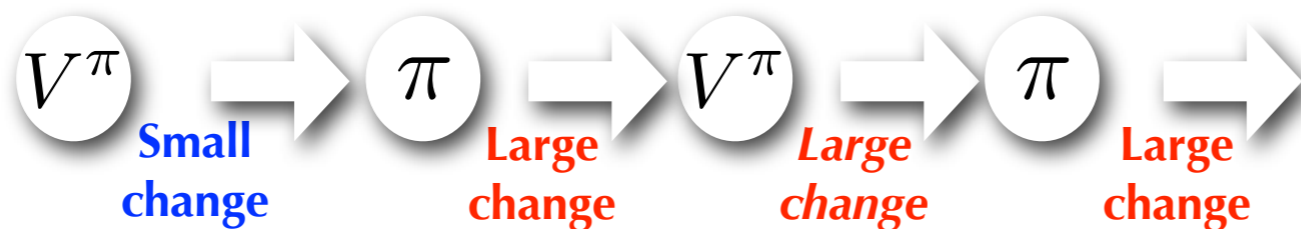


Greedy vs Incremental



Greedy Updates:

$$\theta_{\pi'} = \operatorname{argmax}_{\tilde{\theta}} E_{\pi_{\tilde{\theta}}} \{Q^{\pi}(x, u)\}$$



**potentially
unstable learning
process with large
policy jumps**

Policy Gradient Updates:

$$\theta_{\pi'} = \theta_{\pi} + \alpha \left. \frac{dJ(\theta)}{d\theta} \right|_{\theta=\theta_{\pi}}$$

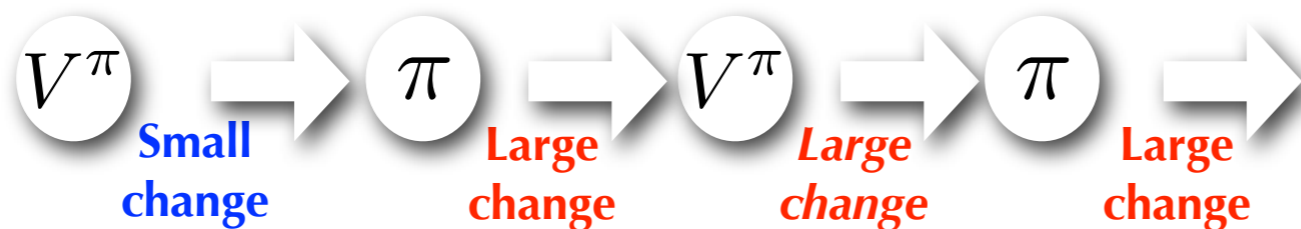


Greedy vs Incremental



Greedy Updates:

$$\theta_{\pi'} = \operatorname{argmax}_{\tilde{\theta}} E_{\pi_{\tilde{\theta}}} \{Q^{\pi}(x, u)\}$$



**potentially
unstable learning
process with large
policy jumps**

Policy Gradient Updates:

$$\theta_{\pi'} = \theta_{\pi} + \alpha \left. \frac{dJ(\theta)}{d\theta} \right|_{\theta=\theta_{\pi}}$$



**stable learning
process with
smooth policy
improvement**



Objective Function



- Goal: Optimize the expected return





Objective Function

- Goal: Optimize the expected return

$$J(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\boldsymbol{x}) \int_{\mathbb{U}} \pi(\boldsymbol{u}|\boldsymbol{x}) r(\boldsymbol{x}, \boldsymbol{u}) d\boldsymbol{u} d\boldsymbol{x},$$





Objective Function

- Goal: Optimize the expected return

$$J(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\boldsymbol{x}) \int_{\mathbb{U}} \pi(\boldsymbol{u}|\boldsymbol{x}) r(\boldsymbol{x}, \boldsymbol{u}) d\boldsymbol{u} d\boldsymbol{x},$$

State distribution







Objective Function

- Goal: Optimize the expected return

$$J(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) r(\mathbf{x}, \mathbf{u}) d\mathbf{u} d\mathbf{x},$$

State distribution  Policy (we can choose it) 





Objective Function

- Goal: Optimize the expected return

$$J(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) r(\mathbf{x}, \mathbf{u}) d\mathbf{u} d\mathbf{x},$$

State distribution Policy (we can choose it) Reward

Three red arrows point from the labels below to the corresponding terms in the equation: one from 'State distribution' to $d^{\pi}(\mathbf{x})$, one from 'Policy (we can choose it)' to $\pi(\mathbf{u}|\mathbf{x})$, and one from 'Reward' to $r(\mathbf{x}, \mathbf{u})$.





Objective Function

- Goal: Optimize the expected return

$$J(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) r(\mathbf{x}, \mathbf{u}) d\mathbf{u} d\mathbf{x},$$

State distribution Policy (we can choose it) Reward

$$= E \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \right\}$$



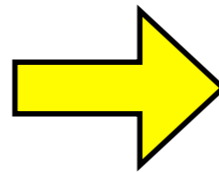
Gradient-based Policy Iteration



Actor: Policy Evaluation

*Estimate
Gradient*

$$\mathbf{g}_t = \nabla J(\boldsymbol{\theta})$$



Gradient-based Policy Iteration



Actor: Policy Evaluation

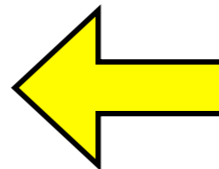
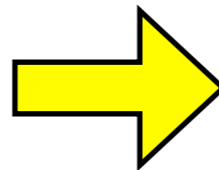
*Estimate
Gradient*

$$\mathbf{g}_t = \nabla J(\boldsymbol{\theta})$$

Critic: Policy Improvement

*Update
Parameters*

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha_t \mathbf{g}_t$$





Policy Gradient Methods

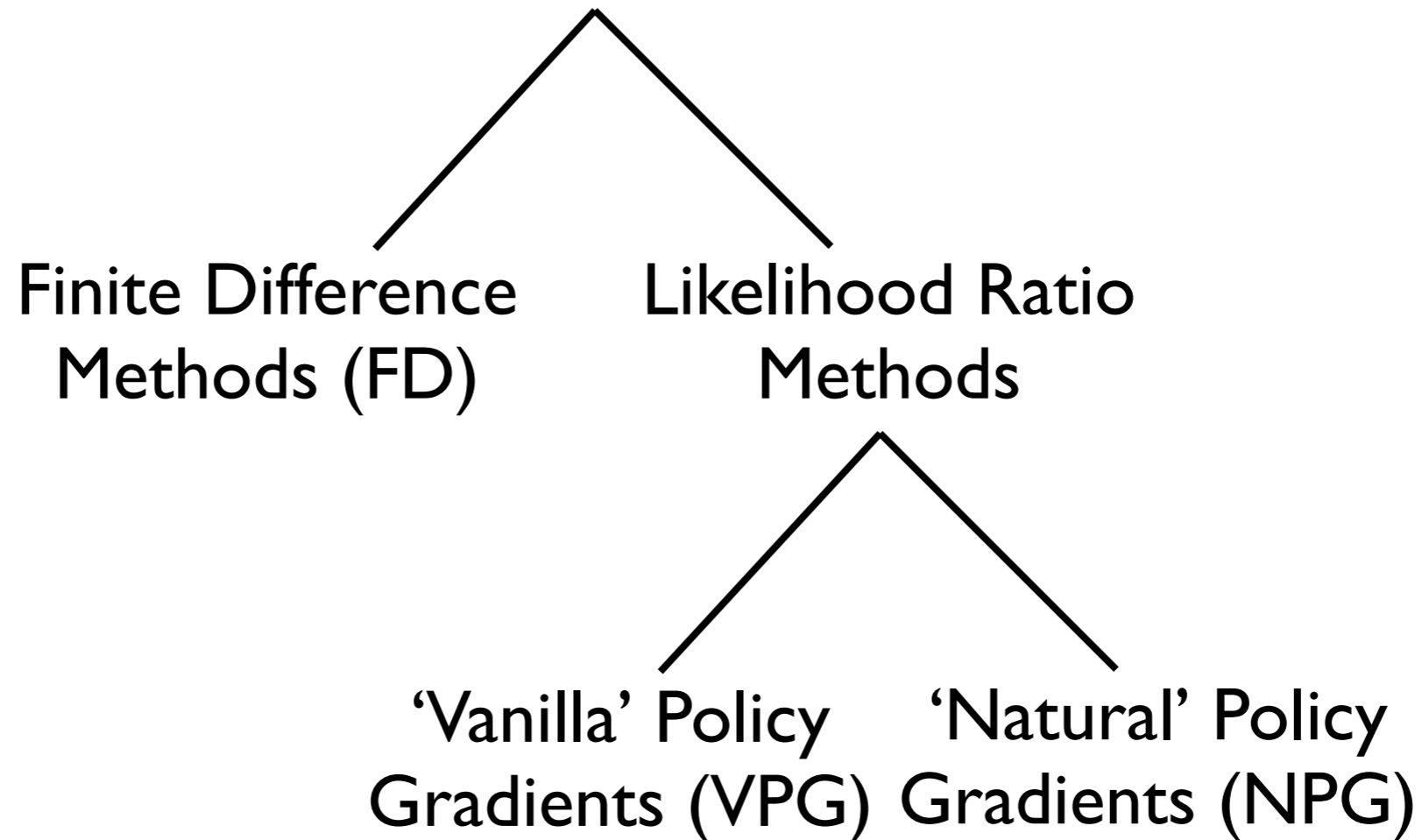


Many related approaches exist in the literature, e.g., Mean-Value Differentiation, Model-based approaches, DDP, Frequency-based approaches, etc.





Policy Gradient Methods



Many related approaches exist in the literature, e.g., Mean-Value Differentiation, Model-based approaches, DDP, Frequency-based approaches, etc.



Black-Box Approaches



11

A large class of algorithms includes Kiefer-Wolfowitz procedure, Robbins-Monroe, Simultaneous Perturbation Stochastic Approximation SPSA, ...



Black-Box Approaches



I. Perturb the parameters of your policy:

$$\theta + \delta\theta \rightarrow$$



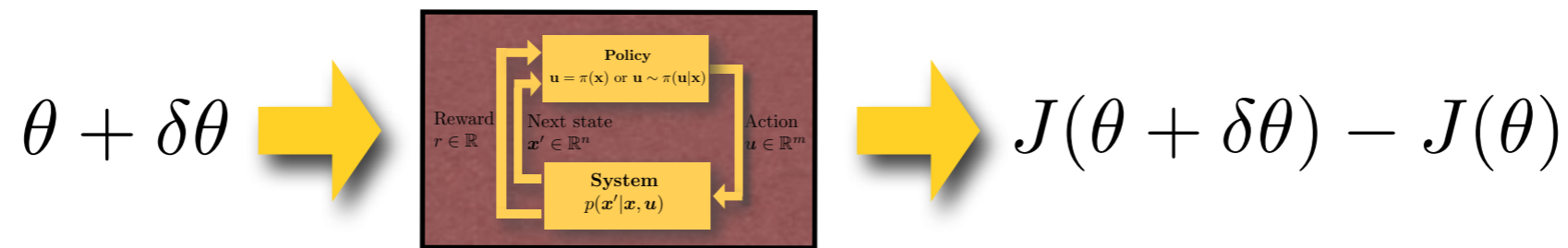
A large class of algorithms includes Kiefer-Wolfowitz procedure, Robbins-Monroe, Simultaneous Perturbation Stochastic Approximation SPSA, ...



Black-Box Approaches



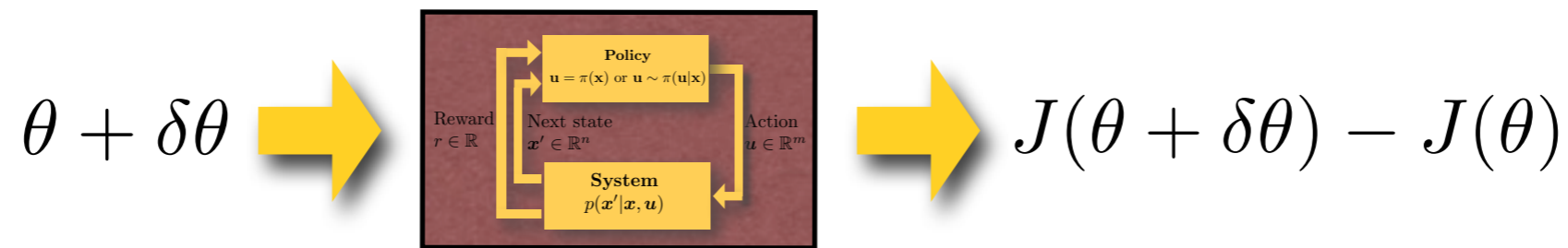
I. Perturb the parameters of your policy:



Black-Box Approaches



1. Perturb the parameters of your policy:



2. Gradient estimation by regression:

$$\mathbf{g}_{\text{FD}} = (\Delta\Theta^T \Delta\Theta)^{-1} \Delta\Theta^T \Delta J.$$

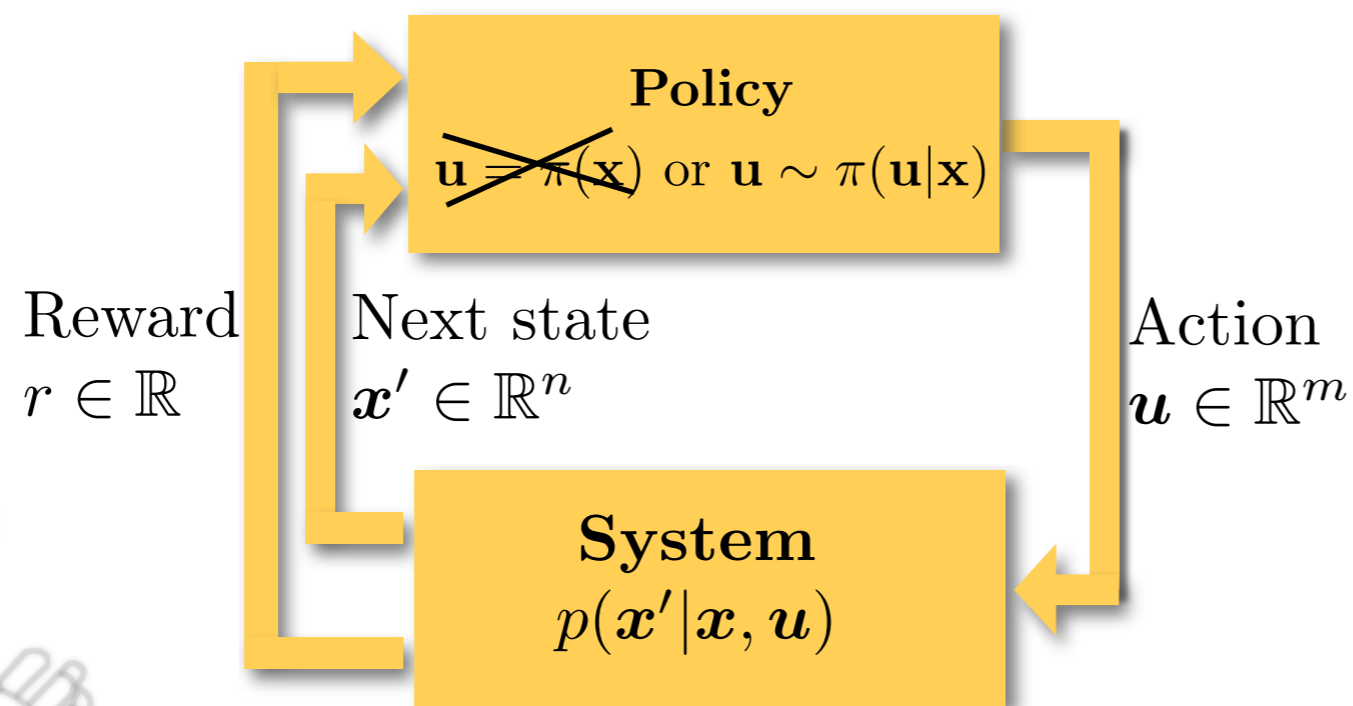
A large class of algorithms includes Kiefer-Wolfowitz procedure, Robbins-Monroe, Simultaneous Perturbation Stochastic Approximation SPSA, ...





Whitebox Approaches

Whitebox Approach: Use an explorative, stochastic policy and make use of the knowledge of your policy.



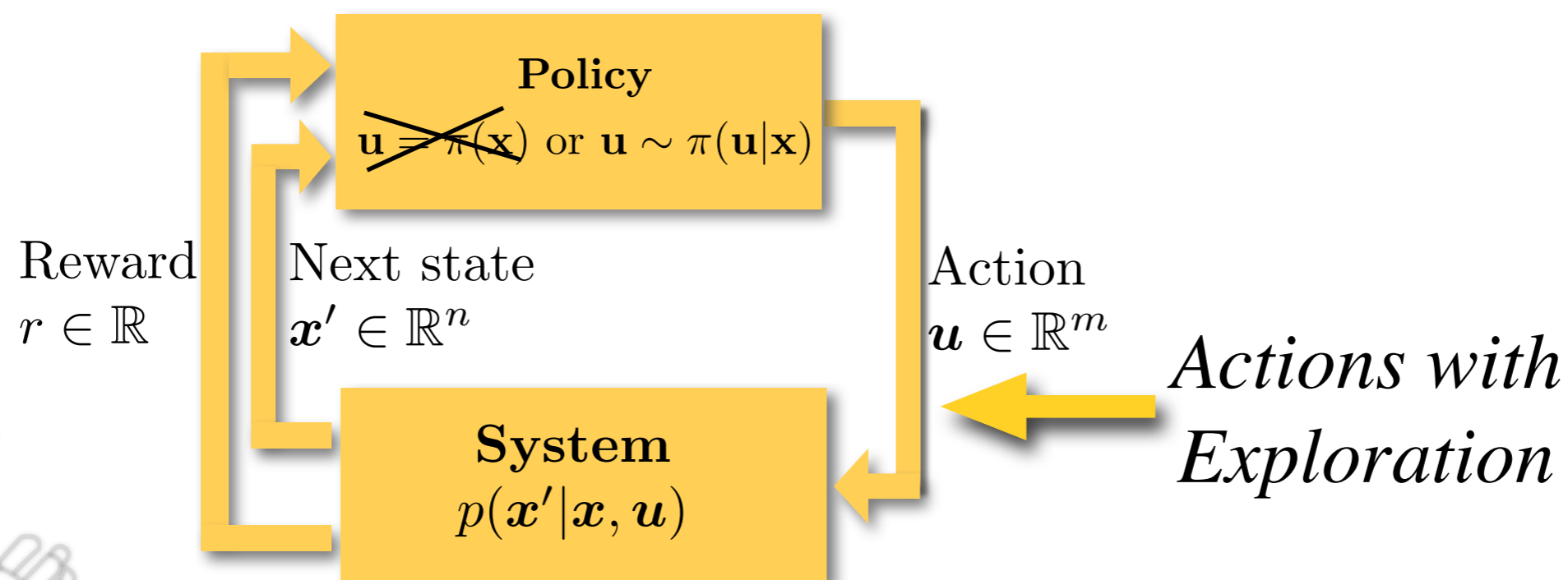
Many related approaches in the RL literature starting from Werbos (1971), Hasdorff (1976), Williams (1988), ...



Whitebox Approaches



Whitebox Approach: Use an explorative, stochastic policy and make use of the knowledge of your policy.



Many related approaches in the RL literature starting from Werbos (1971), Hasdorff (1976), Williams (1988), ...



Likelihood Ratio Gradient



For a cost function

$$J(\theta) = \int_{\mathbb{T}} p_{\theta}(\tau|\pi) R(\tau) d\tau$$

we have the gradient

$$\nabla J(\theta) = \nabla \int_{\mathbb{T}} p_{\theta}(\tau|\pi) R(\tau) d\tau = \int_{\mathbb{T}} \nabla p_{\theta}(\tau|\pi) R(\tau) d\tau$$

Using the trick

$$\nabla p_{\theta}(\tau|\pi) = p_{\theta}(\tau|\pi) \nabla \log p_{\theta}(\tau|\pi)$$

we obtain

$$\begin{aligned} \nabla J(\theta) &= \int_{\mathbb{T}} p_{\theta}(\tau|\pi) \nabla \log p_{\theta}(\tau|\pi) R(\tau) d\tau \\ &= E\{\nabla \log p_{\theta}(\tau|\pi) R(\tau)\} \\ &\approx \frac{1}{K} \sum_{k=1}^K \nabla \log p_{\theta}(\tau_k|\pi) R(\tau_k) \end{aligned}$$

Needs
only
samples!





Likelihood Ratio Gradient

Why is this cool?

