

Tutorial on Policy Gradient Methods

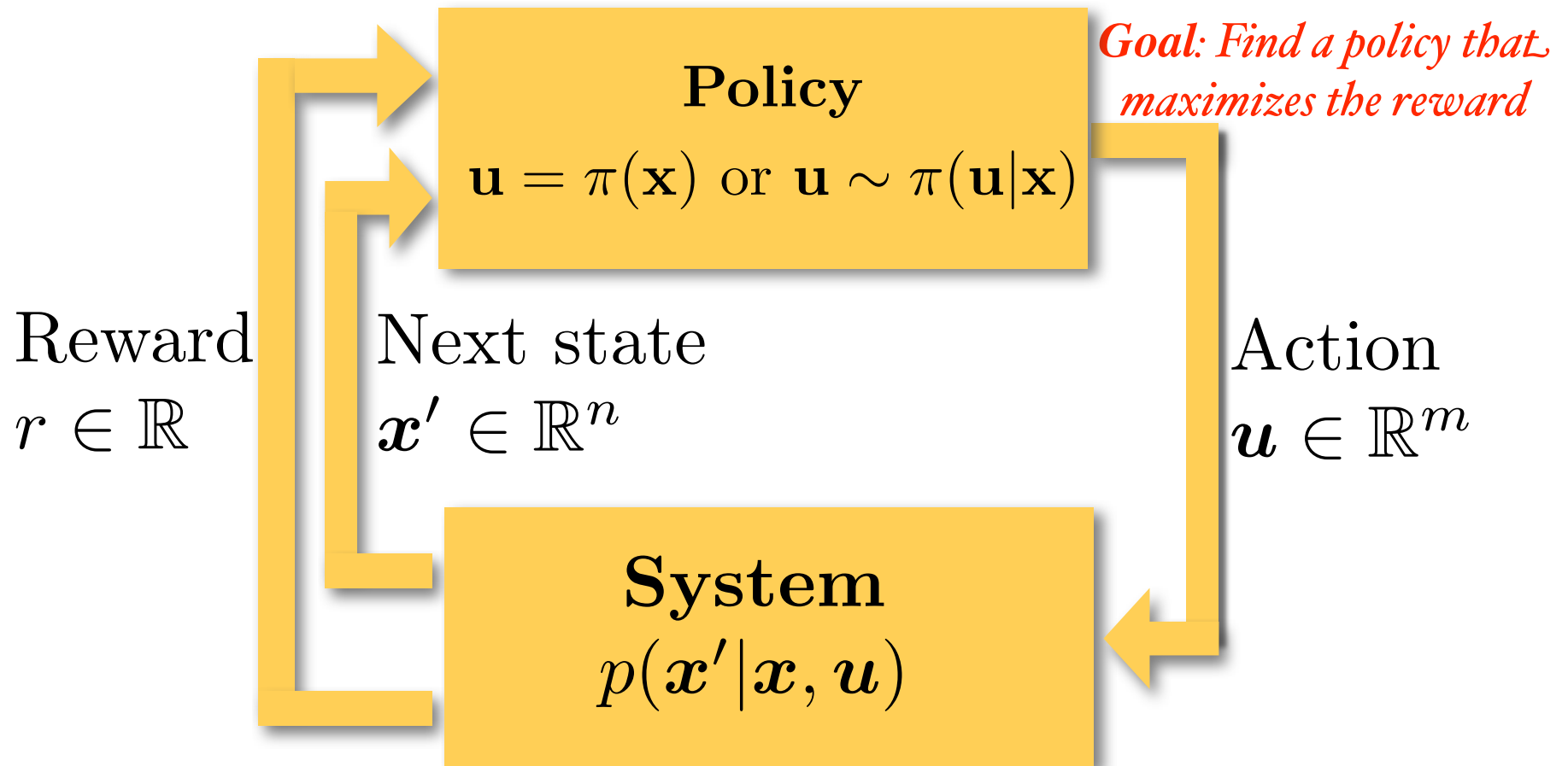
Jan Peters





- ▶ **1. Reinforcement Learning**
2. Finite Difference vs Likelihood-Ratio Policy Gradients
3. Likelihood-Ratio Policy Gradients
4. Conclusion

General Setup



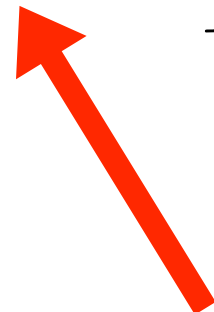
1. Reinforcement Learning

Goal of RL



What does maximizing your rewards mean?

$$J(\pi) = \frac{1}{T} \sum_{t=1}^T r(\mathbf{x}_t, \mathbf{u}_t, t) \rightarrow E\{r(\mathbf{x}, \mathbf{u}, t)\}$$



Find a policy that maximizes the rewards!

