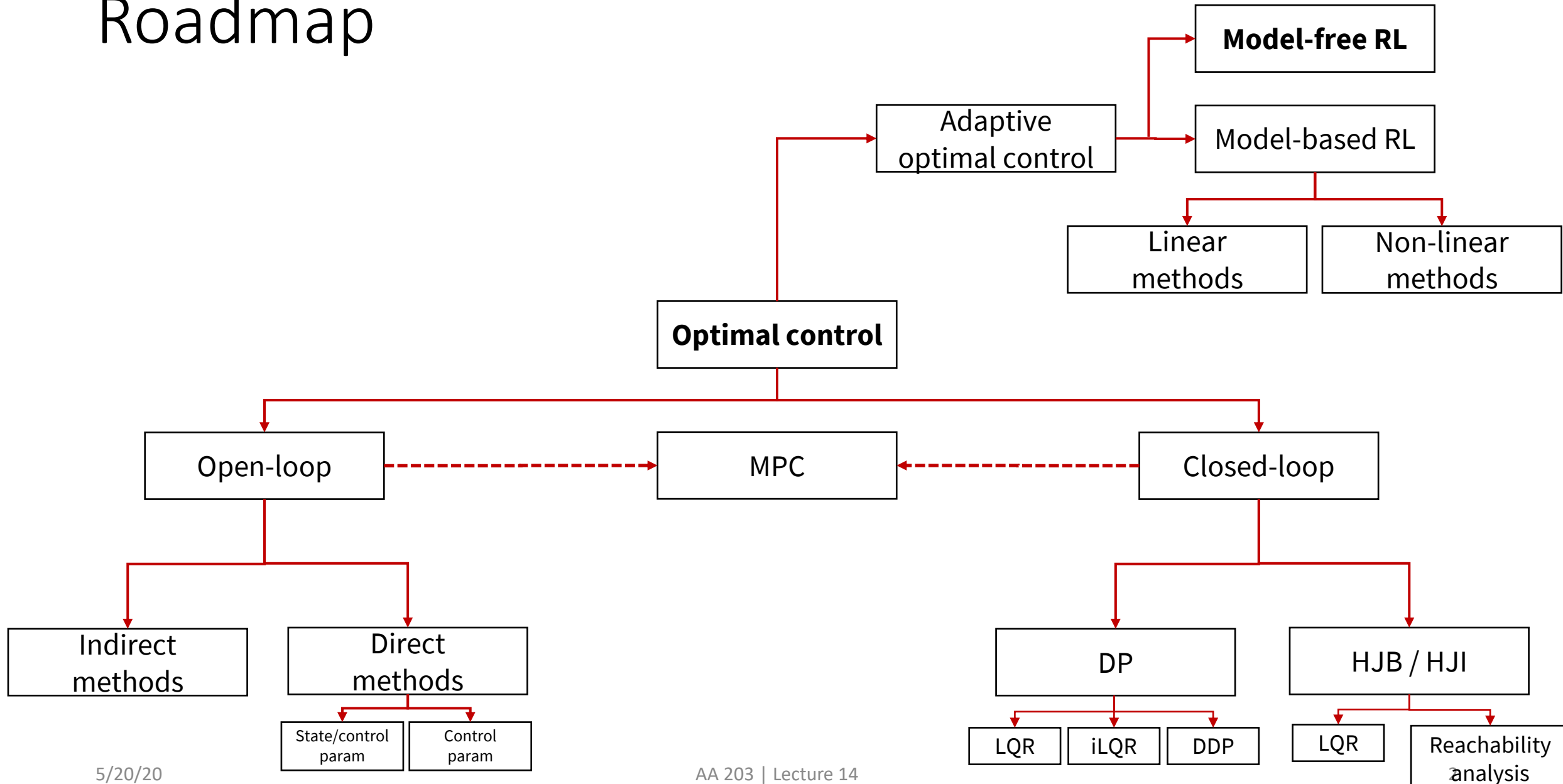


# AA203

# Optimal and Learning-based Control

Actor-critic and advanced policy gradient

# Roadmap



# Agenda

- Review model-free RL
- Introduce variance reduction methods for policy gradient estimation
- Brief survey of the modern model-free RL landscape
- Readings:
  - R. Sutton and A. Barto. *Reinforcement Learning: An Introduction*, 2018.