Because: The definition of a path probability

$$p(\boldsymbol{\tau}) = p(\mathbf{x}_1) \prod_{t=1}^T p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) \pi(\mathbf{u}_t | \mathbf{x}_t)$$

implies

$$\log p(\boldsymbol{\tau}) = \sum_{t=1}^T \log \pi(\mathbf{u}_t | \mathbf{x}_t) + \text{const}$$

Hence, we can get the derivative of the distribution without a model of the system:

$$\nabla \log p(\boldsymbol{\tau}) = \sum_{t=1}^T \nabla \log \pi(\mathbf{u}_t | \mathbf{x}_t)$$



Likelihood Ratio Gradient



As a result:

$$\begin{aligned}\nabla J(\theta) &= E \left\{ \sum_{t=1}^T \nabla \log \pi(\mathbf{u}_t | \mathbf{x}_t) R(\boldsymbol{\tau}) \right\} \\ &= E \left\{ \sum_{t=1}^T \nabla \log \pi(\mathbf{u}_t | \mathbf{x}_t) \sum_{h=t}^T r(\mathbf{x}_h, \mathbf{u}_h) \right\} \\ &= E \left\{ \sum_{t=1}^T \nabla \log \pi(\mathbf{u}_t | \mathbf{x}_t) Q^\pi(\mathbf{x}_t, \mathbf{u}_t) \right\}\end{aligned}$$





Outline of the Lecture

1. Introduction with Policy Gradients

▶ 2. Recent Advances in Policy Gradients

3. Probabilistic Policy Search with EM-like Approaches

4. Conclusion

Likelihood Ratio Approach: Policy Gradient Theorem



17

Originally discovered: Aleksandrov, 1968; Glynn, 1986
Examples: episodic REINFORCE, SRV, GPOMDP

Likelihood Ratio Approach: Policy Gradient Theorem



According to the policy gradient theorem, the gradient can be computed as

$$\nabla_{\theta} J(\theta) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \nabla_{\theta} \pi(\mathbf{u}|\mathbf{x}) (Q^{\pi}(\mathbf{x}, \mathbf{u}) - b^{\pi}(\mathbf{x})) d\mathbf{u} d\mathbf{x}.$$



Likelihood Ratio Approach: Policy Gradient Theorem



According to the policy gradient theorem, the gradient can be computed as

$$\nabla_{\theta} J(\theta) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \nabla_{\theta} \pi(\mathbf{u}|\mathbf{x}) (Q^{\pi}(\mathbf{x}, \mathbf{u}) - b^{\pi}(\mathbf{x})) d\mathbf{u} d\mathbf{x}.$$

Gradient of the
expected return

Derivative
only of the
policy

State-action
value function

Arbitrary baseline
function

Likelihood Ratio Approach: Policy Gradient Theorem



According to the policy gradient theorem, the gradient can be computed as

$$\nabla_{\theta} J(\theta) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \nabla_{\theta} \pi(\mathbf{u}|\mathbf{x}) (Q^{\pi}(\mathbf{x}, \mathbf{u}) - b^{\pi}(\mathbf{x})) d\mathbf{u} d\mathbf{x}.$$

Gradient of the
expected return

Derivative
only of the
policy

State-action
value function

Arbitrary baseline
function

Problems: High Variance, dependence on the baseline, slow convergence!

Originally discovered: Aleksandrov, 1968; Glynn, 1986
Examples: episodic REINFORCE, SRV, GPOMDP

Compatible Function Approximation





Compatible Function Approximation

The state-action value function can be replaced by

$$Q^\pi(\mathbf{x}, \mathbf{u}) \equiv f_{\mathbf{w}}^\pi(\mathbf{x}, \mathbf{u}) = \frac{d \log \pi(\mathbf{u} | \mathbf{x})^T}{d\boldsymbol{\theta}} \mathbf{w}$$

without biasing the gradient.

Compatible Function Approximation



The state-action value function can be replaced by

$$Q^\pi(\mathbf{x}, \mathbf{u}) \equiv f_{\mathbf{w}}^\pi(\mathbf{x}, \mathbf{u}) = \frac{d \log \pi(\mathbf{u} | \mathbf{x})}{d \boldsymbol{\theta}}^T \mathbf{w}$$

State-action
value function

Compatible function
approximation

Log-policy
derivative

Parameters
of the
function
approximator

without biasing the gradient.

Compatible Function Approximation



The state-action value function can be replaced by

$$Q^\pi(\mathbf{x}, \mathbf{u}) \equiv f_{\mathbf{w}}^\pi(\mathbf{x}, \mathbf{u}) = \frac{d \log \pi(\mathbf{u} | \mathbf{x})}{d \boldsymbol{\theta}}^T \mathbf{w}$$

State-action
value function

Compatible function
approximation

Log-policy
derivative

Parameters
of the
function
approximator

without biasing the gradient.

Thus, the gradient becomes

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^\pi(\mathbf{x}) \int_{\mathbb{U}} \nabla_{\boldsymbol{\theta}} \pi(\mathbf{u} | \mathbf{x}) (f_{\mathbf{w}}^\pi(\mathbf{x}, \mathbf{u}) - b^\pi(\mathbf{x})) d\mathbf{u} d\mathbf{x}.$$

All-Action Gradient



19

(Peters et al. 2003, 2005)



All-Action Gradient

By integrating over all possible actions in a state, the baseline can be integrated out, and the gradient becomes

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \int_{\mathbb{X}} d^{\pi}(\boldsymbol{x}) \int_{\mathbb{U}} \nabla_{\boldsymbol{\theta}} \pi(\boldsymbol{u}|\boldsymbol{x}) (f_w^{\pi}(\boldsymbol{x}, \boldsymbol{u}) - b(\boldsymbol{x})) d\boldsymbol{u} d\boldsymbol{x}, \\ &= \int_{\mathbb{X}} d^{\pi}(\boldsymbol{x}) \int_{\mathbb{U}} \pi(\boldsymbol{u}|\boldsymbol{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\boldsymbol{u}|\boldsymbol{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\boldsymbol{u}|\boldsymbol{x})^T \boldsymbol{w} d\boldsymbol{u} d\boldsymbol{x}, \\ &= \boldsymbol{F}(\boldsymbol{\theta}) \boldsymbol{w}.\end{aligned}$$





All-Action Gradient

By integrating over all possible actions in a state, the baseline can be integrated out, and the gradient becomes

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \nabla_{\theta} \pi(\mathbf{u}|\mathbf{x}) (f_w^{\pi}(\mathbf{x}, \mathbf{u}) - b(\mathbf{x})) d\mathbf{u} d\mathbf{x}, \\ &= \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) \nabla_{\theta} \log \pi(\mathbf{u}|\mathbf{x}) \nabla_{\theta} \log \pi(\mathbf{u}|\mathbf{x})^T \mathbf{w} d\mathbf{u} d\mathbf{x}, \\ &= \mathbf{F}(\theta) \mathbf{w}.\end{aligned}$$





All-Action Gradient

By integrating over all possible actions in a state, the baseline can be integrated out, and the gradient becomes

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \nabla_{\theta} \pi(\mathbf{u}|\mathbf{x}) (f_w^{\pi}(\mathbf{x}, \mathbf{u}) - b(\mathbf{x})) d\mathbf{u} d\mathbf{x}, \\ &= \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) \nabla_{\theta} \log \pi(\mathbf{u}|\mathbf{x}) \nabla_{\theta} \log \pi(\mathbf{u}|\mathbf{x})^T \mathbf{w} d\mathbf{u} d\mathbf{x}, \\ &= \mathbf{F}(\theta) \mathbf{w}.\end{aligned}$$

All Action Matrix

Parameters

Natural Gradients





Natural Gradients

A more efficient gradient in learning problems is the natural gradient (Amari, 1998)

$$\tilde{\nabla}_{\theta} J(\theta) = G^{-1}(\theta) \nabla_{\theta} J(\theta)$$





Natural Gradients

A more efficient gradient in learning problems is the natural gradient (Amari, 1998)

Natural gradient



$$\tilde{\nabla}_{\theta} J(\theta) = G^{-1}(\theta) \nabla_{\theta} J(\theta)$$



Inverse of the Fisher Information Matrix



'Vanilla' gradient





Natural Gradients

A more efficient gradient in learning problems is the natural gradient (Amari, 1998)

Natural gradient



$$\tilde{\nabla}_{\theta} J(\theta) = G^{-1}(\theta) \nabla_{\theta} J(\theta)$$



Inverse of the Fisher Information Matrix



'Vanilla' gradient

where the policy gradient $\nabla J(\theta)$ is given by the policy gradient theorem.



Natural Gradients

A more efficient gradient in learning problems is the natural gradient (Amari, 1998)

Natural gradient



$$\tilde{\nabla}_{\theta} J(\theta) = G^{-1}(\theta) \nabla_{\theta} J(\theta)$$



Inverse of the Fisher Information Matrix



'Vanilla' gradient

where the policy gradient $\nabla J(\theta)$ is given by the policy gradient theorem.

But how can we obtain the Fisher information matrix $G(\theta)$??

Fisher Information



So how does the All-Action Matrix



21

(Peters et al., 2003; 2005; Bagnell et al., 2003)

Fisher Information



So how does the All-Action Matrix

$$\mathbf{F}(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}|\mathbf{x}) d\mathbf{u}d\mathbf{x}.$$



Fisher Information



So how does the All-Action Matrix

$$\mathbf{F}(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}|\mathbf{x}) d\mathbf{u} d\mathbf{x}.$$

relate to the Fisher Information Matrix



Fisher Information



So how does the All-Action Matrix

$$\mathbf{F}(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}|\mathbf{x}) d\mathbf{u} d\mathbf{x}.$$

relate to the Fisher Information Matrix

$$\mathbf{G}(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log (d^{\pi}(\mathbf{x})\pi(\mathbf{u}|\mathbf{x})) \nabla_{\boldsymbol{\theta}} \log (d^{\pi}(\mathbf{x})\pi(\mathbf{u}|\mathbf{x})) d\mathbf{u} d\mathbf{x}.$$



Fisher Information

So how does the All-Action Matrix

$$\mathbf{F}(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}|\mathbf{x}) d\mathbf{u}d\mathbf{x}.$$

relate to the Fisher Information Matrix

$$\mathbf{G}(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log (d^{\pi}(\mathbf{x})\pi(\mathbf{u}|\mathbf{x})) \nabla_{\boldsymbol{\theta}} \log (d^{\pi}(\mathbf{x})\pi(\mathbf{u}|\mathbf{x})) d\mathbf{u}d\mathbf{x}.$$

While Kakade (2002) suggested that \mathbf{F} is an ‘average of point Fisher information matrices’, we could prove that

Fisher Information



So how does the All-Action Matrix

$$\mathbf{F}(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}|\mathbf{x}) d\mathbf{u} d\mathbf{x}.$$

relate to the Fisher Information Matrix

$$\mathbf{G}(\boldsymbol{\theta}) = \int_{\mathbb{X}} d^{\pi}(\mathbf{x}) \int_{\mathbb{U}} \pi(\mathbf{u}|\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log (d^{\pi}(\mathbf{x})\pi(\mathbf{u}|\mathbf{x})) \nabla_{\boldsymbol{\theta}} \log (d^{\pi}(\mathbf{x})\pi(\mathbf{u}|\mathbf{x})) d\mathbf{u} d\mathbf{x}.$$

While Kakade (2002) suggested that \mathbf{F} is an ‘average of point Fisher information matrices’, we could prove that

$$\mathbf{F} = \mathbf{G}.$$

Natural Policy Gradients





Natural Policy Gradients

Thus, the gradient simplifies to

$$\tilde{\nabla}_{\theta} J(\theta) = G^{-1}(\theta) \nabla_{\theta} J(\theta) = G^{-1}(\theta) F(\theta) w = w,$$



Natural Policy Gradients



Thus, the gradient simplifies to

$$\tilde{\nabla}_{\theta} J(\theta) = G^{-1}(\theta) \nabla_{\theta} J(\theta) = G^{-1}(\theta) F(\theta) w = w,$$

and the policy parameter update becomes

$$\theta_{t+1} = \theta_t + \alpha_t w_t.$$

Important: The gradient estimation simplifies to determining the parameters of the compatible function approximation.

Are they useful?



Are they useful?



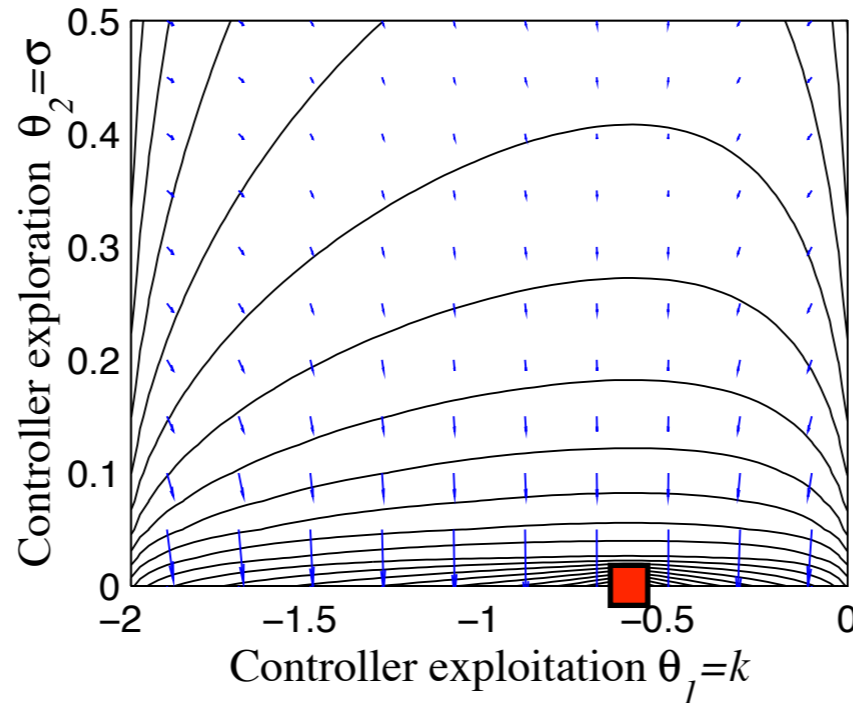
Linear Quadratic Regulation

$$x_{t+1} = Ax_t + Bu_t$$

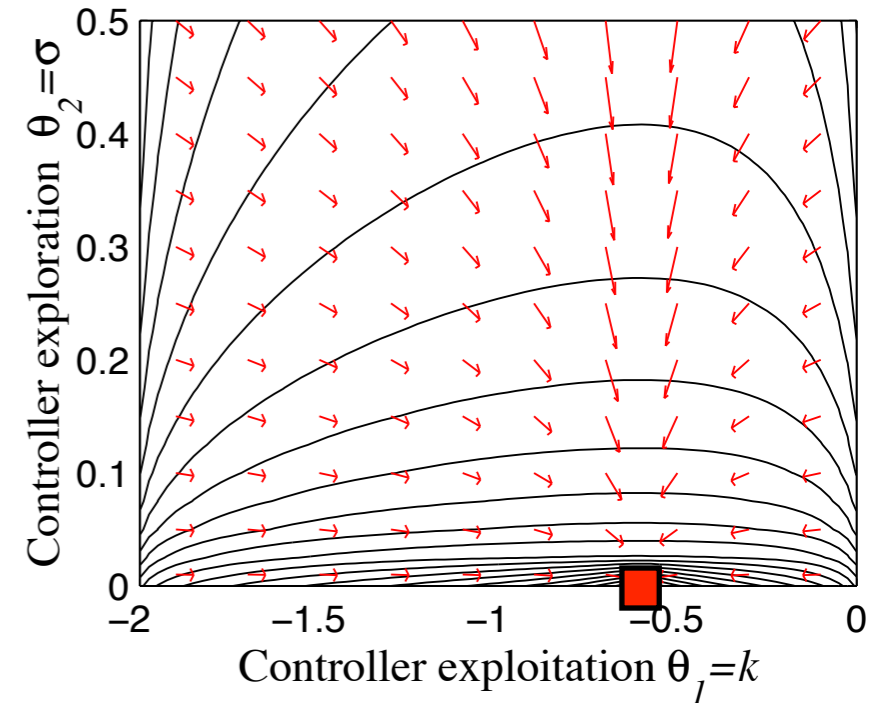
$$u_t \sim \pi(u|x_t) = \mathcal{N}(u|kx_t, \sigma)$$

$$r_t = -x_t^T Q x_t - u_t^T R u_t$$

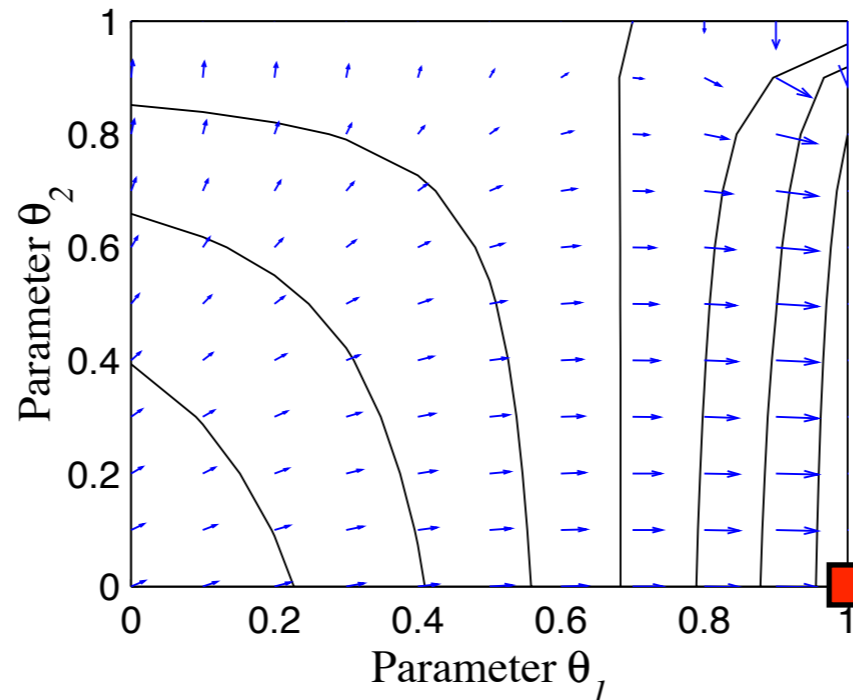
(a) LQR policy gradient



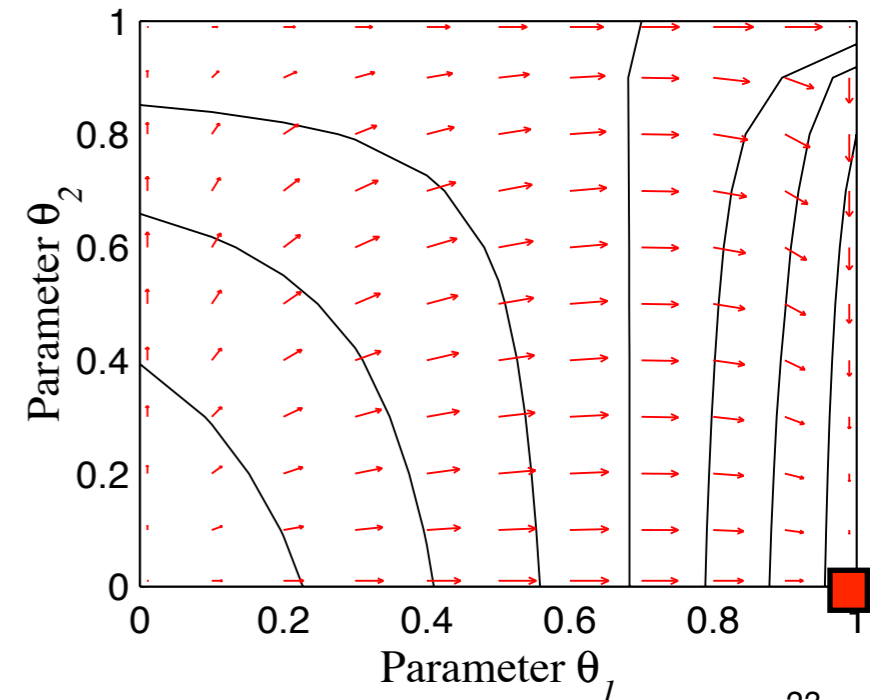
(b) LQR natural gradient



(c) Two state policy gradient

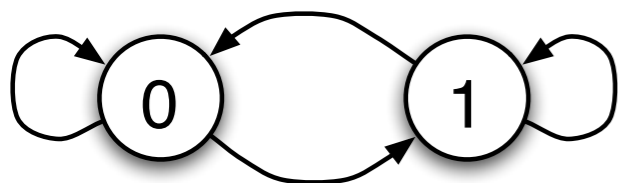


(d) Two state natural gradient



Two-State Problem

$$u = 0, r = 0$$



$$u = 0$$

$$u = 1$$

$$u = 1$$

$$r = 1$$

$$r = 0$$

$$r = 2$$

Can the Compatible FA be learned?





Can the Compatible FA be learned?