A policy tells you which actions to use for each state!

1. Reinforcement Learning

Policy Search vs Value Function Methods



Value Function View!

Critic: Policy Evaluation

$$Q(\mathbf{x}_t, \mathbf{u}_t, t) = E \left\{ \sum_{\tau=t}^T r(\mathbf{x}_\tau, \mathbf{u}_\tau, \tau) \mid \mathbf{x}_t, \mathbf{u}_t \right\}$$



Actor: Compute Optimal Policy

$$\mathbf{u}_t = \pi(\mathbf{x}_t, t) = \operatorname{argmax} Q(\mathbf{x}_t, \mathbf{u}, t)$$

Policy Search View!

Critic: Policy Sensitivity

$$J(\pi) = E \left\{ \sum_{t=0}^T r(\mathbf{x}_t, \mathbf{u}_t, t) \right\}$$



Actor: Policy Improvement

$$\pi' = \operatorname{argmax}_{\pi'} \{ J(\pi') - J(\pi) \}$$

1. Reinforcement Learning

Greedy vs Gradients



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*

2. Value Function Methods

Outline



1. Reinforcement Learning for Motor Control?

▶ **2. Finite Difference vs Likelihood-Ratio Policy Gradients**

3. Likelihood-Ratio Policy Gradients

4. Conclusion

Why Policy Gradient Methods?



Why Policy Gradients?

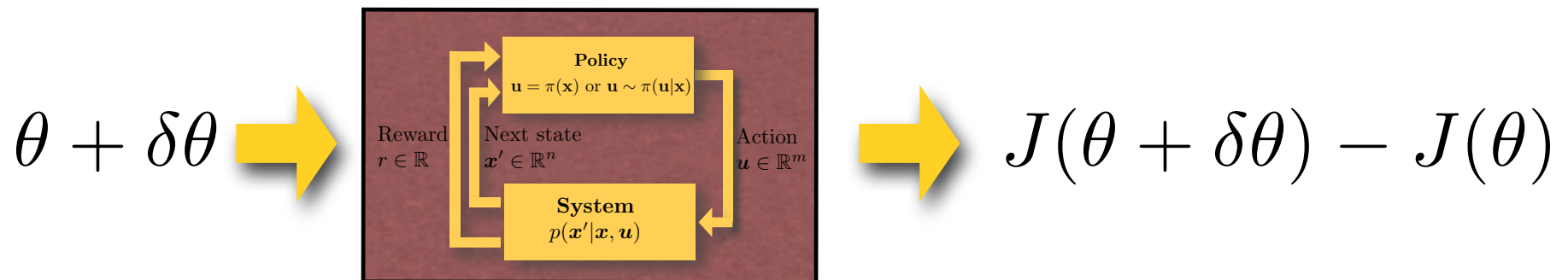
- Smooth changes in the parameters result into stability.
- Prior information can be incorporated with ease.
- Works with incomplete information.
- Exploration-Exploitation Dilemma implicitly treated.
- Is unbiased!
- Only requires much fewer samples.

3. FD vs LR gradients



Blackbox-Approach

Perturb the Parameters of your Policy



$$\frac{dJ}{d\theta} \approx \frac{J(\theta + \delta\theta) - J(\theta)}{\delta\theta}$$

Finite Difference Gradients



Why use Finite Difference Gradients?

- Only needs a black box!
- Works on any parameterization and deterministic policy.
- Fast to estimate for deterministic systems.
- State of the art in the simulation community

Why not?

- Parameter perturbation can destroy your robot.
- Exploration is hard to include.
- For stochastic systems it is very **slow**.

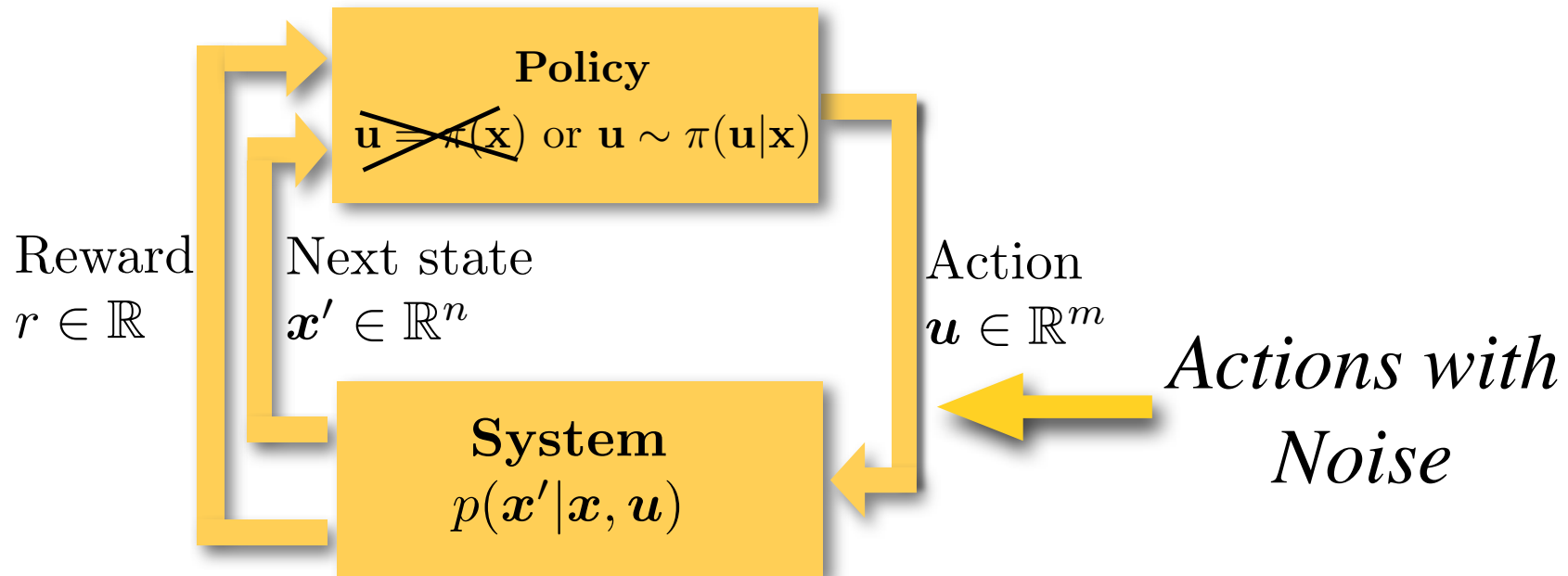
3. FD vs LR gradients

Likelihood Ratio Gradients



Whitebox-Approach

Perturb the actions using a stochastic policy



3. FD vs LR gradients



Whitebox-Approach: Likelihood Ratio Trick

$$\frac{d}{d\theta} J(\theta) = \frac{d}{d\theta} \int_{\mathcal{U}} \pi(u) r(u) du, \quad (1)$$

$$= \int_{\mathcal{U}} \frac{d\pi(u)}{d\theta} r(u) du, \quad (2)$$

$$= \int_{\mathcal{U}} \pi(u) \frac{1}{\pi(u)} \frac{d\pi(u)}{d\theta} r(u) du, \quad (3)$$

$$= \int_{\mathcal{U}} \pi(u) \frac{d \log \pi(u)}{d\theta} r(u) du, \quad (4)$$

$$= E \left\{ \frac{d \log \pi(u)}{d\theta} r(u) \right\} \approx \sum_{i=1}^N \frac{d \log \pi(u_i)}{d\theta} r(u_i) \quad (5)$$

Likelihood Ratio Gradients



Why use Likelihood Ratio Gradients?

- Fastest Gradient Method!
- We know the policy derivative - thus they are more efficient than for the simulation community.
- Only perturb the motor command -> policies will remain stable!

Why not?

- Stochastic policy.
- Injection of noise into the system.
- Theory much more complex!

3. FD vs LR gradients

Outline



1. Reinforcement Learning for Motor Control?
2. Finite Difference vs Likelihood-Ratio Policy Gradients
- ▶ **3. Likelihood-Ratio Policy Gradients**
4. Conclusion

Goal of RL Revisited



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

$$= (1 - \gamma) E \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \right\}$$

4. Likelihood Ratio Gradients

Policy Gradient Methods



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, very slow convergence, dependence on baseline!

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

4. Likelihood Ratio Gradients

Policy Gradient Methods



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, very slow convergence, dependence on baseline!

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

4. Likelihood Ratio Gradients

Compatible Function Approximation



The state-action value function can be replaced by

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

State-action
value function

Compatible function
approximation

Log-policy
derivative

Parameters
of the
function
approximator

without biasing the gradient.