The compatible function approximation is mean-zero! Thus, it can only represent the Advantage Function:

$$f_w^\pi(\mathbf{x}, \mathbf{u}) = Q^\pi(\mathbf{x}, \mathbf{u}) - V^\pi(\mathbf{x}) = A^\pi(\mathbf{x}, \mathbf{u}).$$





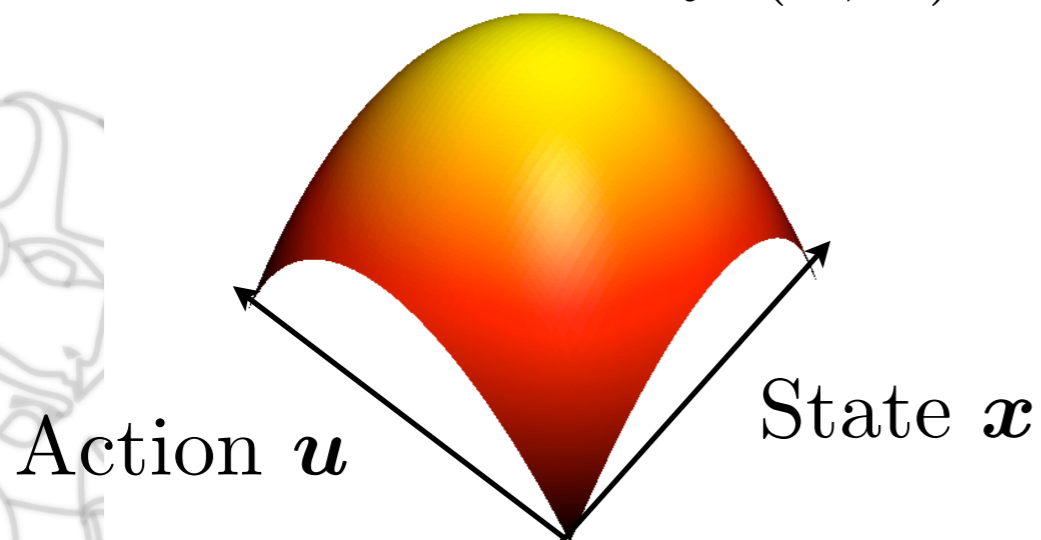
Can the Compatible FA be learned?

The compatible function approximation is mean-zero! Thus, it can only represent the Advantage Function:

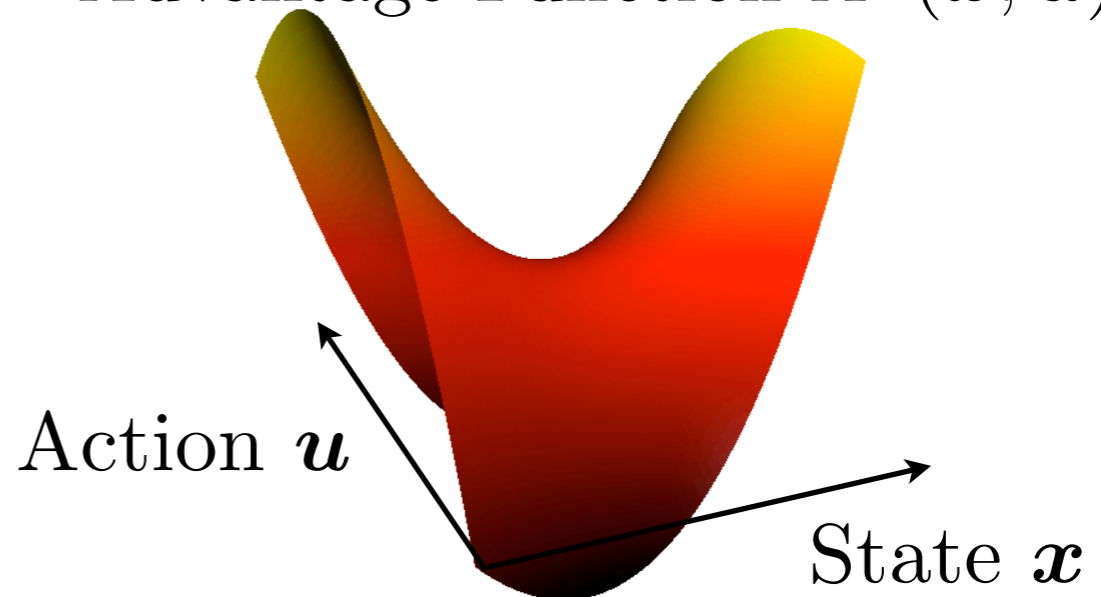
$$f_w^\pi(\mathbf{x}, \mathbf{u}) = Q^\pi(\mathbf{x}, \mathbf{u}) - V^\pi(\mathbf{x}) = A^\pi(\mathbf{x}, \mathbf{u}).$$

The advantage function is very different from the value functions

Value Function $Q^\pi(\mathbf{x}, \mathbf{u})$



Advantage Function $A^\pi(\mathbf{x}, \mathbf{u})$





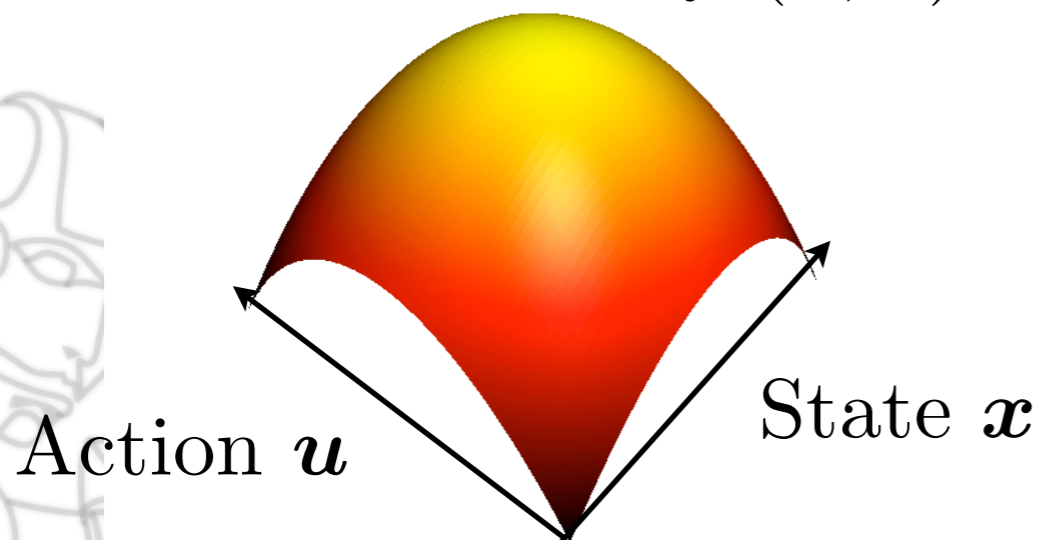
Can the Compatible FA be learned?

The compatible function approximation is mean-zero! Thus, it can only represent the Advantage Function:

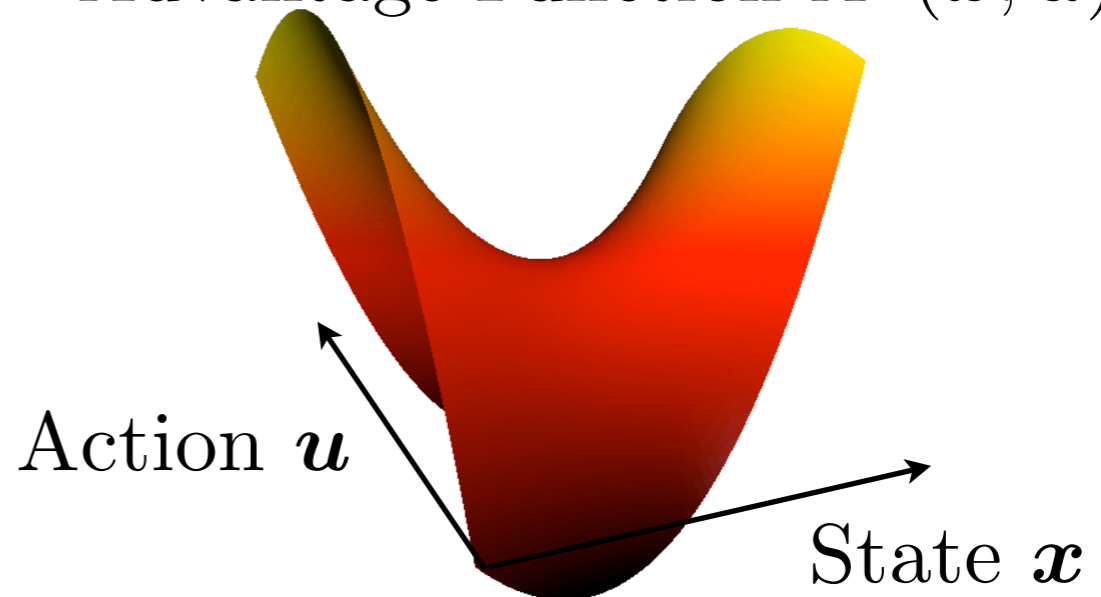
$$f_{\mathbf{w}}^{\pi}(\mathbf{x}, \mathbf{u}) = Q^{\pi}(\mathbf{x}, \mathbf{u}) - V^{\pi}(\mathbf{x}) = A^{\pi}(\mathbf{x}, \mathbf{u}).$$

The advantage function is very different from the value functions

Value Function $Q^{\pi}(\mathbf{x}, \mathbf{u})$



Advantage Function $A^{\pi}(\mathbf{x}, \mathbf{u})$



Traditional value function learning methods such as Temporal Difference learning cannot be applied.

The Compatible FA can be learned!





The Compatible FA can be learned!

We cannot apply traditional methods directly on





The Compatible FA can be learned!

We cannot apply traditional methods directly on

$$f_w^\pi(\mathbf{x}, \mathbf{u}) = Q^\pi(\mathbf{x}, \mathbf{u}) - V^\pi(\mathbf{x}) = A^\pi(\mathbf{x}, \mathbf{u}).$$

But when we add further function approximation





The Compatible FA can be learned!

We cannot apply traditional methods directly on

$$f_w^\pi(\mathbf{x}, \mathbf{u}) = Q^\pi(\mathbf{x}, \mathbf{u}) - V^\pi(\mathbf{x}) = A^\pi(\mathbf{x}, \mathbf{u}).$$

But when we add further function approximation

$$V^\pi(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{v}$$

into the Bellman equation





The Compatible FA can be learned!

We cannot apply traditional methods directly on

$$f_w^\pi(\mathbf{x}, \mathbf{u}) = Q^\pi(\mathbf{x}, \mathbf{u}) - V^\pi(\mathbf{x}) = A^\pi(\mathbf{x}, \mathbf{u}).$$

But when we add further function approximation

$$V^\pi(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{v}$$

into the Bellman equation

$$V^\pi(\mathbf{x}_t) + \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}_t | \mathbf{x}_t)^T \mathbf{w} = r(\mathbf{x}_t, \mathbf{u}_t) + \gamma V^\pi(\mathbf{x}_{t+1}) + \epsilon_t$$

we get a **linear regression problem** which can be solved with appropriate regression techniques, e.g., Boyan's (1996) LSTD(λ) algorithm.



The Compatible FA can be learned!

We cannot apply traditional methods directly on

$$f_w^\pi(\mathbf{x}, \mathbf{u}) = Q^\pi(\mathbf{x}, \mathbf{u}) - V^\pi(\mathbf{x}) = A^\pi(\mathbf{x}, \mathbf{u}).$$

But when we add further function approximation

$$V^\pi(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{v}$$

into the Bellman equation

$$V^\pi(\mathbf{x}_t) + \nabla_{\theta} \log \pi(\mathbf{u}_t | \mathbf{x}_t)^T \mathbf{w} = r(\mathbf{x}_t, \mathbf{u}_t) + \gamma V^\pi(\mathbf{x}_{t+1}) + \epsilon_t$$

we get a **linear regression problem** which can be solved with appropriate regression techniques, e.g., Boyan's (1996) LSTD(λ) algorithm.

➔ **Allows the derivation of many well-known old reinforcement learning algorithms, e.g., Sutton et al. (1983) Actor-Critic and Bradtke & Barto's (1993)**

25 LQR-Q-Learning.

What about this additional FA?





What about this additional FA?

...but in many cases, we don't have a good additional function approximations!





What about this additional FA?

...but in many cases, we don't have a good additional function approximations!

For one rollout, if we sum up the Bellman Equations





What about this additional FA?

...but in many cases, we don't have a good additional function approximations!

For one rollout, if we sum up the Bellman Equations

and eliminate the values of the intermediary states, we obtain



What about this additional FA?

...but in many cases, we don't have a good additional function approximations!

For one rollout, if we sum up the Bellman Equations

$$\begin{aligned} V^\pi(\mathbf{x}_0) + \nabla \log \pi(\mathbf{u}_0|\mathbf{x}_0) &= r(\mathbf{x}_0, \mathbf{u}_0) + \gamma V^\pi(\mathbf{x}_1) \\ V^\pi(\mathbf{x}_1) + \nabla \log \pi(\mathbf{u}_1|\mathbf{x}_1) &= r(\mathbf{x}_1, \mathbf{u}_1) + \gamma V^\pi(\mathbf{x}_0) \\ &\vdots \\ V^\pi(\mathbf{x}_T) + \nabla \log \pi(\mathbf{u}_T|\mathbf{x}_T) &= r(\mathbf{x}_T, \mathbf{u}_T) + \gamma V^\pi(\mathbf{x}_{T+1}) \end{aligned}$$

and eliminate the values of the intermediary states, we obtain

$$\underbrace{V^\pi(\mathbf{x}_0)}_J + \underbrace{\left(\sum_{t=0}^T \nabla_{\theta} \log \pi(\mathbf{u}_t|\mathbf{x}_t) \right)^T}_{\varphi_i} \mathbf{w} = \underbrace{\sum_{t=0}^T \gamma^t r(\mathbf{x}_t, \mathbf{u}_t)}_{R_i} + \underbrace{\gamma^{T+1} V^\pi(\mathbf{x}_{T+1})}_0$$

26

ONE offset parameter suffices as additional function approximation!

(Peters et al. 2003, 2005)

Episodic Natural Actor-Critic



Episodic Natural Actor-Critic



Critic: Episodic Evaluation

$$\Phi = \begin{bmatrix} \varphi_1 & \varphi_2 & \cdots & \varphi_N \\ 1 & 1 & \cdots & 1 \end{bmatrix}^T$$

$$\mathbf{R} = [R_1, R_2^T, \dots, R_N^T]^T$$

$$\begin{bmatrix} w \\ J \end{bmatrix} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{R}$$

Episodic Natural Actor-Critic



Critic: Episodic Evaluation

Sufficient
Statistics

$$\Phi = \begin{bmatrix} \varphi_1 & \varphi_2 & \cdots & \varphi_N \\ 1 & 1 & \cdots & 1 \end{bmatrix}^T$$

$$\mathbf{R} = [R_1, R_2^T, \dots, R_N^T]^T$$

$$\begin{bmatrix} w \\ J \end{bmatrix} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{R}$$

Episodic Natural Actor-Critic



Critic: Episodic Evaluation

Sufficient
Statistics

$$\Phi = \begin{bmatrix} \varphi_1 & \varphi_2 & \cdots & \varphi_N \\ 1 & 1 & \cdots & 1 \end{bmatrix}^T$$

$$\mathbf{R} = [R_1, R_2^T, \dots, R_N^T]^T$$

Linear
Regression

$$\begin{bmatrix} w \\ J \end{bmatrix} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{R}$$

Episodic Natural Actor-Critic



Critic: Episodic Evaluation

Sufficient
Statistics

$$\Phi = \begin{bmatrix} \varphi_1 & \varphi_2 & \cdots & \varphi_N \\ 1 & 1 & \cdots & 1 \end{bmatrix}^T$$

$$\mathbf{R} = [R_1, R_2^T, \dots, R_N^T]^T$$

Linear
Regression

$$\begin{bmatrix} w \\ J \end{bmatrix} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{R}$$

Episodic Natural Actor-Critic



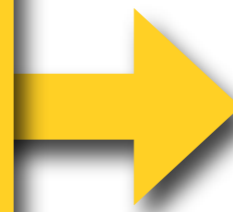
Sufficient
Statistics

$$\Phi = \begin{bmatrix} \varphi_1 & \varphi_2 & \cdots & \varphi_N \\ 1 & 1 & \cdots & 1 \end{bmatrix}^T$$

$$\mathbf{R} = [R_1, R_2^T, \dots, R_N^T]^T$$

Linear
Regression

$$\begin{bmatrix} \mathbf{w} \\ J \end{bmatrix} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{R}$$



**Actor: Natural
Policy Gradient
Improvement**

$$\theta_{t+1} = \theta_t + \alpha_t \mathbf{w}_t.$$

Episodic Natural Actor-Critic



Sufficient
Statistics

$$\Phi = \begin{bmatrix} \varphi_1 & \varphi_2 & \cdots & \varphi_N \\ 1 & 1 & \cdots & 1 \end{bmatrix}^T$$

$$\mathbf{R} = [R_1, R_2^T, \dots, R_N^T]^T$$

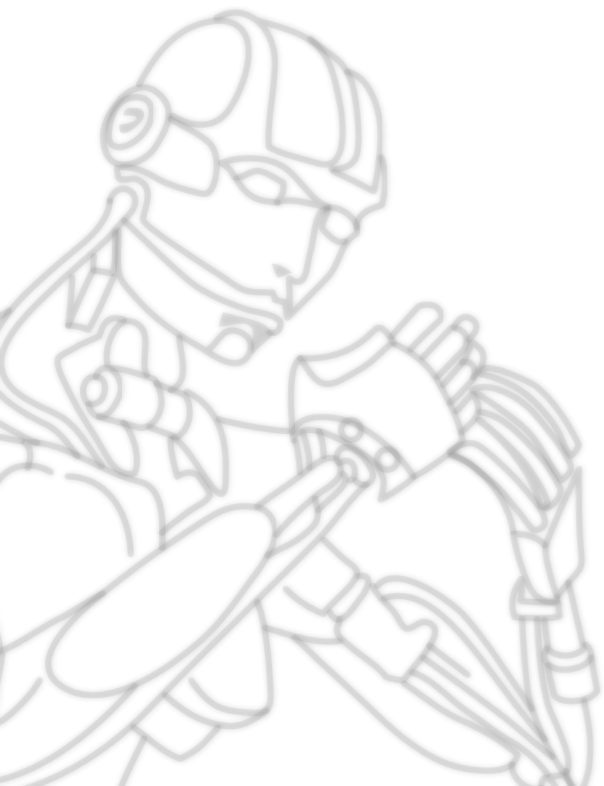
Linear
Regression

$$\begin{bmatrix} \mathbf{w} \\ J \end{bmatrix} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{R}$$

Actor: Natural
Policy Gradient
Improvement

$$\theta_{t+1} = \theta_t + \alpha_t \mathbf{w}_t.$$

Important Points



Important Points



Points worth highlighting:



Important Points



Points worth highlighting:

- ★ Natural policy gradients are *independent* of the chosen policy parameterization!



Important Points



Points worth highlighting:

- ★ Natural policy gradients are *independent* of the chosen policy parameterization!
- ★ They correspond to *steepest descent in policy space* and not in the parameter space.



Important Points



Points worth highlighting:

- ★ Natural policy gradients are *independent* of the chosen policy parameterization!
- ★ They correspond to *steepest descent in policy space* and not in the parameter space.
- ★ *Convergence to a local minimum* is guaranteed!



Important Points



Points worth highlighting:

- ★ Natural policy gradients are *independent* of the chosen policy parameterization!
 - ★ They correspond to *steepest descent in policy space* and not in the parameter space.
 - ★ *Convergence to a local minimum* is guaranteed!
- ? ...but we still need to estimate the natural gradient!

Important Points



Points worth highlighting:

- ★ Natural policy gradients are *independent* of the chosen policy parameterization!
 - ★ They correspond to *steepest descent in policy space* and not in the parameter space.
 - ★ *Convergence to a local minimum* is guaranteed!
- ? ...but we still need to estimate the natural gradient!