Thus, the policy gradient becomes

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

(Sutton et al., 2000; Konda & Tsitsiklis, 2000)

4. Likelihood Ratio Gradients

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

(Peters et al., 2003)

4. Likelihood Ratio Gradients

Natural Gradients



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

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

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

4. Likelihood Ratio Gradients

All Action = 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

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

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

4. Likelihood Ratio Gradients

Natural Gradient Updates



As $\mathbf{G} = \mathbf{F}$, the gradient simplifies to

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

and the policy parameter update becomes

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

Important: The estimation of the gradient has simplified upon estimating the compatible function approximation / critic!!!

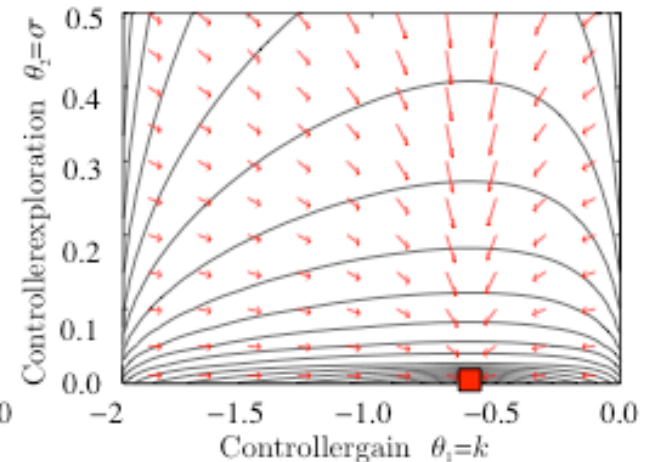
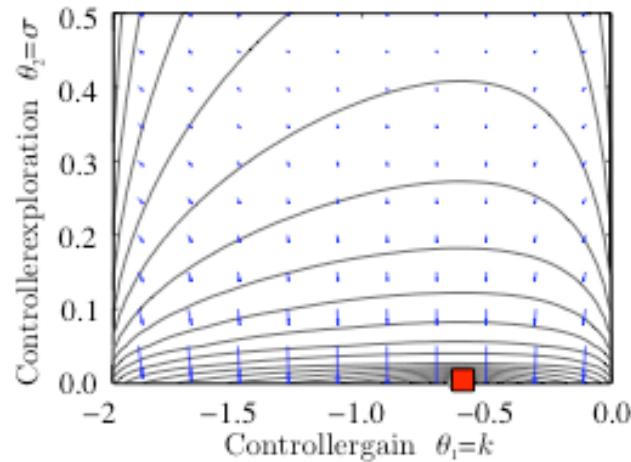
(Kakade, 2002; Peters et al., 2003, 2005; Bagnell&Schneider, 2003)

4. Likelihood Ratio Gradients

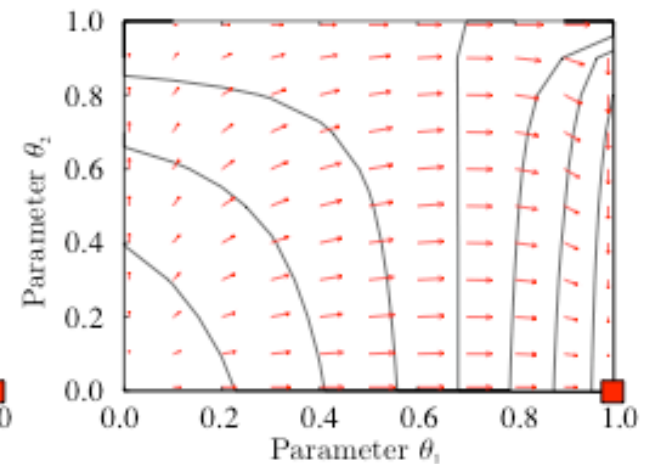
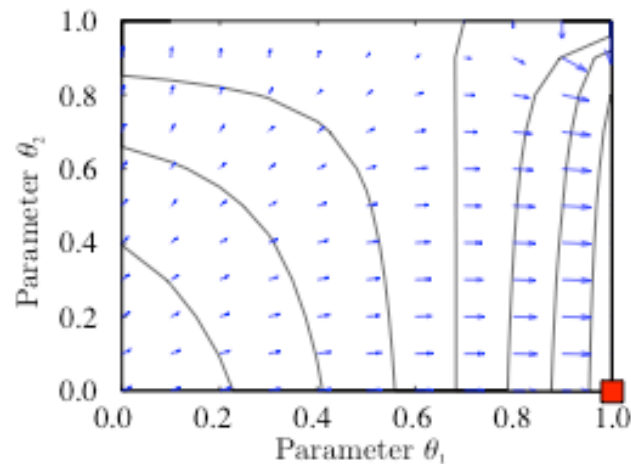
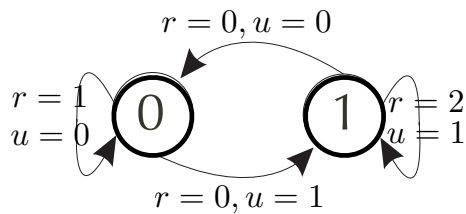
Natural Policy Gradients



Linear
Quadratic
Regulation



Two-State
Problem



4. Likelihood Ratio Gradients

Compatible Function Approximation



To obtain the natural gradient

$$\tilde{\nabla}_{\theta} J(\theta) = w$$

we need to estimate the compatible function approximation

$$f_w^{\pi}(x, u) = \frac{d \log \pi(u|x)}{d\theta}^T w$$

This function approximation is mean zero! Therefore it can ONLY represent the Advantage Function

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

4. Likelihood Ratio Gradients

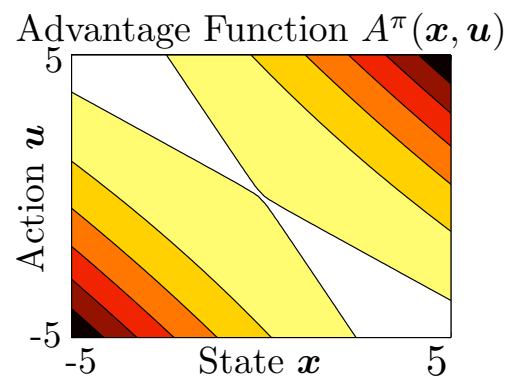
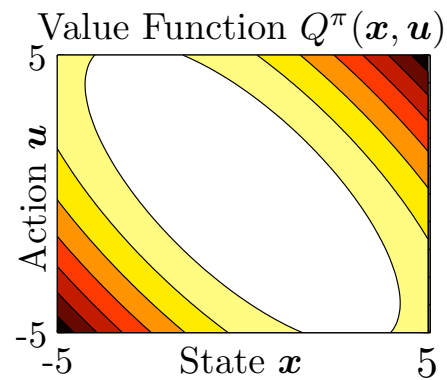
Compatible Function Approximation



The advantage function

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

is very different from the value functions!



...and we cannot directly do Temporal Difference Learning on this representation!