Benchmarking on Cart-Pole Regulation

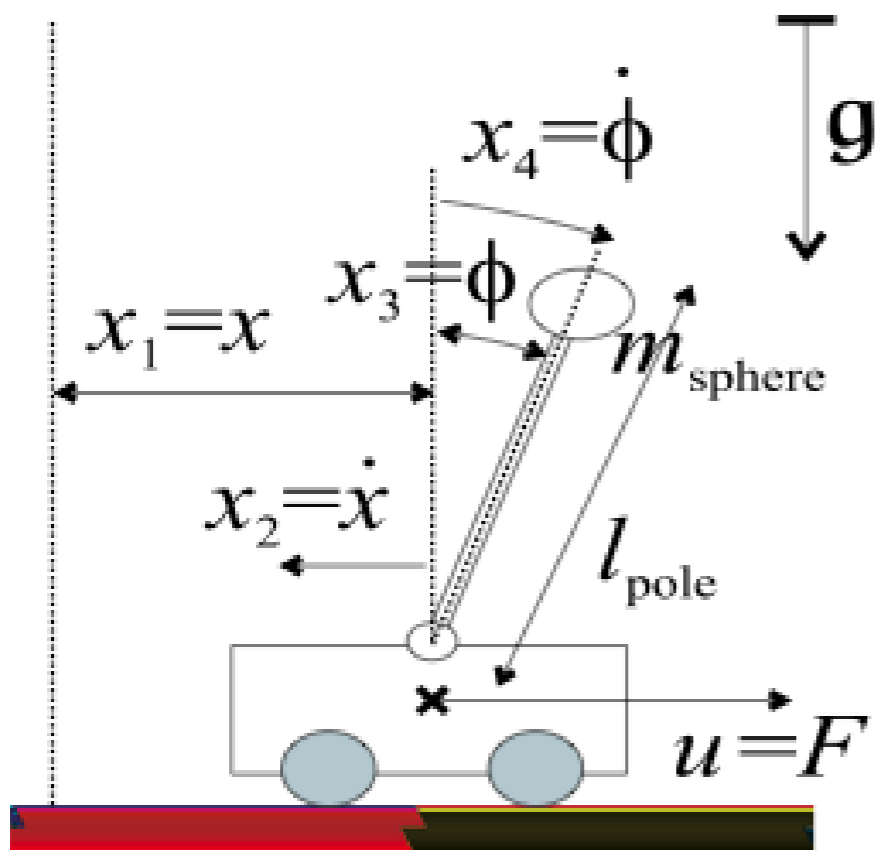


- standard benchmark

- maximize time inside the target area:

$$r(\mathbf{x}, \mathbf{u}) = \begin{cases} 0 & \text{if } |\phi| < 0.05\text{rad}, |x| < 0.05\text{m} \\ -1 & \text{otherwise.} \end{cases}$$

- episodic restarts





Finite Difference Gradients

Algorithm	Fair performance (>-120) after	Good performance (>-80) after	best performance
Finite Difference Gradients with Standard Descent	12,300	Not reached	-84
Finite Difference Gradients with RPROP Rule	7,450	45,650	-76



Vanilla Policy Gradients

Algorithm	Fair performance (> -120) after	Good performance (> -80) after	best performance
Vanilla PG without Baseline	22,200	Not reached	-102
Vanilla PG with Optimal Baseline	1,200	26,450	-76
Vanilla PG with Optimal Baseline and RPROP	450	3,000	-64



Vanilla Policy Gradients

Algorithm	Fair performance (>-120) after	Good performance (>-80) after	best performance
Vanilla PG without Baseline	22,200	Not reached	-102
Vanilla PG with Optimal Baseline	1,200	26,450	-76
Vanilla PG with Optimal Baseline and RPROP	450	3,000	-64

Fastest Initial Improvement



Vanilla Policy Gradients

Algorithm	Fair performance (>-120) after	Good performance (>-80) after	best performance
Vanilla PG without Baseline	22,200	Not reached	-102
Vanilla PG with Optimal Baseline	1,200	26,450	-76
Vanilla PG with Optimal Baseline and RPROP	450	3,000	Not optimal -64

Fastest Initial Improvement

Episodic Natural Actor-Critic



Algorithm	Fair performance (>-120) after	Good performance (>-80) after	best performance
Episodic Natural Actor-Critic	750	5,050	-55
Episodic Natural Actor-Critic with RPROP	Not reached	Not reached	-130

Episodic Natural Actor-Critic



Algorithm	Fair performance (>-120) after	Good performance (>-80) after	best performance Best Final Performance
Episodic Natural Actor-Critic	750	5,050	-55
Episodic Natural Actor-Critic with RPROP	Not reached	Not reached	-130

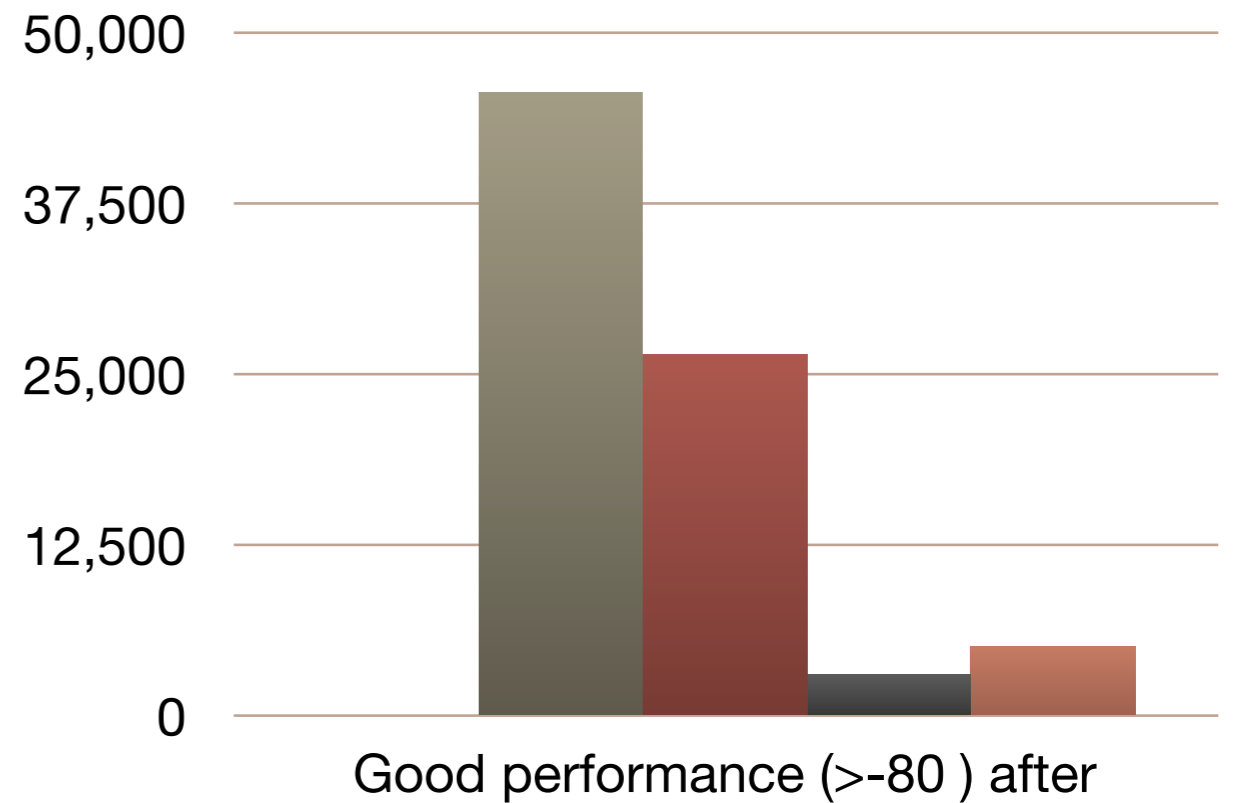
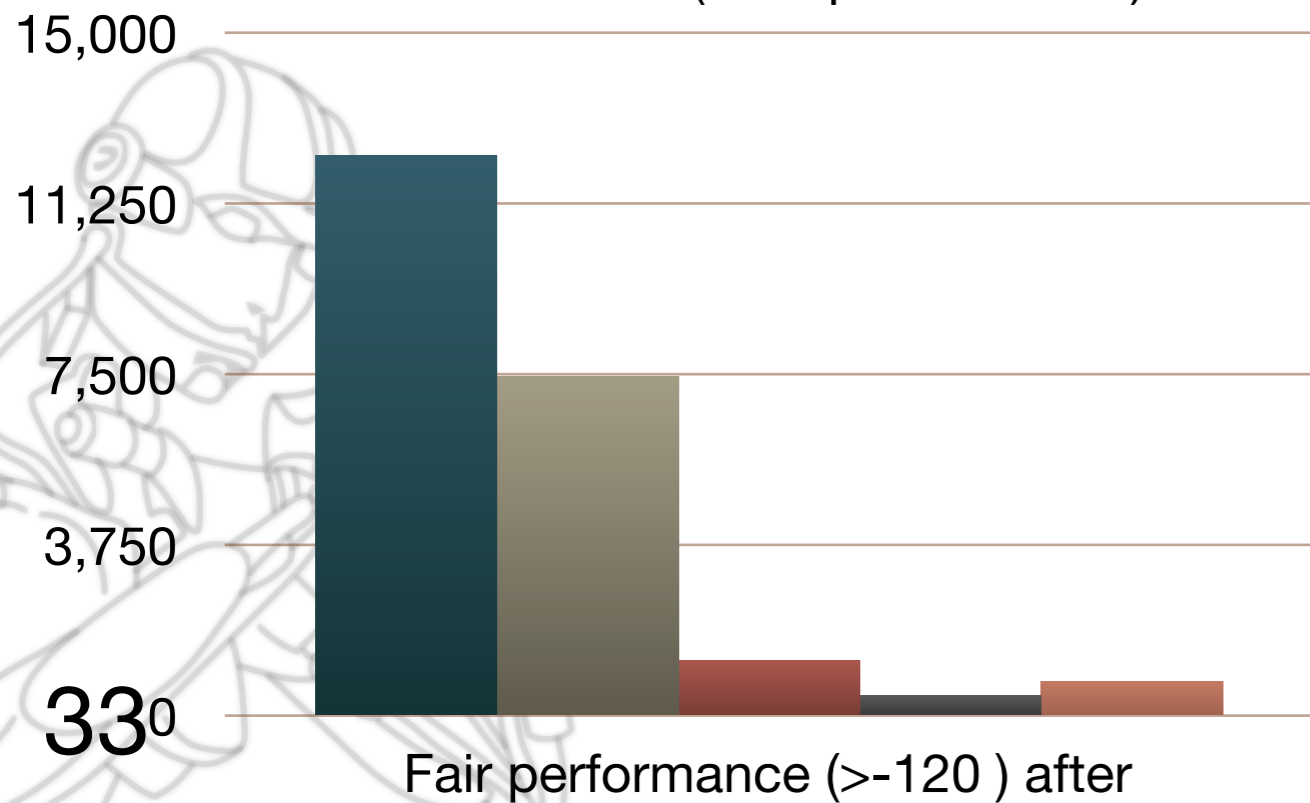
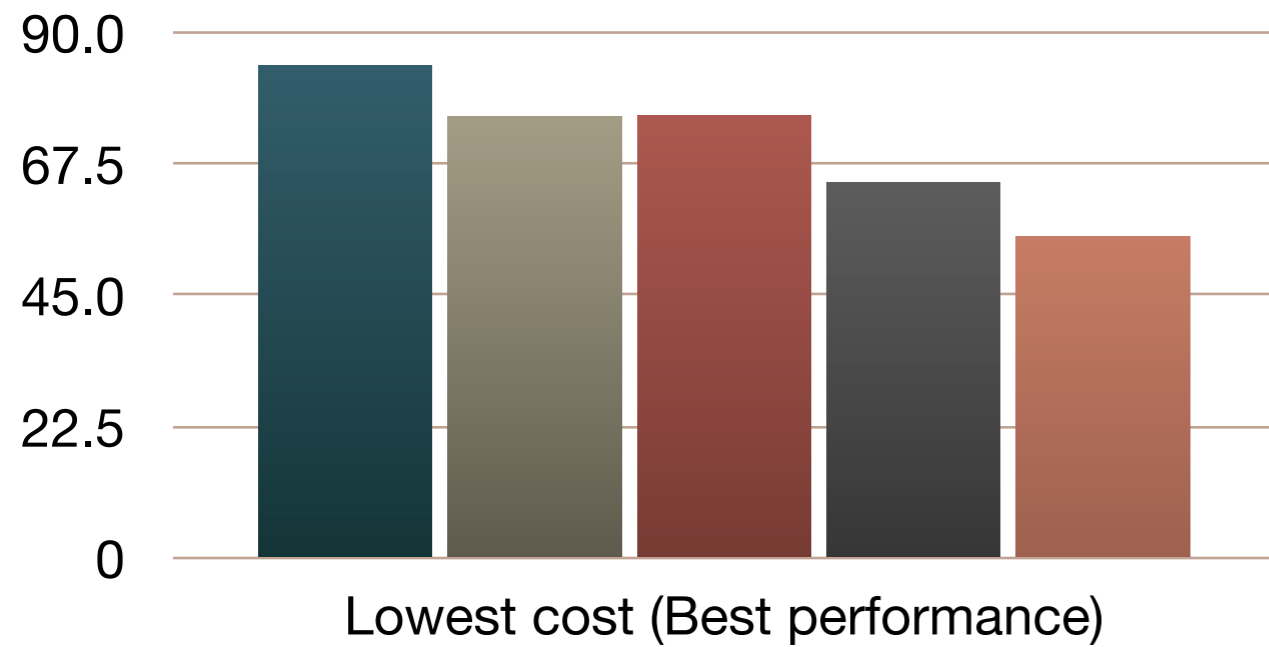
Episodic Natural Actor-Critic



Algorithm	Fair performance (>-120) after	Good performance (>-80) after	best performance Best Final
Episodic Natural Actor-Critic	750	5,050	Performance -55
Episodic Natural Actor-Critic with RPROP	Not reached	Not reached	-130

RPROP Updates do not seem to be compatible with the Episodic Natural Actor-Critic

Comparison of the Results



Natural Actor-Critic



Given: A parameterized stochastic policy (e.g., Gaussian)



Natural Actor-Critic



Given: A parameterized stochastic policy (e.g., Gaussian)

1. Perform trajectories and collect data.



Natural Actor-Critic



Given: A parameterized stochastic policy (e.g., Gaussian)

1. Perform trajectories and collect data.
2. Estimate the (natural) gradient using the compatible function approximation.



Natural Actor-Critic



Given: A parameterized stochastic policy (e.g., Gaussian)

1. Perform trajectories and collect data.
2. Estimate the (natural) gradient using the compatible function approximation.
3. Update the policy with gradient descent.

Natural Actor-Critic



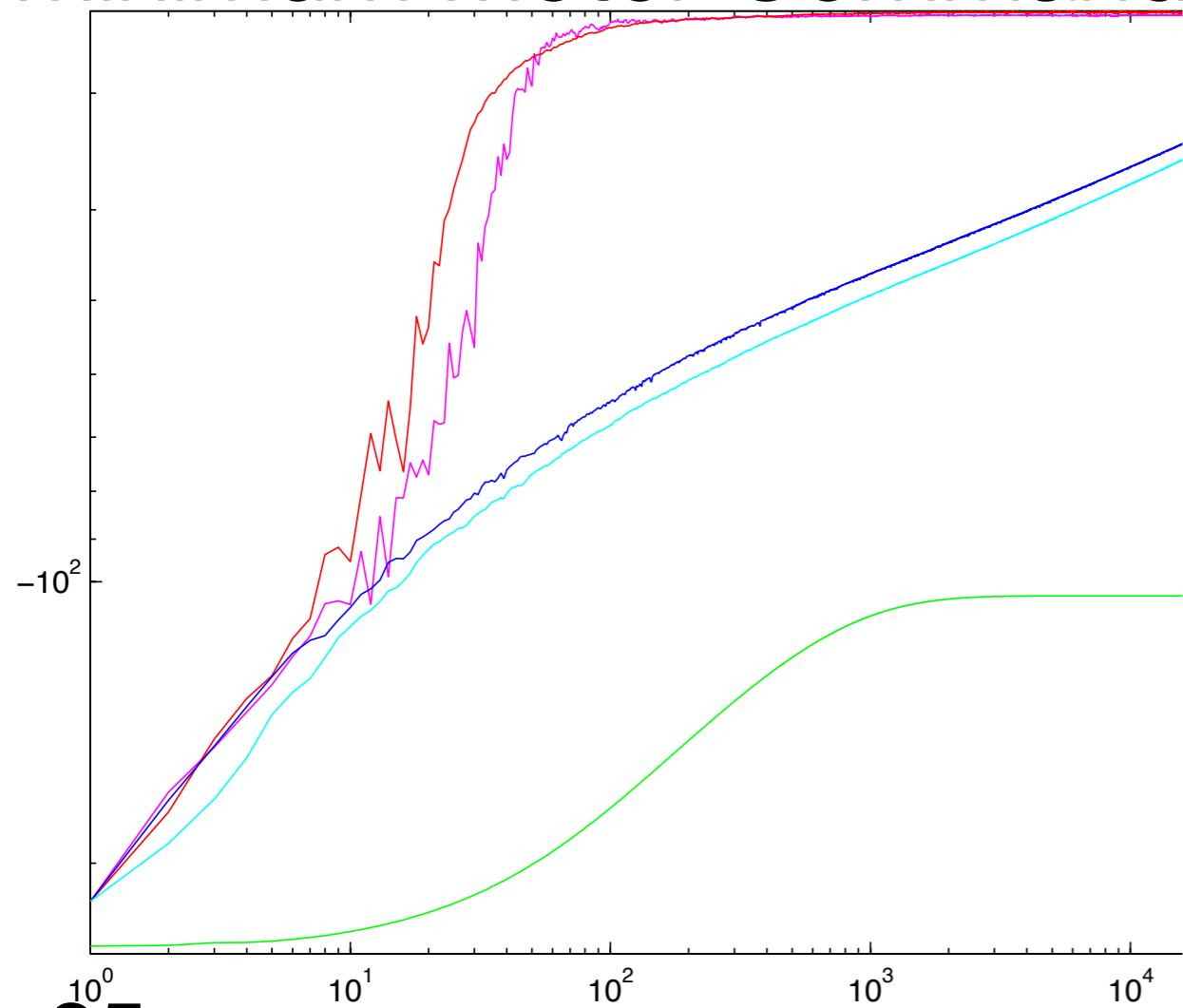
Given: A parameterized stochastic policy (e.g., Gaussian)

1. Perform trajectories and collect data.
2. Estimate the (natural) gradient using the compatible function approximation.
3. Update the policy with gradient descent.
4. Return to 1.

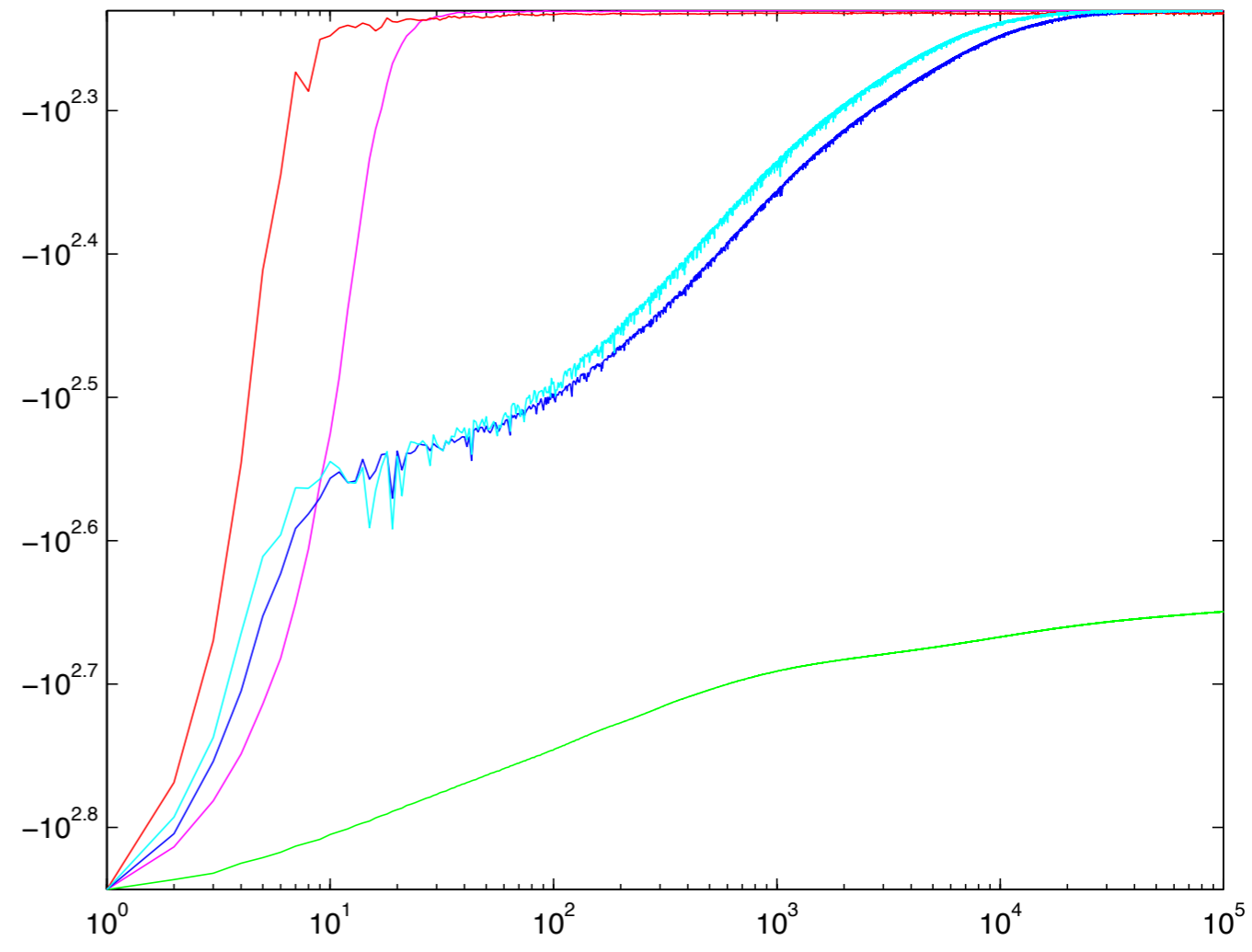
Improving MPs



Minimum Motor Command



Two Goals Policies



Learning T-Ball



Learning T-Ball



1) Teach motor primitives by imitation



Learning T-Ball



- 1) Teach motor primitives by imitation
- 2) Improve movement by Episodic Natural-Actor Critic



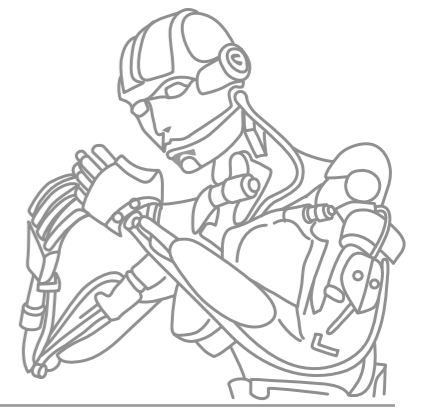
Learning T-Ball



- 1) Teach motor primitives by imitation
- 2) Improve movement by Episodic Natural-Actor Critic

*Good
performance
often after
150-300
trials.*





Outline of the Lecture

1. Introduction with Policy Gradients

2. Recent Advances in Policy Gradients

▶ 3. Probabilistic Policy Search with EM-like Approaches

4. Conclusion



Objective & Assumptions

Objective: maximize expected return

$$J(\boldsymbol{\theta}) = \int_{\mathbb{T}} p(\boldsymbol{\tau}) R(\boldsymbol{\tau}) d\boldsymbol{\tau}$$

Assumptions: Markovian & accumulated reward

path distribution

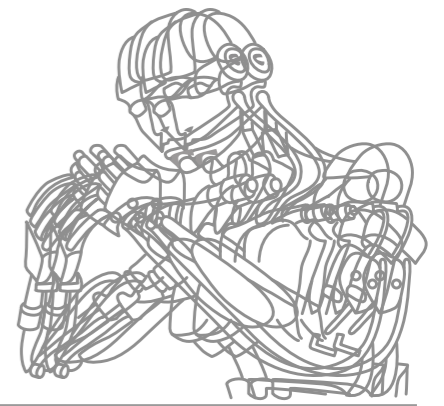
$$p(\boldsymbol{\tau}) = p(\mathbf{x}_1) \prod_{t=1}^T p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) \pi(\mathbf{u}_t | \mathbf{x}_t)$$

return

$$R(\boldsymbol{\tau}) = \frac{1}{T} \sum_{t=1}^T r(\mathbf{x}_t, \mathbf{u}_t)$$



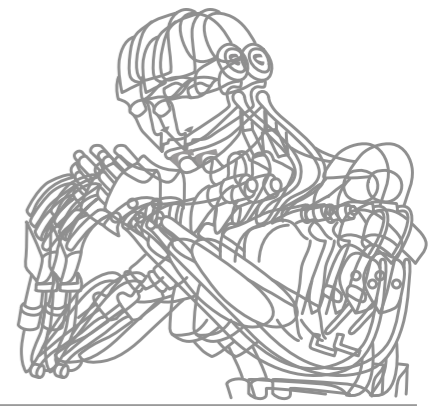
Success Matching Principle



“When learning from a set of their own trials in iterated decision problems, humans attempt to match not the best taken action but the reward-weighted frequency of their actions and outcomes” (Arrow, 1958).

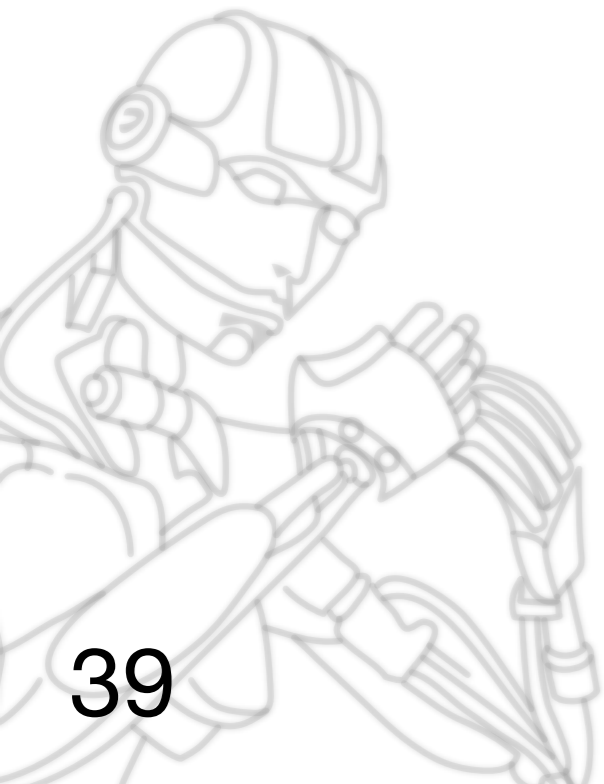


Success Matching Principle



“When learning from a set of their own trials in iterated decision problems, humans attempt to match not the best taken action but the reward-weighted frequency of their actions and outcomes” (Arrow, 1958).

Thus, why don't we create policies such that $\pi'(\mathbf{u}|\mathbf{x})$ matches $\pi(\mathbf{u}|\mathbf{x})r(\mathbf{x}, \mathbf{u})$?
(Dayan & Hinton, 1998)

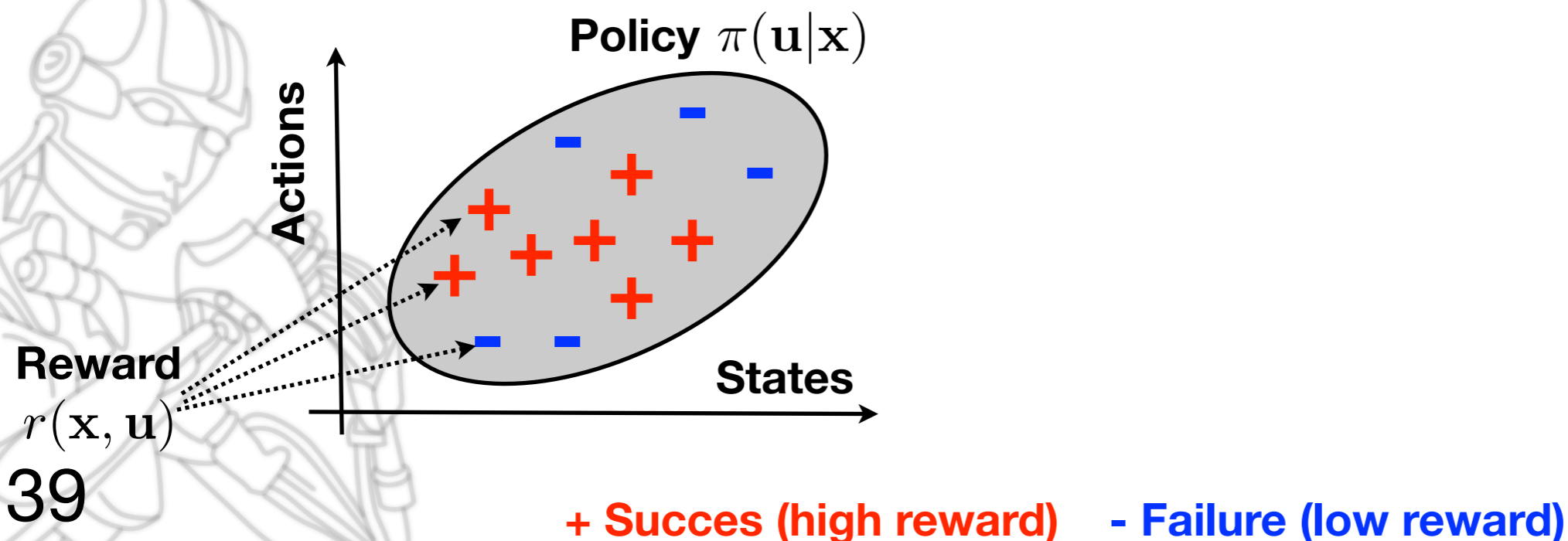


Success Matching Principle



“When learning from a set of their own trials in iterated decision problems, humans attempt to match not the best taken action but the reward-weighted frequency of their actions and outcomes” (Arrow, 1958).

Thus, why don't we create policies such that $\pi'(\mathbf{u}|\mathbf{x})$ matches $\pi(\mathbf{u}|\mathbf{x})r(\mathbf{x}, \mathbf{u})$?
(Dayan & Hinton, 1998)

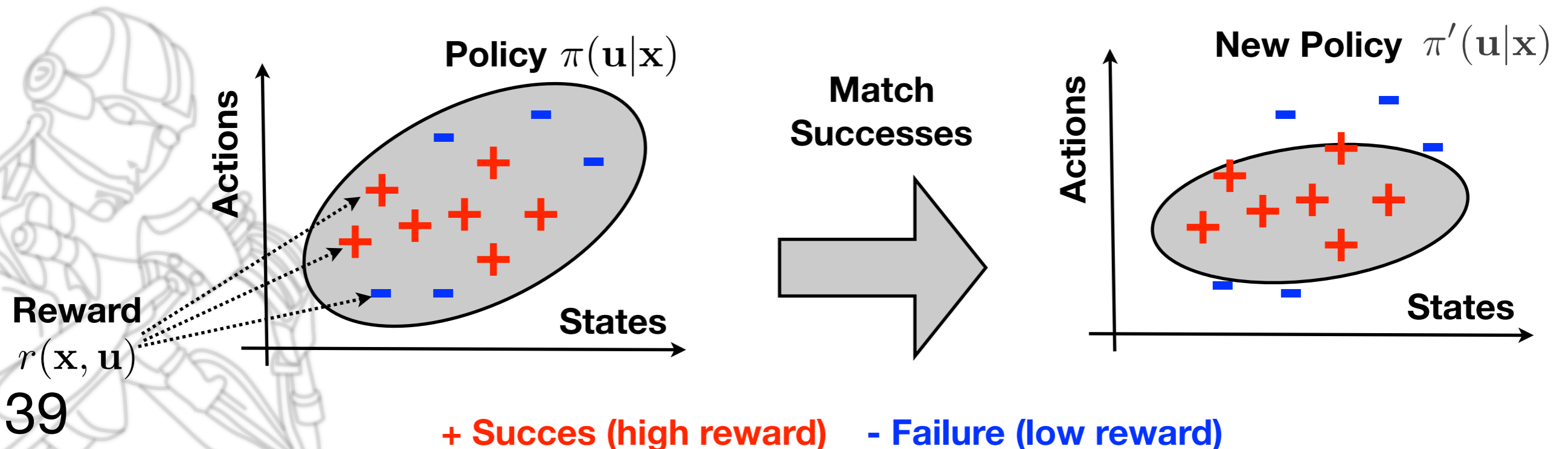


Success Matching Principle

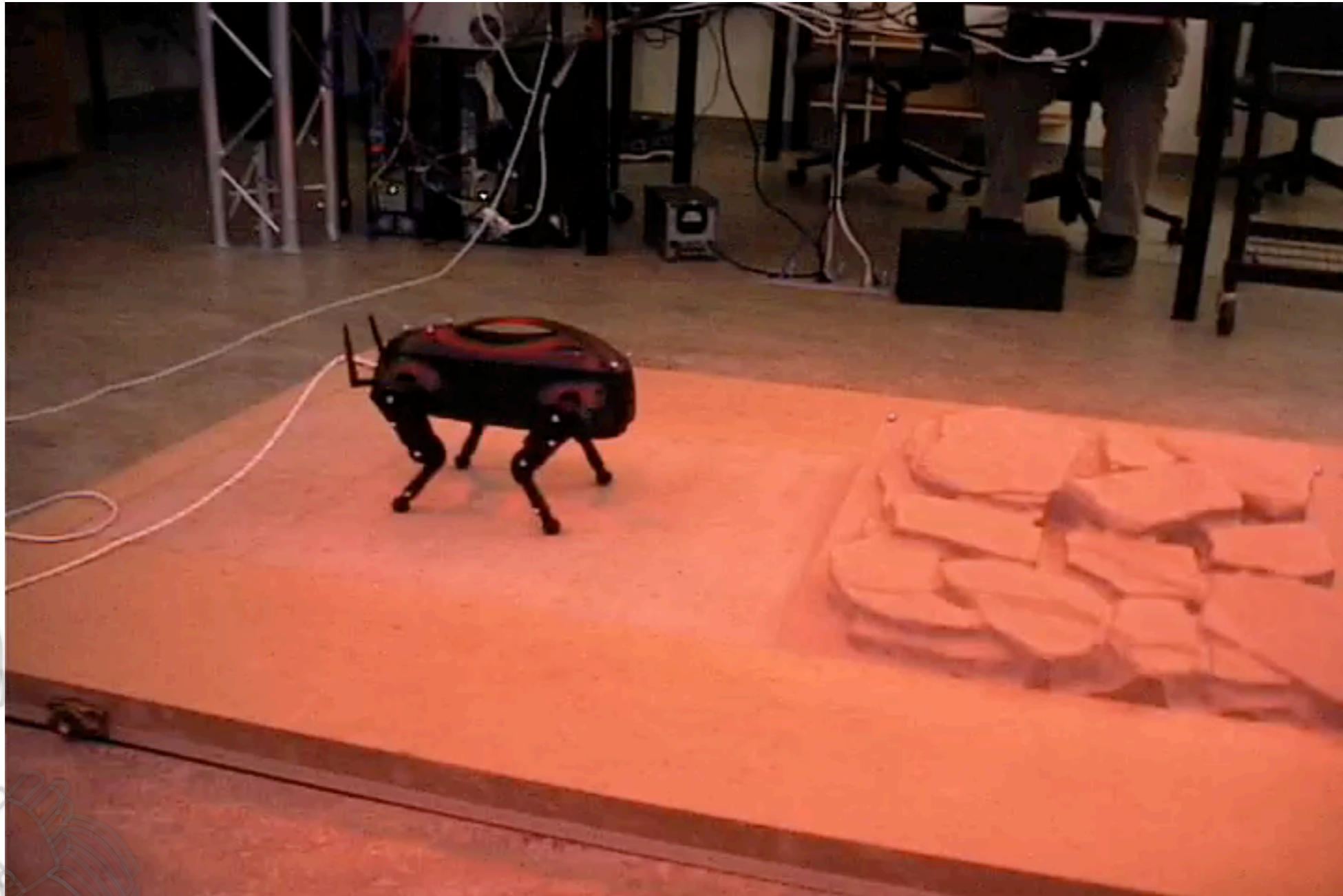


“When learning from a set of their own trials in iterated decision problems, humans attempt to match not the best taken action but the reward-weighted frequency of their actions and outcomes” (Arrow, 1958).

Thus, why don't we create policies such that $\pi'(\mathbf{u}|\mathbf{x})$ matches $\pi(\mathbf{u}|\mathbf{x})r(\mathbf{x}, \mathbf{u})$?
(Dayan & Hinton, 1998)



Selecting Footholds



Match successful footholds!

From Success Matching to Reward-Weighted Regression



Matching successful actions corresponds to minimizing the Kullback-Leibler 'distance'

$$D(r(\mathbf{x}, \mathbf{u})\pi(\mathbf{u}|\mathbf{x})||\pi'(\mathbf{u}|\mathbf{x})) \rightarrow \min$$

or

$$D(p(\boldsymbol{\tau}|\pi)R(\boldsymbol{\tau})||p(\boldsymbol{\tau}|\pi')) \rightarrow \min$$

➔ This minimization can be shown to correspond to optimizing a lower bound on the expected return!



Basic Intuition



42



Basic Intuition



- Lower Bound on Expected Return



Basic Intuition



- Lower Bound on Expected Return
 - reward is an improper probability distribution



Basic Intuition



- Lower Bound on Expected Return
 - reward is an improper probability distribution
 - log-likelihood \rightarrow log(expected return)
(Dayan & Hinton, Neural Computation 1997; Peters & Schaal, ICML 2007)

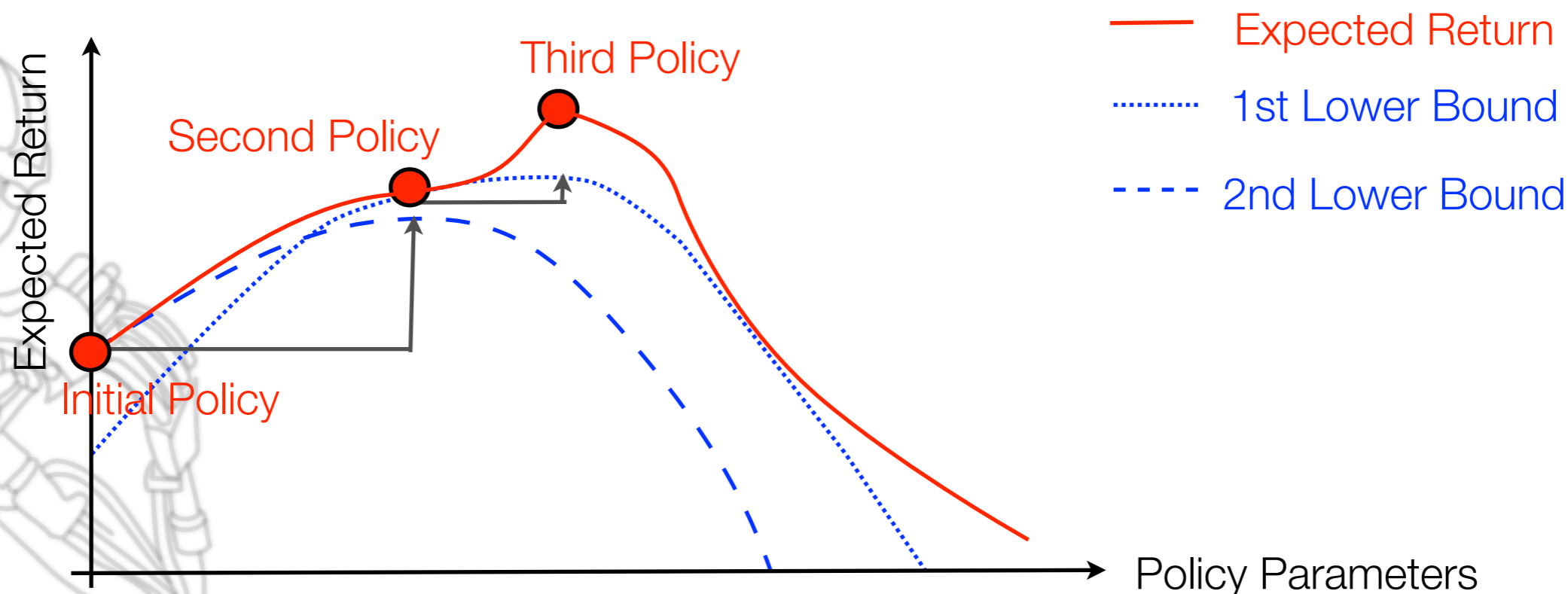




Basic Intuition

- Lower Bound on Expected Return
 - reward is an improper probability distribution
 - log-likelihood \rightarrow log(expected return)
(Dayan & Hinton, Neural Computation 1997; Peters & Schaal, ICML 2007)

$$\log J(\theta') \geq \int_{\mathbb{T}} p_{\theta}(\tau) R(\tau) \log \frac{p_{\theta'}(\tau)}{p_{\theta}(\tau)} d\tau + \text{const} = L_{\theta}(\theta')$$