4. Likelihood Ratio Gradients

Natural Actor-Critic



We cannot do TD learning with

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

But when we add further basis function approximators

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

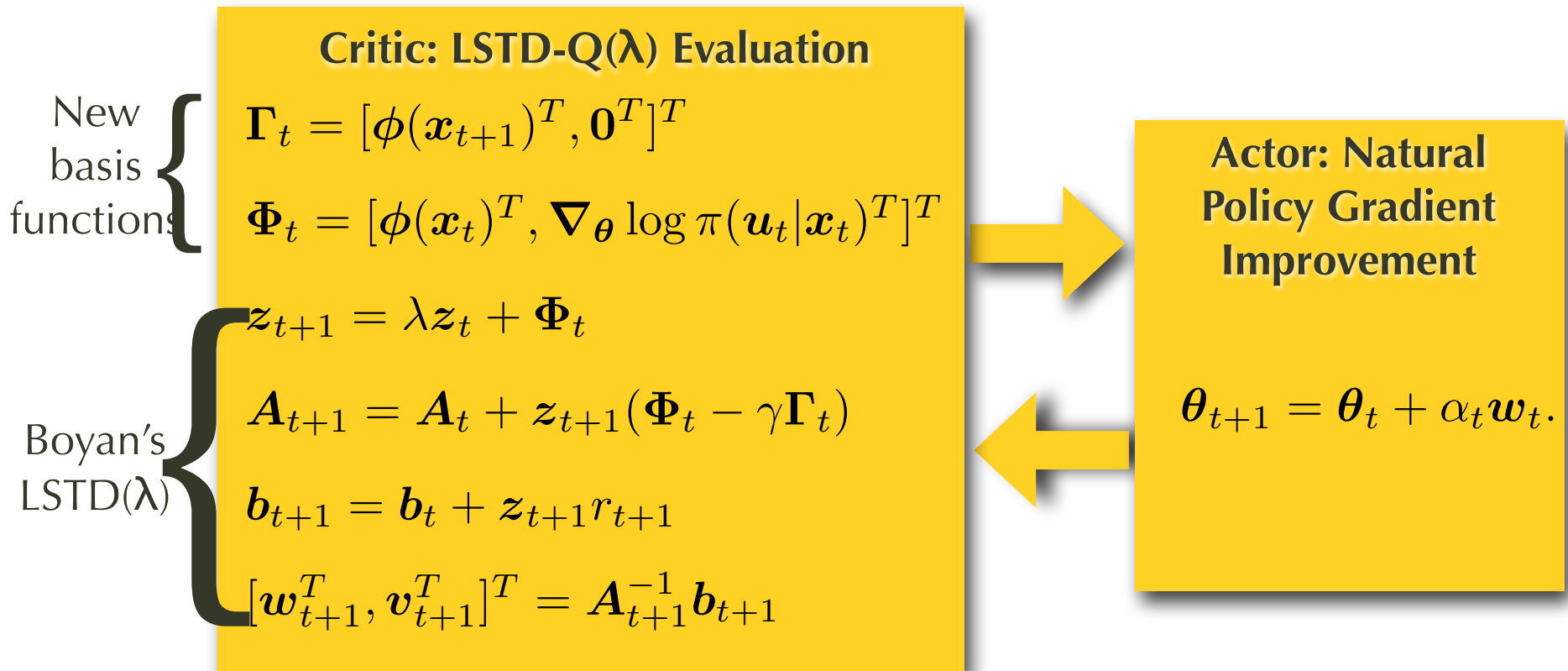
into the Bellman equation

$$V^\pi(x_t) + \nabla_{\theta} \log \pi(u_t|x_t)^T w = r(x_t, u_t) + \gamma V^\pi(x_{t+1}) + \epsilon_t$$

we get a **linear regression problem which can be solved with the LSTD(λ) algorithm (Boyan, 1996) in one step!**

4. Likelihood Ratio Gradients

Natural Actor-Critic



4. Likelihood Ratio Gradients

Algorithms Derivable from this Framework



Gibbs Policy

$$\pi(u_t|x_t) = e^{\theta_{xu}} / \sum_b e^{\theta_{xb}}$$

Additional Basis Functions

$$\phi(x) = [0, \dots, 0, 1, 0, \dots, 0]^T$$



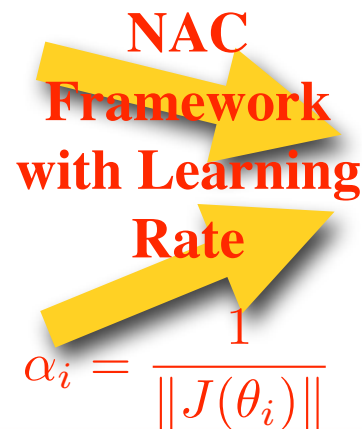
Sutton et al.'s (1983)
Actor Critic

Linear Gauss-Policy

$$\pi(u|x) = \mathcal{N}(u - \theta_{\text{gain}}^T x, \theta_{\text{explore}})$$

Additional Basis Functions

$$\phi(x) = x^T P x + p$$



Bradtke&Bartos (1993)
Q-Learning for LQR

4. Likelihood Ratio Gradients

Episodic Natural Actor Critic



...but in many cases we don't have any good additional basis functions!