Resulting Algorithms



Resulting Algorithms



Policy Gradients: maximize lower bound by following the gradient



Resulting Algorithms



Policy Gradients: maximize lower bound by following the gradient



Resulting Algorithms



Policy Gradients: maximize lower bound by following the gradient





Resulting Algorithms

Policy Gradients: maximize lower bound by following the gradient

EM-like Methods: maximize lower bound by expectation-maximization



Resulting Algorithms

Policy Gradients: maximize lower bound by following the gradient

$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$

EM-like Methods: maximize lower bound by expectation-maximization

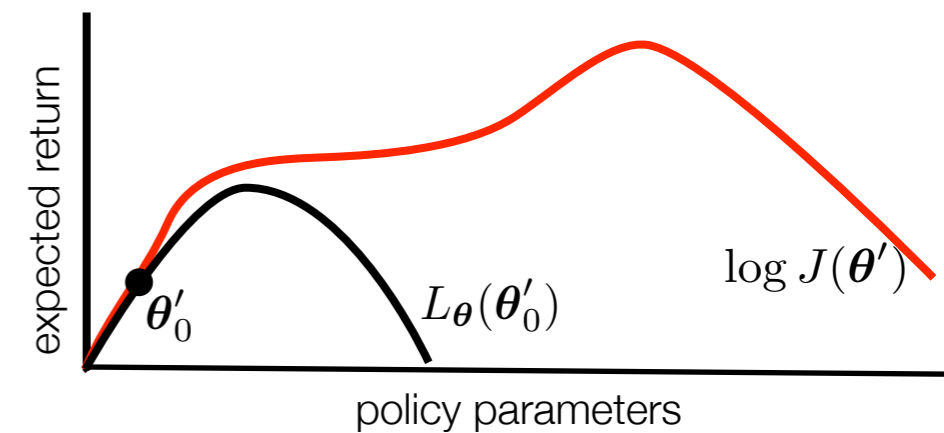
Resulting Algorithms



Policy Gradients: maximize lower bound by following the gradient

$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$



EM-like Methods: maximize lower bound by expectation-maximization

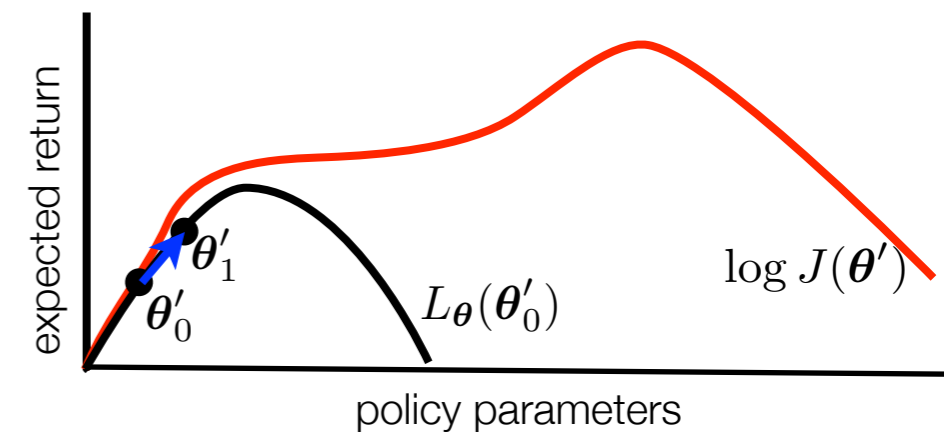
Resulting Algorithms



Policy Gradients: maximize lower bound by following the gradient

$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$



EM-like Methods: maximize lower bound by expectation-maximization

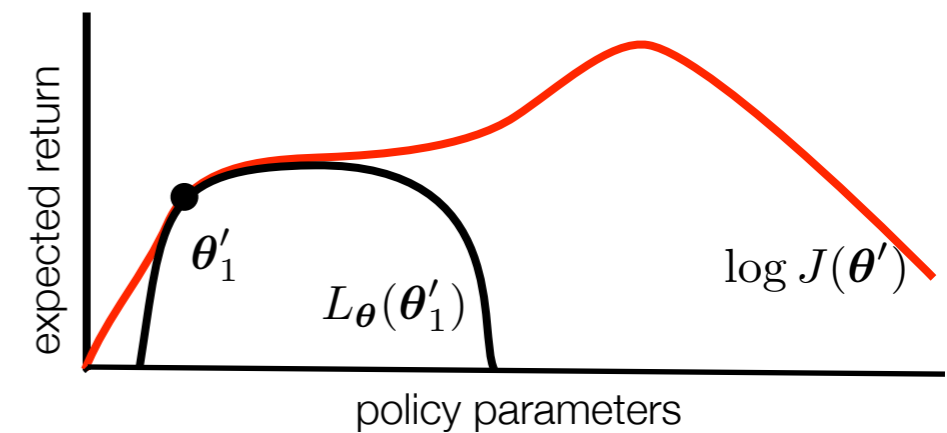
Resulting Algorithms



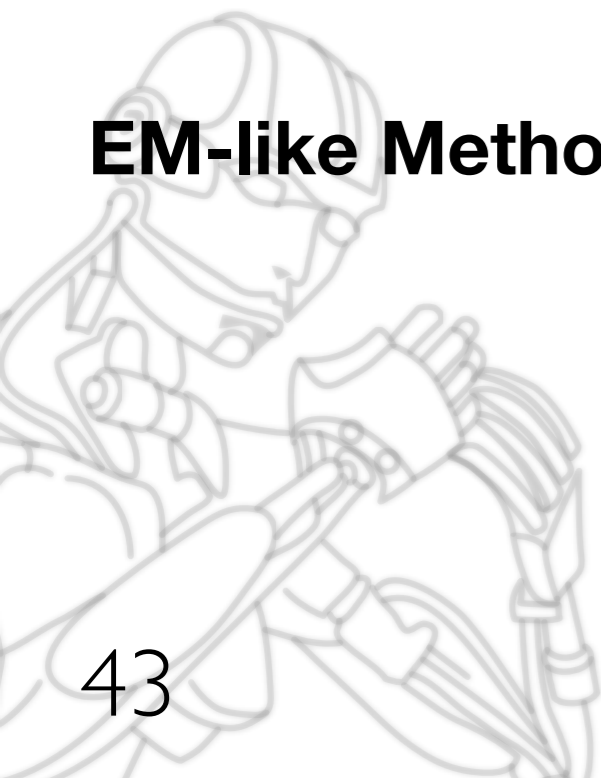
Policy Gradients: maximize lower bound by following the gradient

$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$



EM-like Methods: maximize lower bound by expectation-maximization



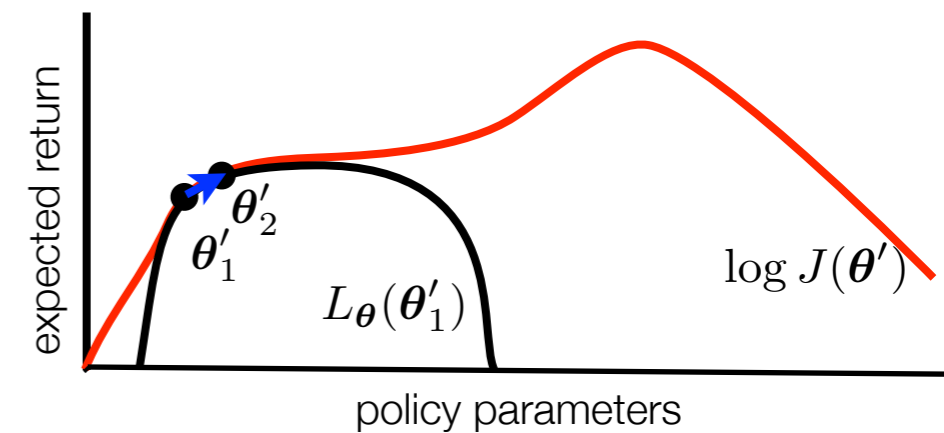
Resulting Algorithms



Policy Gradients: maximize lower bound by following the gradient

$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$



EM-like Methods: maximize lower bound by expectation-maximization

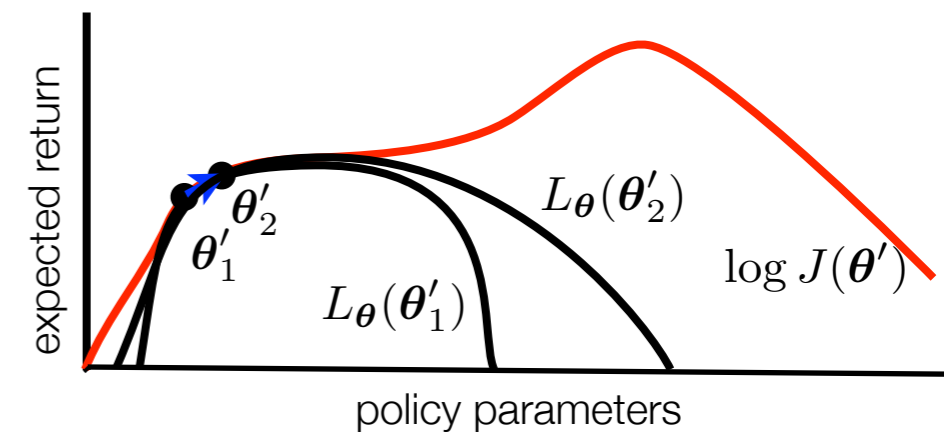
Resulting Algorithms



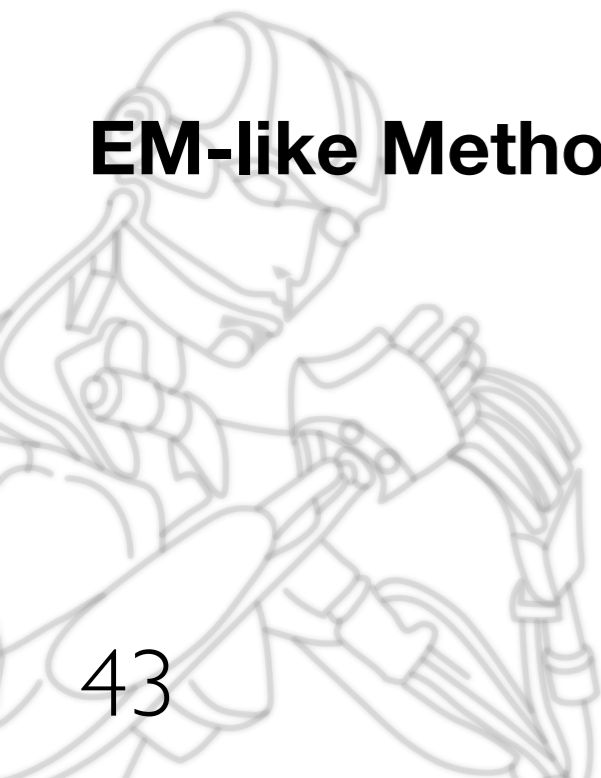
Policy Gradients: maximize lower bound by following the gradient

$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$



EM-like Methods: maximize lower bound by expectation-maximization



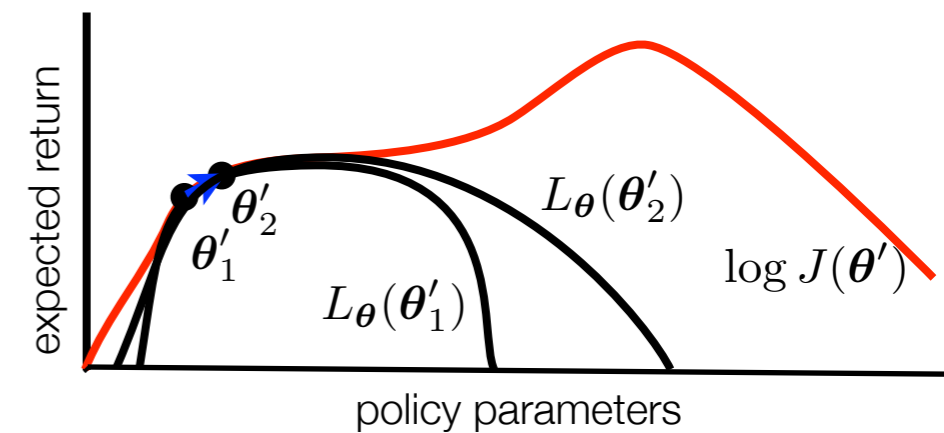
Resulting Algorithms



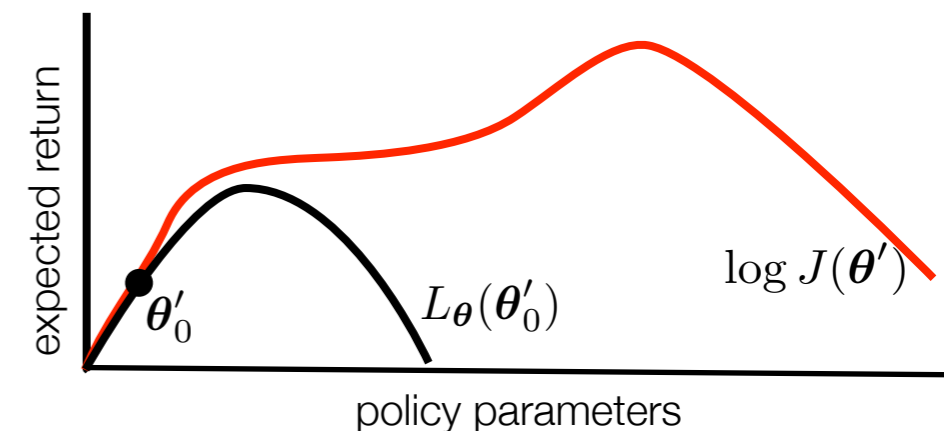
Policy Gradients: maximize lower bound by following the gradient

$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$



EM-like Methods: maximize lower bound by expectation-maximization



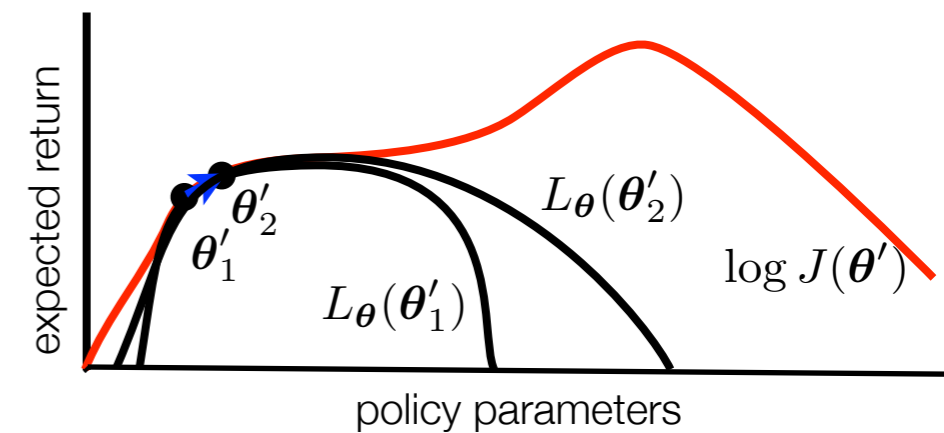
Resulting Algorithms



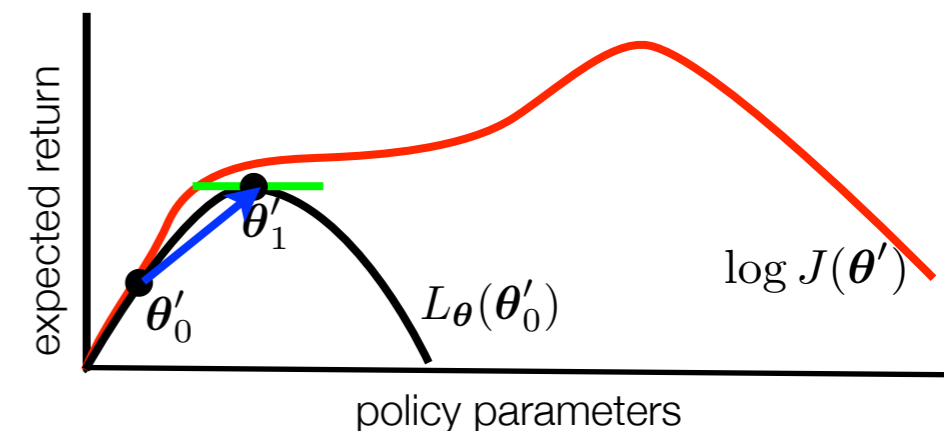
Policy Gradients: maximize lower bound by following the gradient

$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$



EM-like Methods: maximize lower bound by expectation-maximization



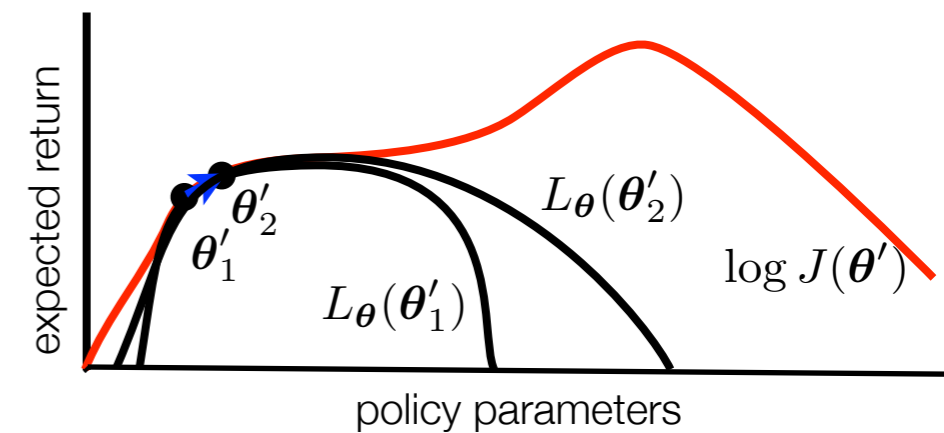
Resulting Algorithms



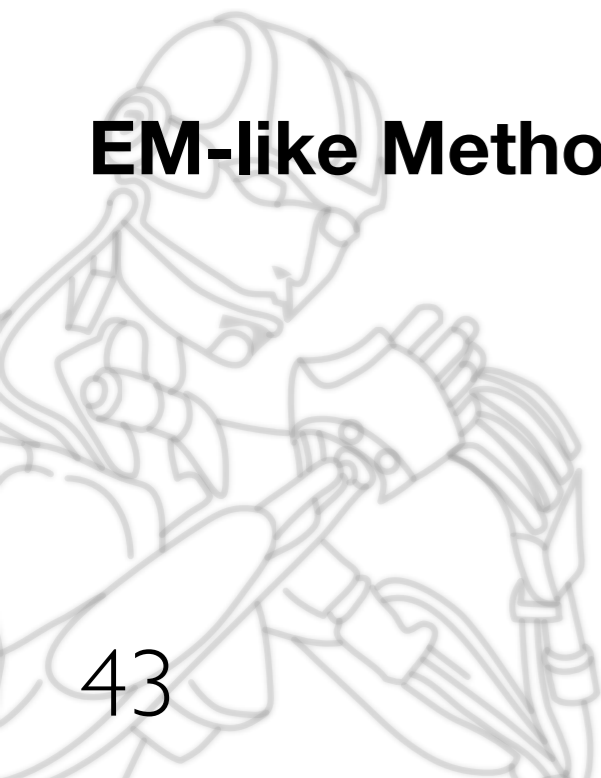
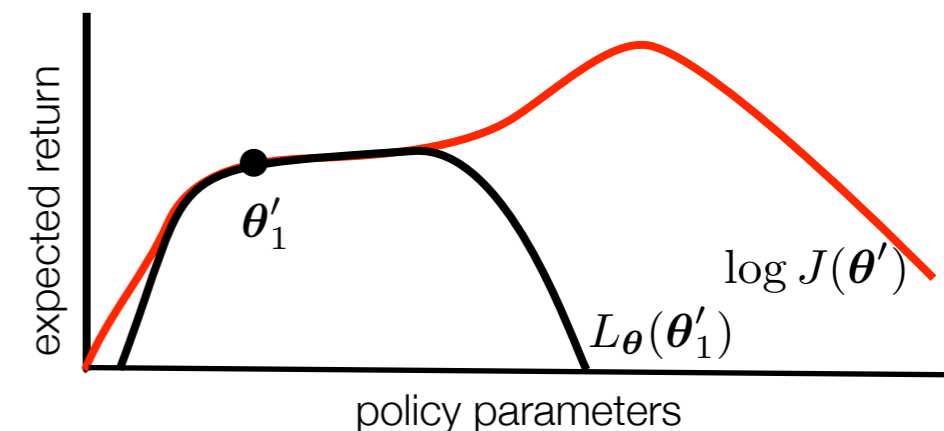
Policy Gradients: maximize lower bound by following the gradient

$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$



EM-like Methods: maximize lower bound by expectation-maximization



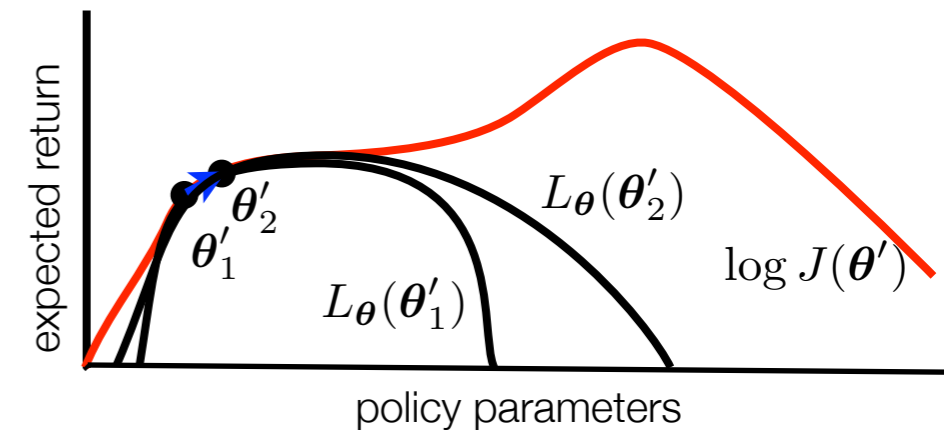
Resulting Algorithms



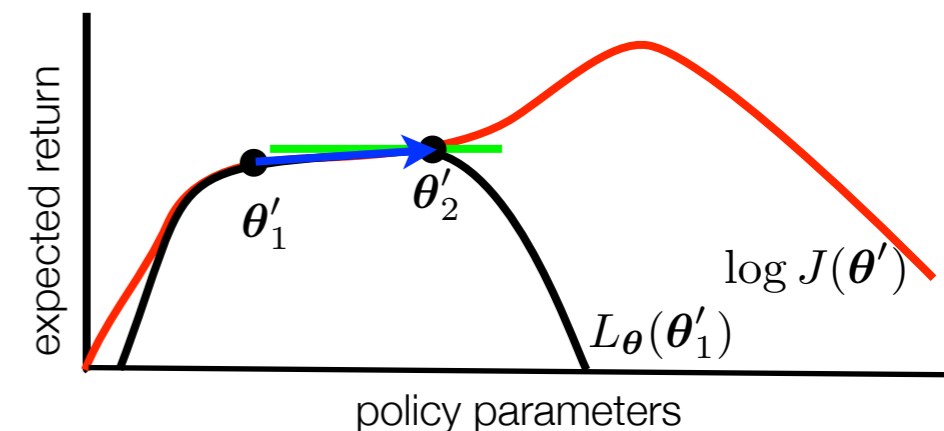
Policy Gradients: maximize lower bound by following the gradient

$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$



EM-like Methods: maximize lower bound by expectation-maximization



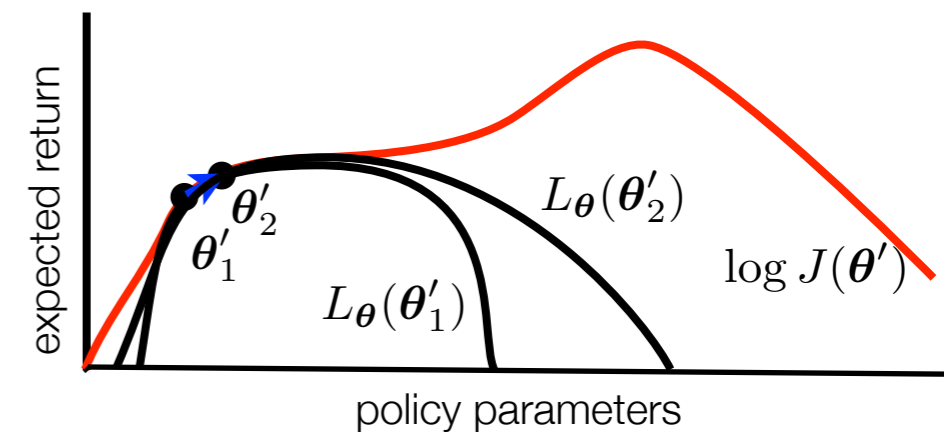
Resulting Algorithms



Policy Gradients: maximize lower bound by following the gradient

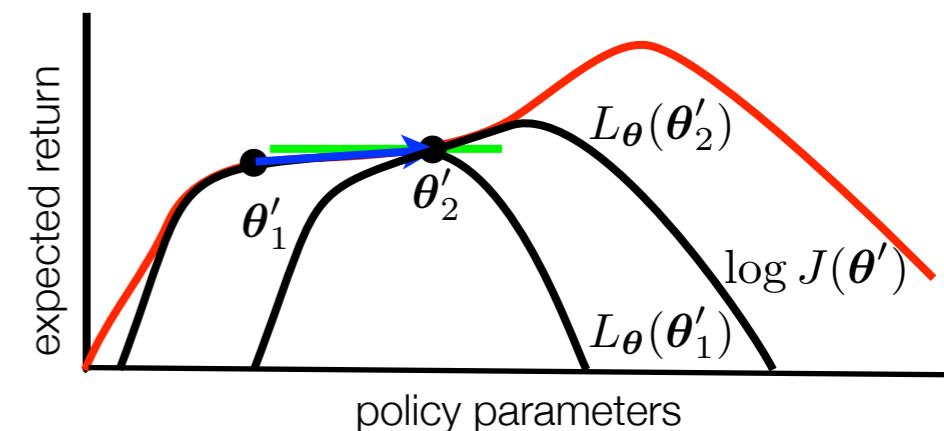
$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$



EM-like Methods: maximize lower bound by expectation-maximization

$$\theta' = \operatorname{argmax}_{\theta'} L_{\theta}(\theta')$$



Stochastic Policies



Use the **Policy**:

$$\mathbf{u} = f(\mathbf{x}) + \epsilon = \boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}) + \epsilon$$

with Gaussian exploration

$$\epsilon \sim \mathcal{N}(0, \sigma^2) \quad \rightarrow \text{episodic Reward Weighted Regression}$$

with State-dependent exploration

$$\epsilon = \boldsymbol{\varepsilon}^T \boldsymbol{\phi}(\mathbf{x}) \quad \text{with} \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2) \quad \rightarrow \text{PoWER}$$

Underactuated Swing-Up



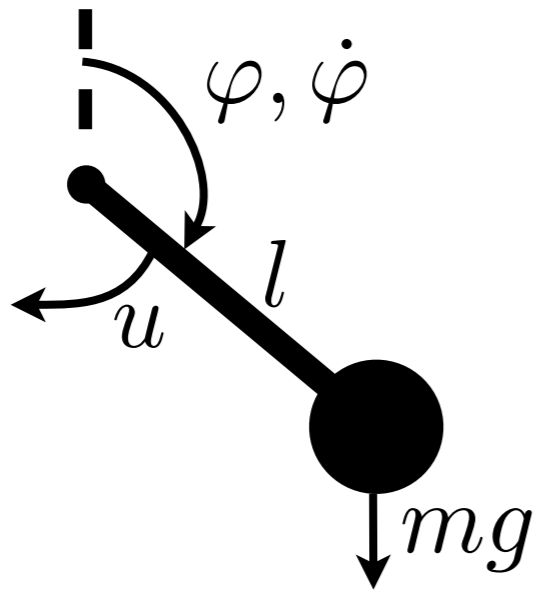
45

(Schaal, NIPS 1997; Atkeson, ICML 1997)

Underactuated Swing-Up



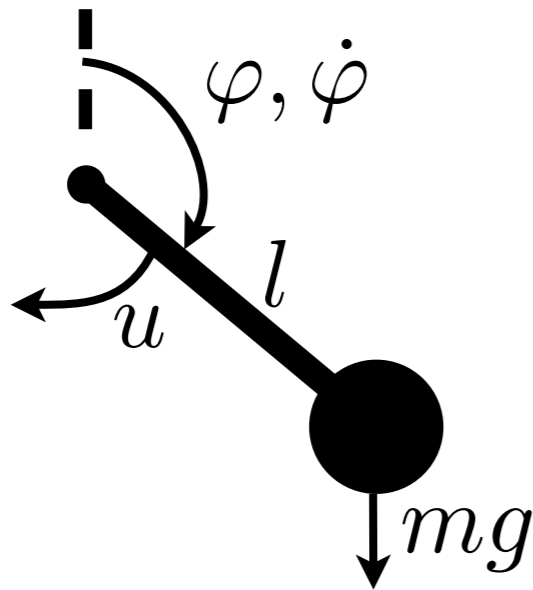
- swing heavy pendulum up



Underactuated Swing-Up



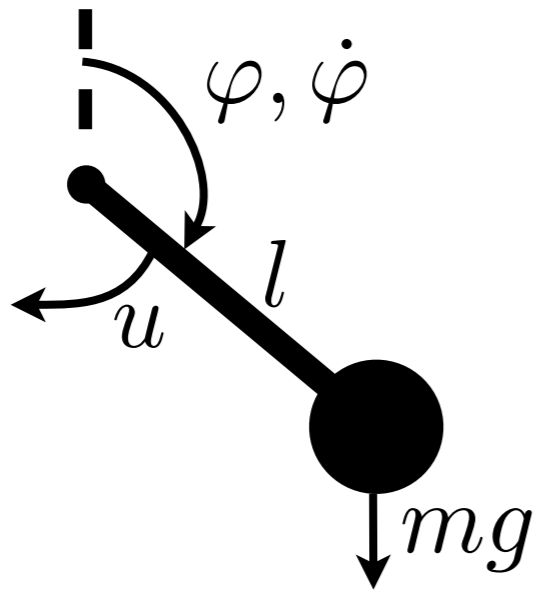
- swing heavy pendulum up



Underactuated Swing-Up



- swing heavy pendulum up



- motor torques limited

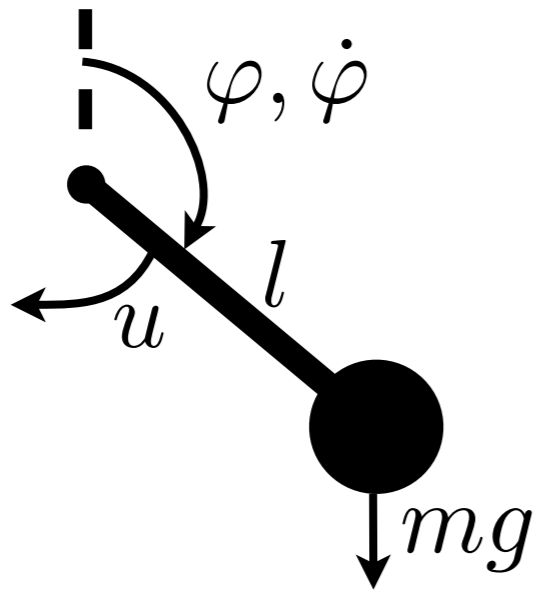
$$|u| \leq u_{max}$$



Underactuated Swing-Up



- swing heavy pendulum up



- motor torques limited

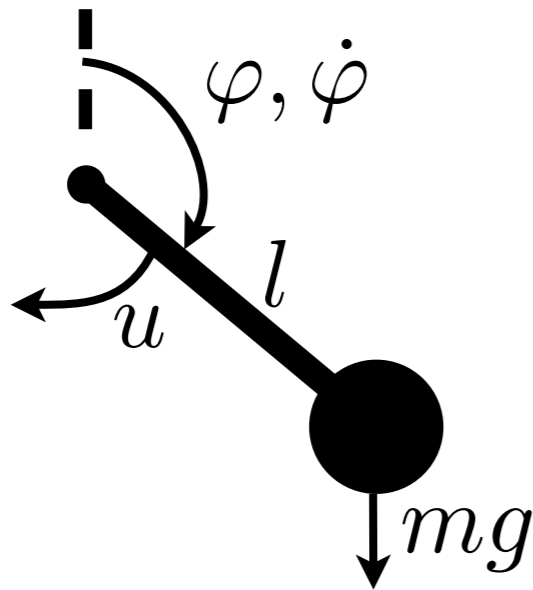
$$|u| \leq u_{max}$$



Underactuated Swing-Up



- swing heavy pendulum up



- motor torques limited

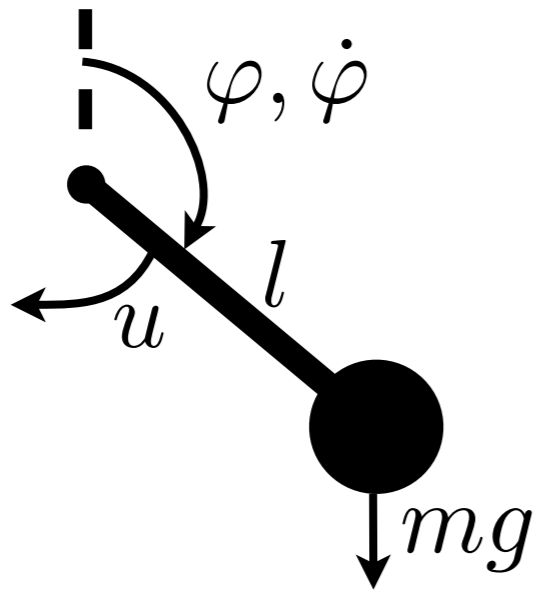
$$|u| \leq u_{max}$$



Underactuated Swing-Up



- swing heavy pendulum up



- motor torques limited

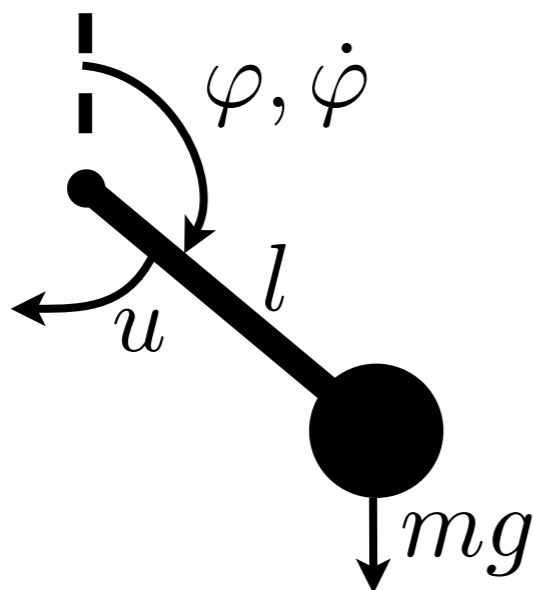
$$|u| \leq u_{max}$$





Underactuated Swing-Up

- swing heavy pendulum up



$$ml^2 \ddot{\varphi} = -\mu \dot{\varphi} + mgl \sin \varphi + u$$
$$\varphi \in [-\pi, \pi]$$

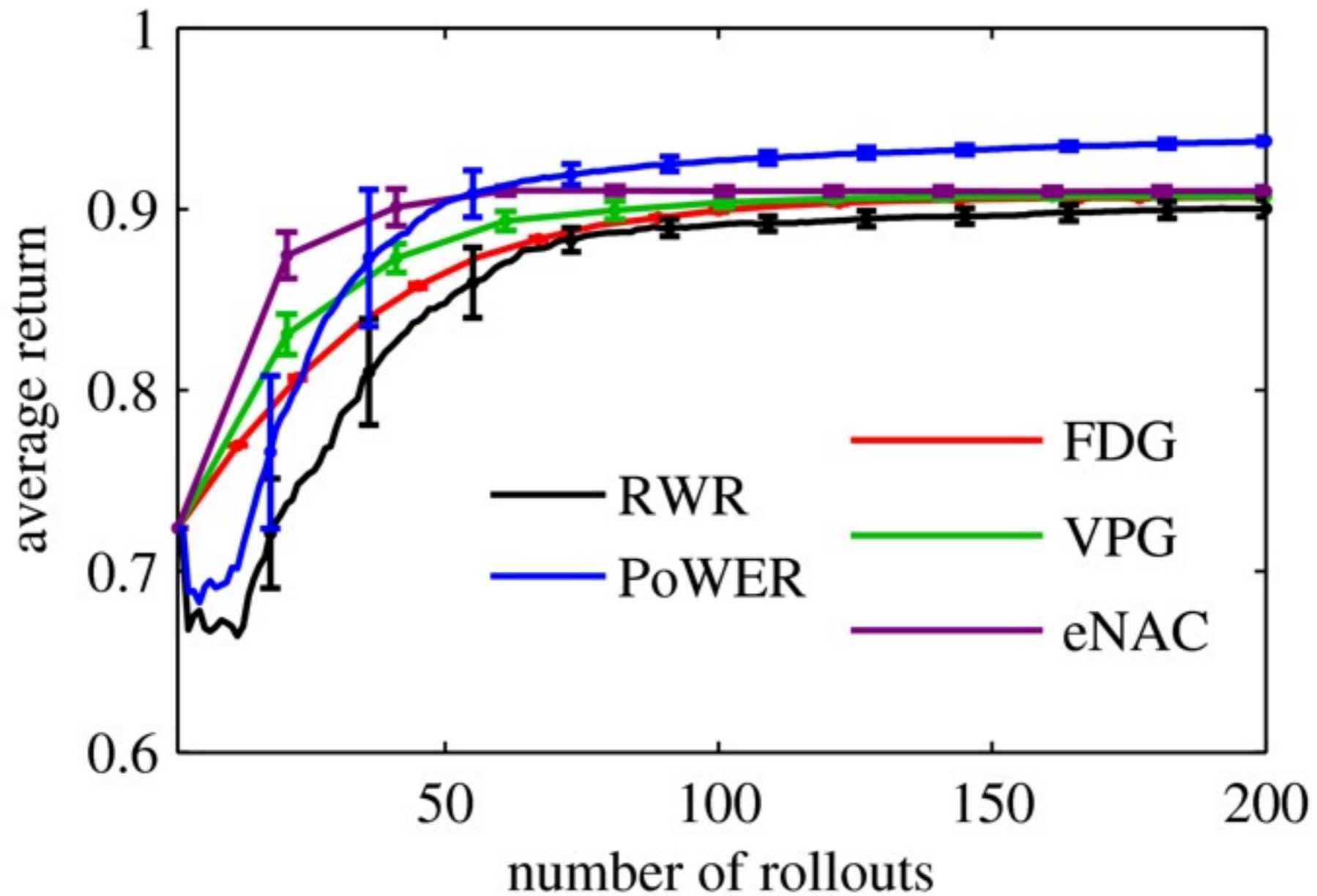
- motor torques limited

$$|u| \leq u_{max}$$

- reward function

$$r = \exp \left(-\alpha \left(\frac{\varphi}{\pi} \right)^2 - \beta \left(\frac{2}{\pi} \right)^2 \log \cos \left(\frac{\pi}{2} \frac{u}{u_{max}} \right) \right)$$

Underactuated Swing-Up



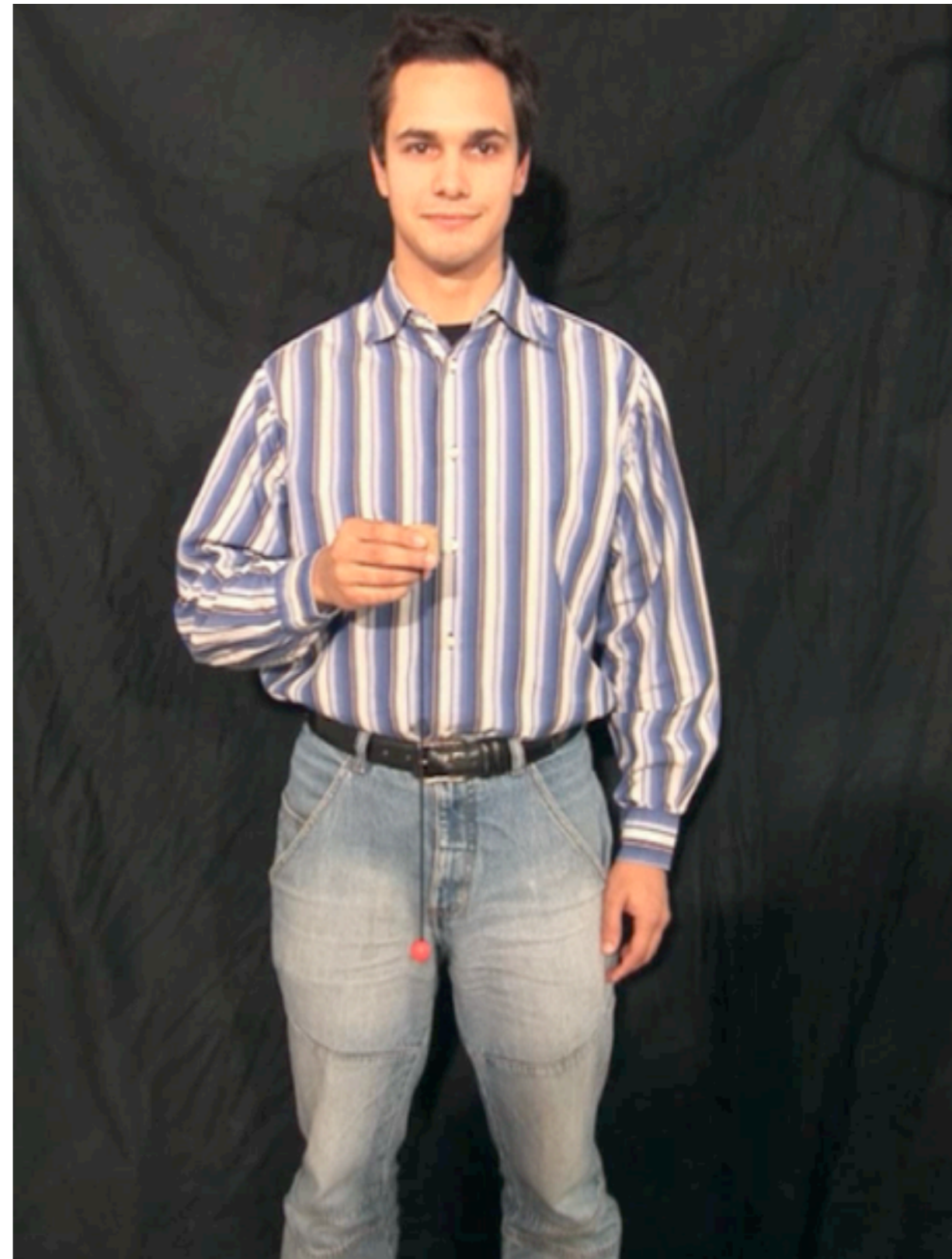
Ball-in-a-Cup



Ball-in-a-Cup



48

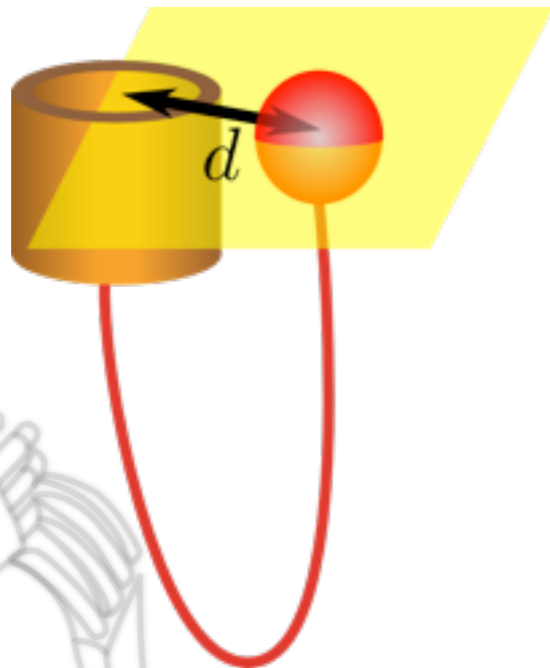


Ball-in-a-Cup



- reward function

$$r_t = \begin{cases} \exp\left(-\alpha\left((x_c - x_b)^2 + (y_c - y_b)^2\right)\right) & \text{if } t = t_c \\ 0 & \text{if } t \neq t_c \end{cases}$$