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

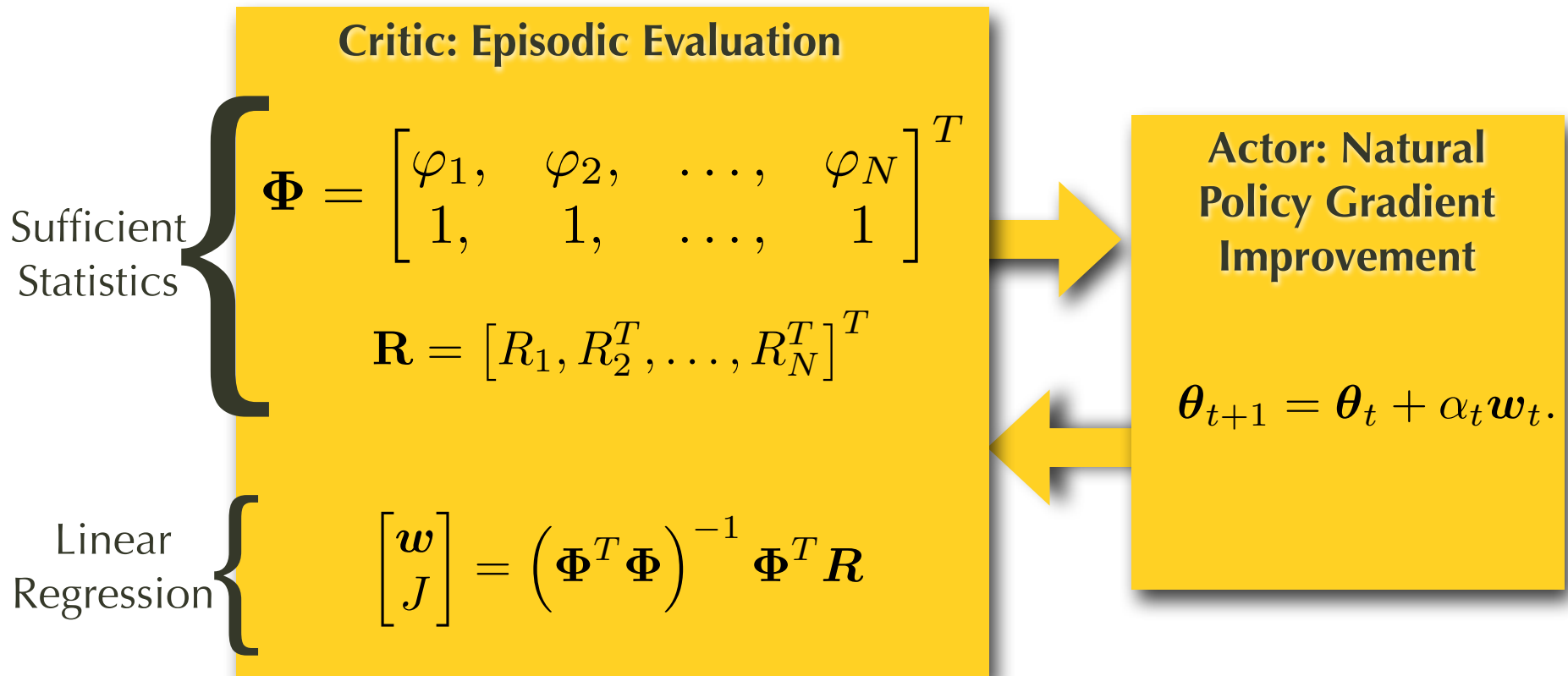
In this case, we can sum up the advantages along a trajectory and obtain one data point for a linear regression problem

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

...and an additional basis function of 1 suffices!

4. Likelihood Ratio Gradients

Episodic Natural Actor-Critic



4. Likelihood Ratio Gradients

Improving Motor Primitives



Ijspeert et al. (2002) suggested a nonlinear dynamics approach for motor primitives in imitation learning:

$$\text{Trajectory Plan Dynamics} \quad \left\{ \begin{array}{l} \dot{z} = \alpha_z (\beta_z (g - y) - z) \\ \dot{y} = \alpha_y (f(x, v) + z) \end{array} \right.$$

where

$$\text{Canonical Dynamics} \quad \left\{ \begin{array}{l} \dot{v} = \alpha_v (\beta_v (g - x) - v) \\ \dot{x} = \alpha_x v \end{array} \right.$$

Local Linear Model Approx.

$$f(x, v) = \frac{\sum_{i=1}^k w_i b_i v}{\sum_{i=1}^k w_i}$$

$$w_i = \exp\left(-\frac{1}{2} d_i (\bar{x} - c_i)^2\right) \quad \text{and} \quad \bar{x} = \frac{x - x_0}{g - x_0}$$

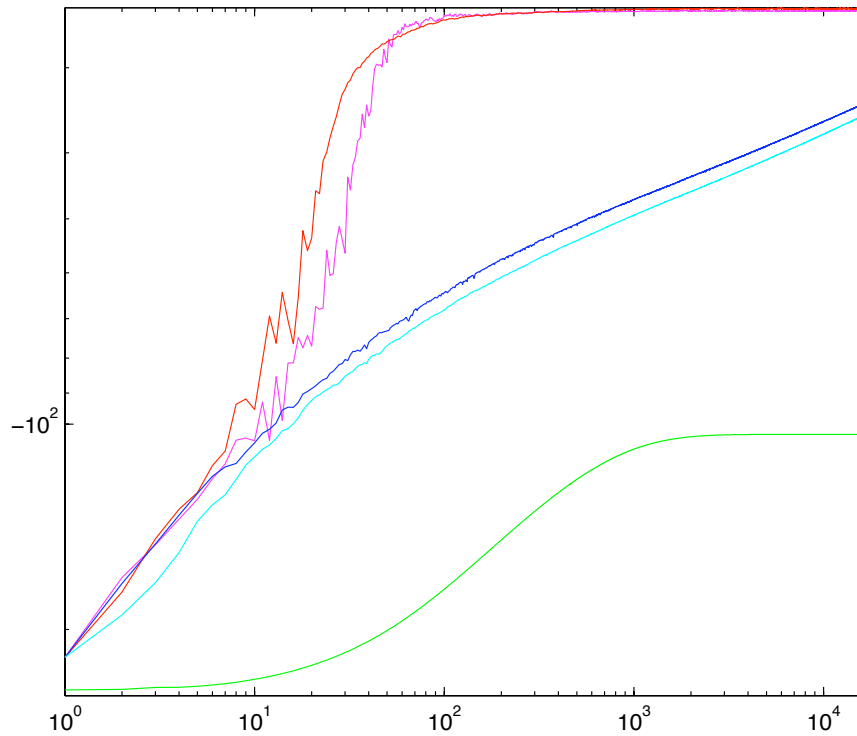
The parameters b can also be improved by Reinforcement Learning