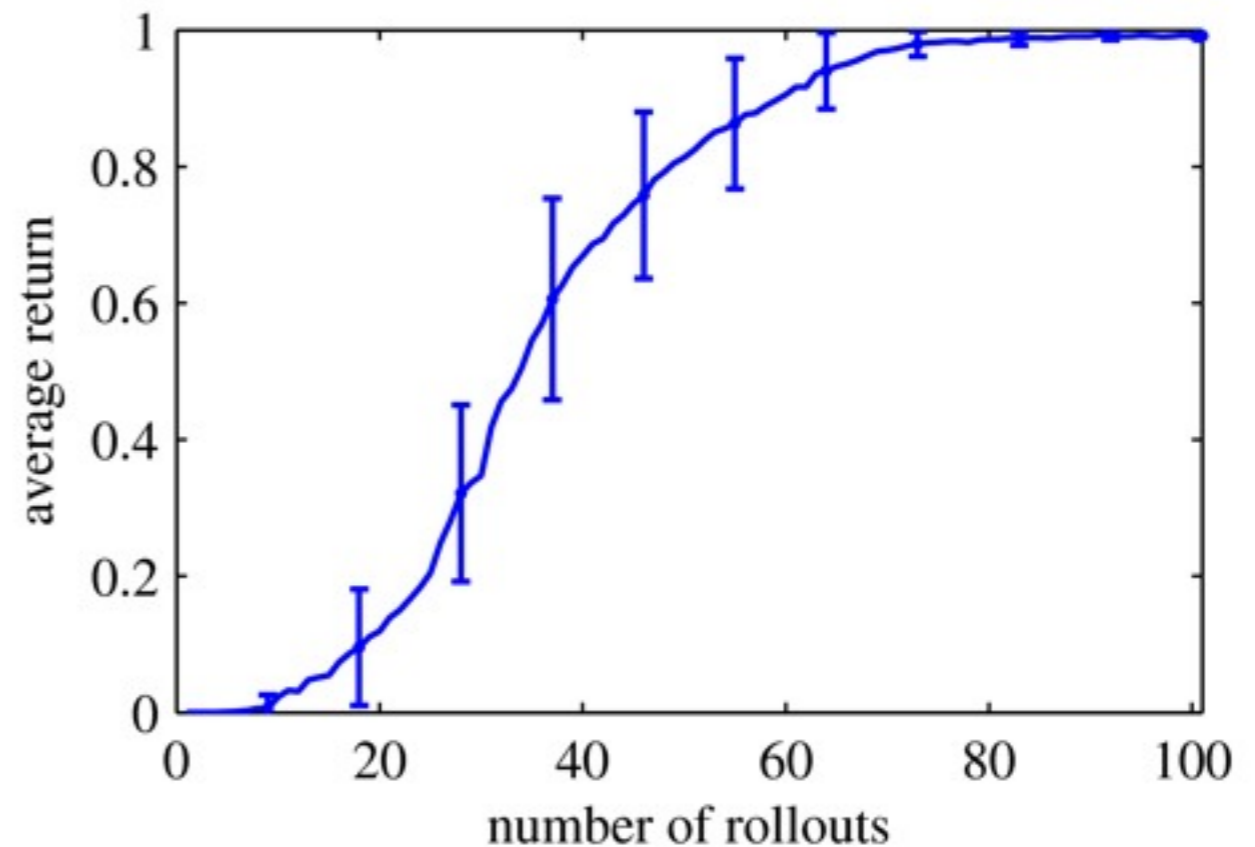
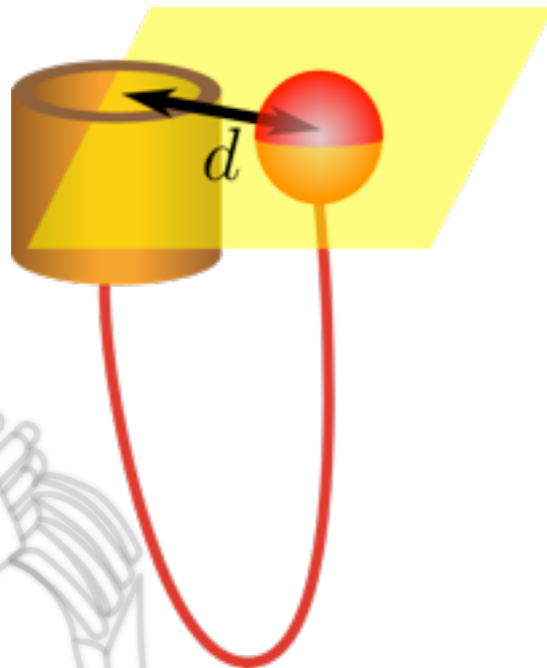


Ball-in-a-Cup



- reward function

$$r_t = \begin{cases} \exp\left(-\alpha\left((x_c - x_b)^2 + (y_c - y_b)^2\right)\right) & \text{if } t = t_c \\ 0 & \text{if } t \neq t_c \end{cases}$$



Cost-regularized Gaussian Processes



The Reward-Weighted Regression required known basis functions.
Using the Kernel-Trick

$$\begin{aligned}\bar{\mathbf{u}}_i &= \phi(\mathbf{x})^T \mathbf{w} = \phi(\mathbf{x})^T \left(\Phi^T \mathbf{R} \Phi + \lambda \mathbf{I} \right)^{-1} \Phi^T \mathbf{R} \mathbf{U}_i \\ &= \phi(\mathbf{x})^T \Phi^T \left(\Phi \Phi^T + \lambda \mathbf{R}^{-1} \right)^{-1} \mathbf{U}_i\end{aligned}$$

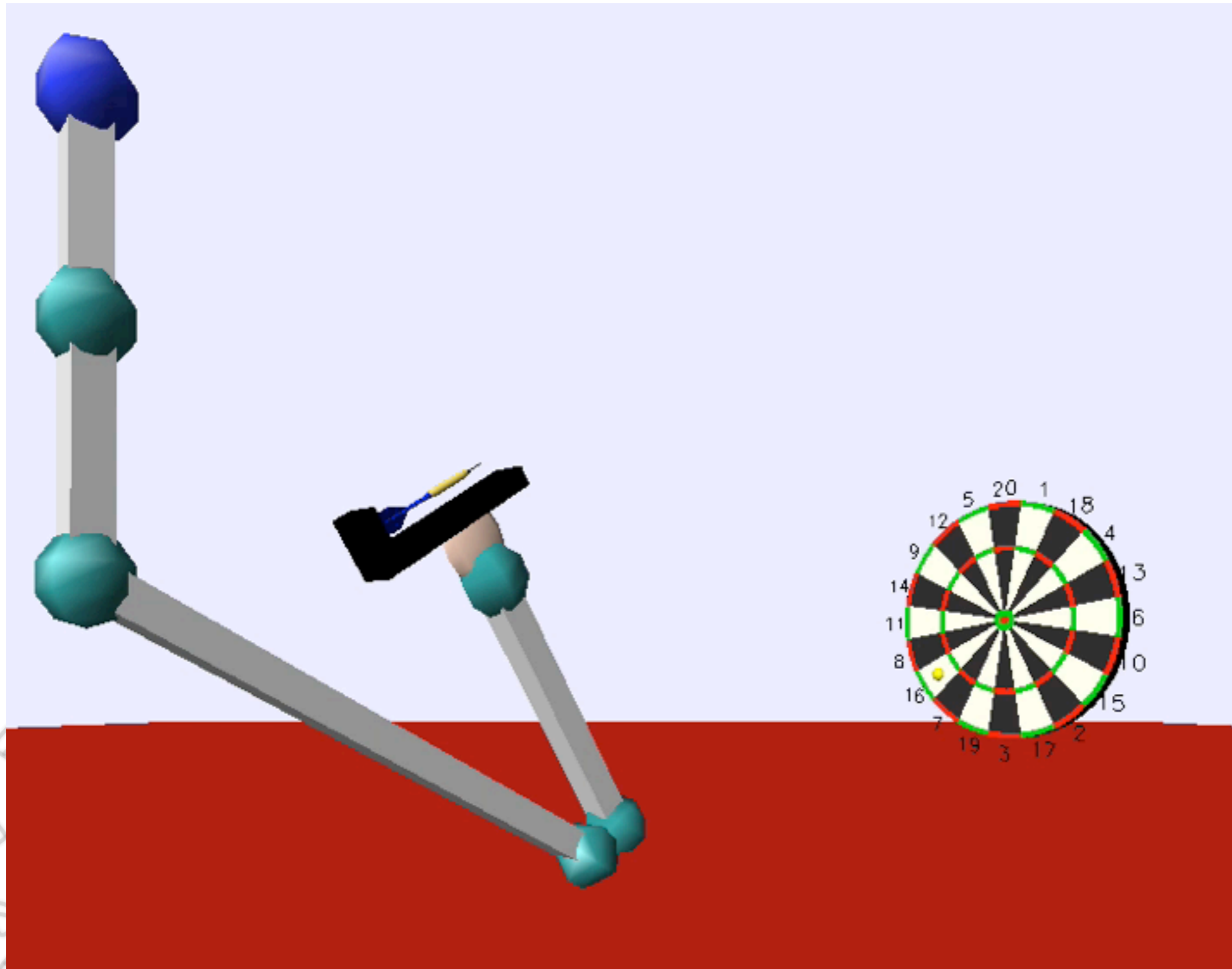
we can turn this in *a cost-regularized Gaussian Process approach*.
The predictive variance acts as a policy

$$\mathbf{u}^j \sim \pi_j(\mathbf{u} | \mathbf{x}^j) = \mathcal{N}(\mathbf{u} | \gamma(\mathbf{x}^j), \sigma^2(\mathbf{x}^j) \mathbf{I})$$

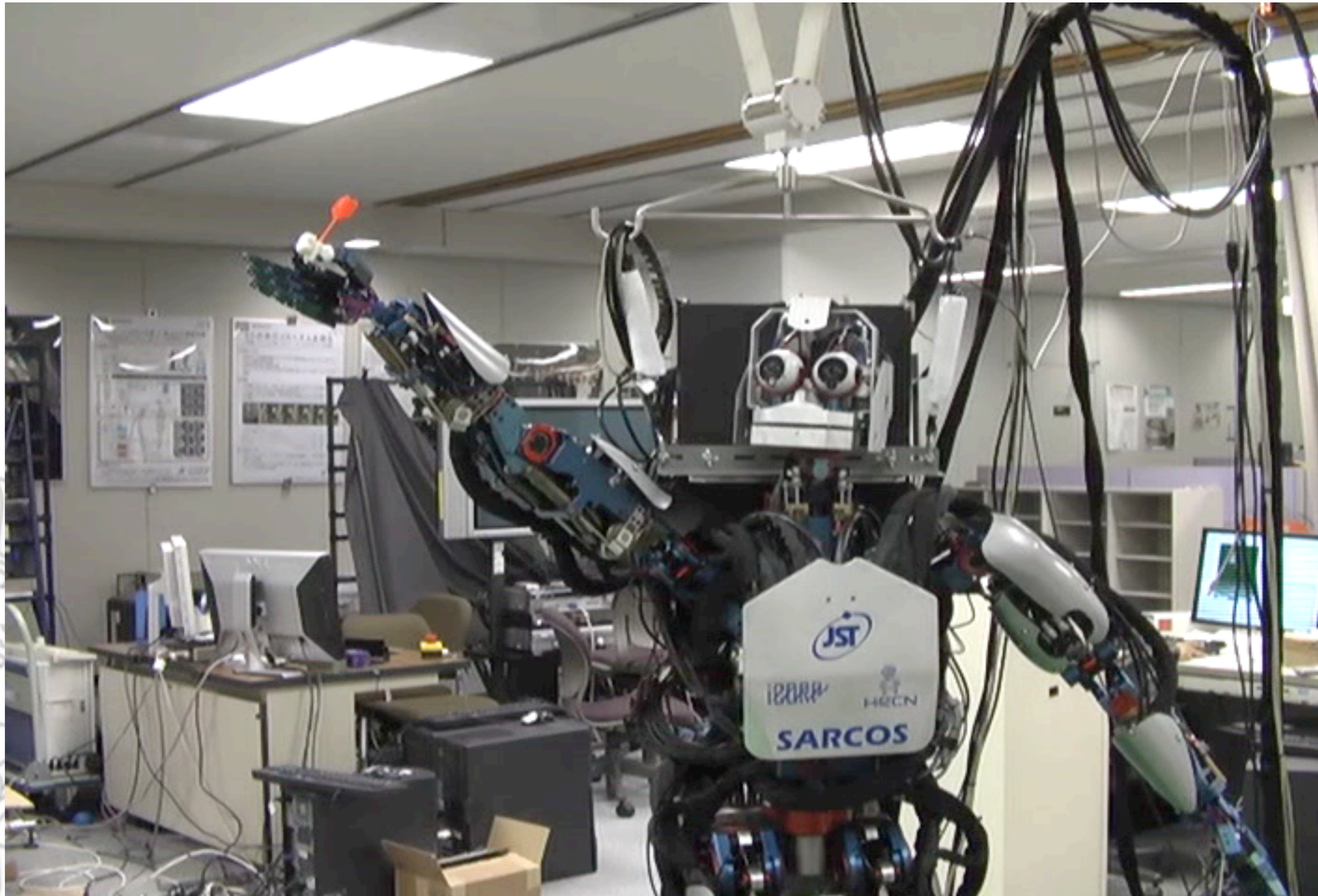
with

$$\begin{aligned}\gamma_i(\mathbf{x}^j) &= \mathbf{k}(\mathbf{x}^j)^T (\mathbf{K} + \lambda \mathbf{C})^{-1} \mathbf{U}_i \\ \sigma^2(\mathbf{x}^j) &= k(\mathbf{x}^j, \mathbf{x}^j) - \mathbf{k}(\mathbf{x}^j)^T (\mathbf{K} + \lambda \mathbf{C})^{-1} \mathbf{k}(\mathbf{x}^j)\end{aligned}$$

Dart-Throwing with Sledge



Dart-Throwing with Fingers



Learning for Table Tennis



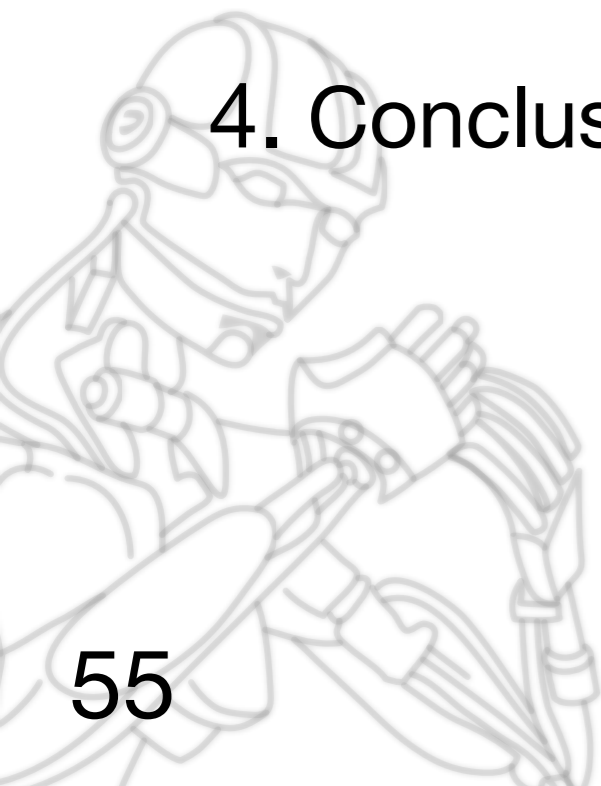
Throwing and Catching





Outline of the Lecture

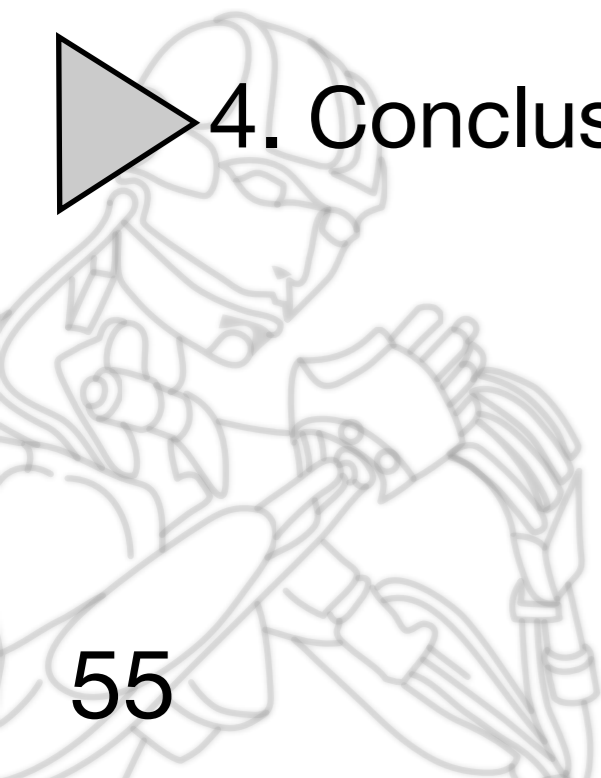
1. Introduction with Policy Gradients
2. Recent Advances in Policy Gradients
3. Probabilistic Policy Search with EM-like Approaches
4. Conclusion





Outline of the Lecture

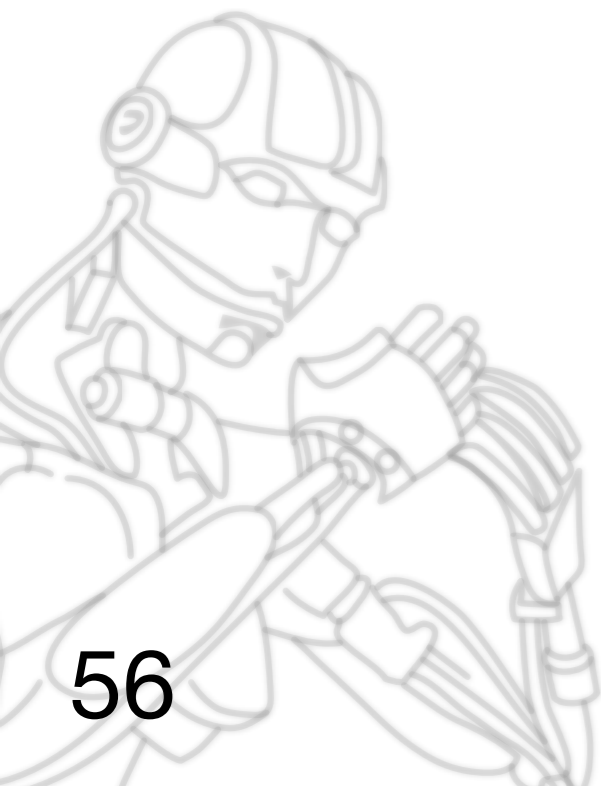
1. Introduction with Policy Gradients
2. Recent Advances in Policy Gradients
3. Probabilistic Policy Search with EM-like Approaches
- ▶ 4. Conclusion





Conclusion

- Policy Search is a powerful and practical alternative to value function and model-based methods.
- Policy gradients have dominated this area for a long time and solidly working methods exist.
- Newer methods focus on probabilistic policy search approaches.





Further Reading

- Peters, J.; Schaal, S. (2008). Reinforcement learning of motor skills with policy gradients, *Neural Networks*, 21, 4, pp.682-97
- Kober, J.; Peters, J. (2011). Policy Search for Motor Primitives in Robotics, *Machine Learning*, 84, 1-2, pp.171-203
- Peters, J.; Muelling, K.; Altun, Y. (2010). Relative Entropy Policy Search, *Proceedings of the Twenty-Fourth National Conference on Artificial Intelligence (AAAI), Physically Grounded AI Track*