4. Evaluations

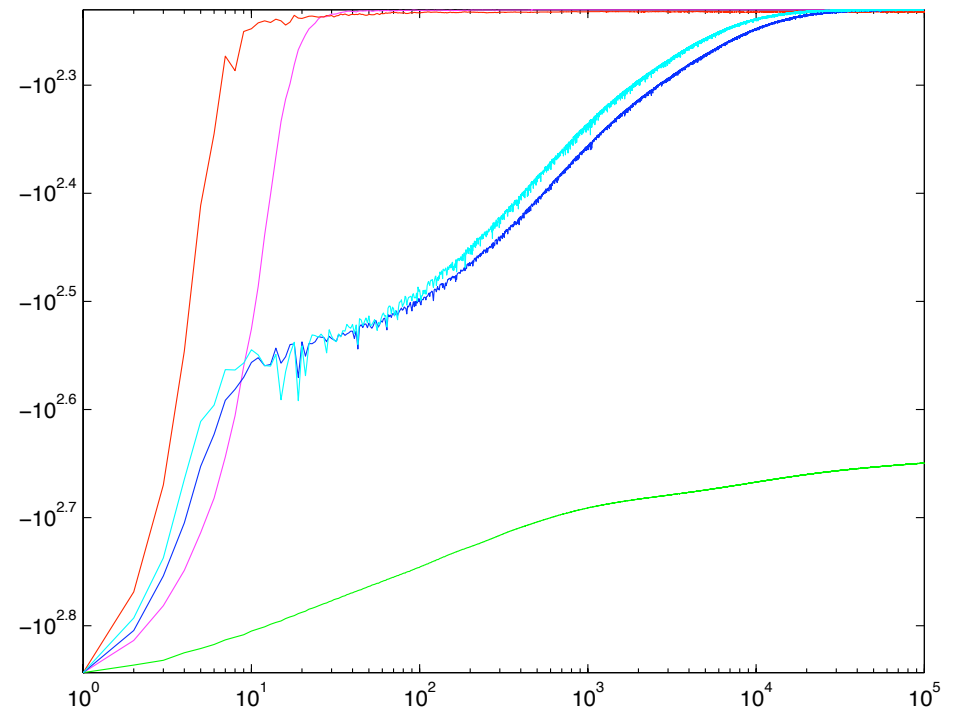
Improving Motor Primitives



Minimum Motor Command



Two Goals Policies



4. Evaluations

Outline



1. Reinforcement Learning for Motor Control?
2. Finite Difference vs Likelihood-Ratio Policy Gradients
3. Likelihood-Ratio Policy Gradients
- ▶ **4. Conclusion**

Conclusions



- *If you can explore your complete state-action space sufficiently ... use value function methods.*
- *If you have access to policy and its derivatives... use likelihood ratio policy gradient methods.*
 - *If you have good additional basis function... use Natural Actor-Critic.*
 - *If not ... use Episodic Natural Actor-Critic.*
 - *If you want to explore a problem fastly ... use 'vanilla' likelihood ratio policy gradient methods.*
- *If you can only access your parameters... use finite difference policy gradient methods.*
- *If you can only access your parameters... use finite difference policy gradient methods.*

6. Conclusion

Projects...



- **Learning Motor Primitives with Reward-Weighted Regression**
 - *This will be a nearly new method... very enthusiastic people needed!*
- ***Applying Policy Gradients (methods of your choice!) to Oscillator Optimization!***
 - *Here we can create several projects if you like :)*

6. Conclusion