# Review: Policy Gradient

For policy  $\pi_\theta$ , define

$$J(\theta) = \mathbb{E}_\tau[r(\tau)]$$

for  $\tau = (x_0, u_0, x_1, \dots, x_N)$ ,  $r(\tau) = \sum_{t \geq 0} r(x_t, u_t)$ .

Using  $\nabla_\theta \pi_\theta(u_t | x_t) = \pi_\theta(u_t | x_t) \nabla_\theta \log \pi_\theta(u_t | x_t)$ , have

$$\nabla_\theta J(\theta) \approx \sum_i r(\tau_i) \nabla_\theta \log \pi_\theta(u_{t,i} | x_{t,i})$$

# Review: Q Learning

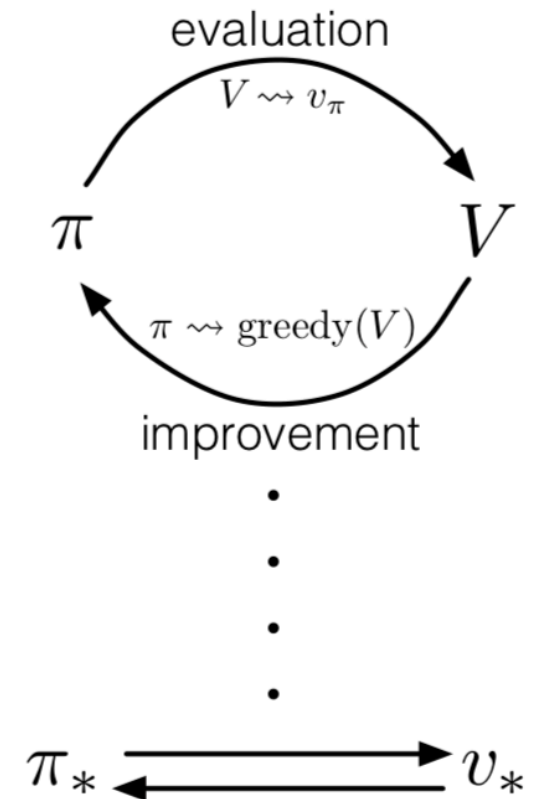
Instead of directly optimize policy, interleave a policy evaluation step,

$$\min_{\theta} \left( r_t + \gamma \max_u Q_{\theta'}(x_{t+1}, u) - Q_{\theta}(x_t, u_t) \right)^2$$

with a greedy policy improvement step,  $\pi(x) = \max_a Q_{\theta}(x, u)$ .

# Generalized policy iteration

- Recall policy iteration from discrete MDP discussion:
  - First, run policy evaluation (dynamic programming) to convergence to get  $V^\pi$
  - Then, improve policy via policy improvement step
- Q learning
  - Policy evaluation via TD error minimization
  - Improvement via maximizing Q function
- Policy gradient
  - Policy evaluation via Monte Carlo rollouts
  - Improvement via gradient step



Sutton and Barto, 2018.

# Time dependency of policy gradient theorem

- Previous estimate of infinite horizon return was

$$J(\theta) = \mathbb{E}_{x_0, u_0, x_1, \dots} \left[ \sum_{t \geq 0} r(x_t, u_t) \right]$$

Equally valid:

$$J(\theta) = \mathbb{E}_{x_\tau, u_\tau, \dots} \left[ \sum_{t \geq \tau} r(x_t, u_t) \right]$$

for  $x_\tau \sim p(\cdot | x_0, \theta)$ .

# REINFORCE

Loop forever:

Generate episode  $x_0, u_0, r_0, x_1, u_1, r_1 \dots$  with  $\pi_\theta$

Loop for all  $t = 0, \dots, N - 1$ :

$$G \leftarrow \sum_{k=t+1}^N r_k$$

$$\theta \leftarrow \theta + \alpha G \nabla_\theta \log \pi_\theta(u_t | x_t)$$



# Adding baselines to policy evaluation

- Monte Carlo policy gradient estimator has **extremely high variance**.
- We want to search for gradient estimators that have lower variance
- Add in **baseline**

$$J(\theta) = \mathbb{E}_{x_\tau, u_\tau, \dots} \left[ \sum_{t \geq \tau} r(x_t, u_t) - b(x_\tau) \right]$$

Use policy gradient theorem to get

$$\nabla_\theta J(\theta) = \mathbb{E}_{x_\tau, u_\tau, \dots} \left[ \left( \sum_{t \geq \tau} r(x_t, u_t) - b(x_\tau) \right) \nabla_\theta \log \pi(u_t | x_t, \theta) \right]$$

# A closer look at the baseline

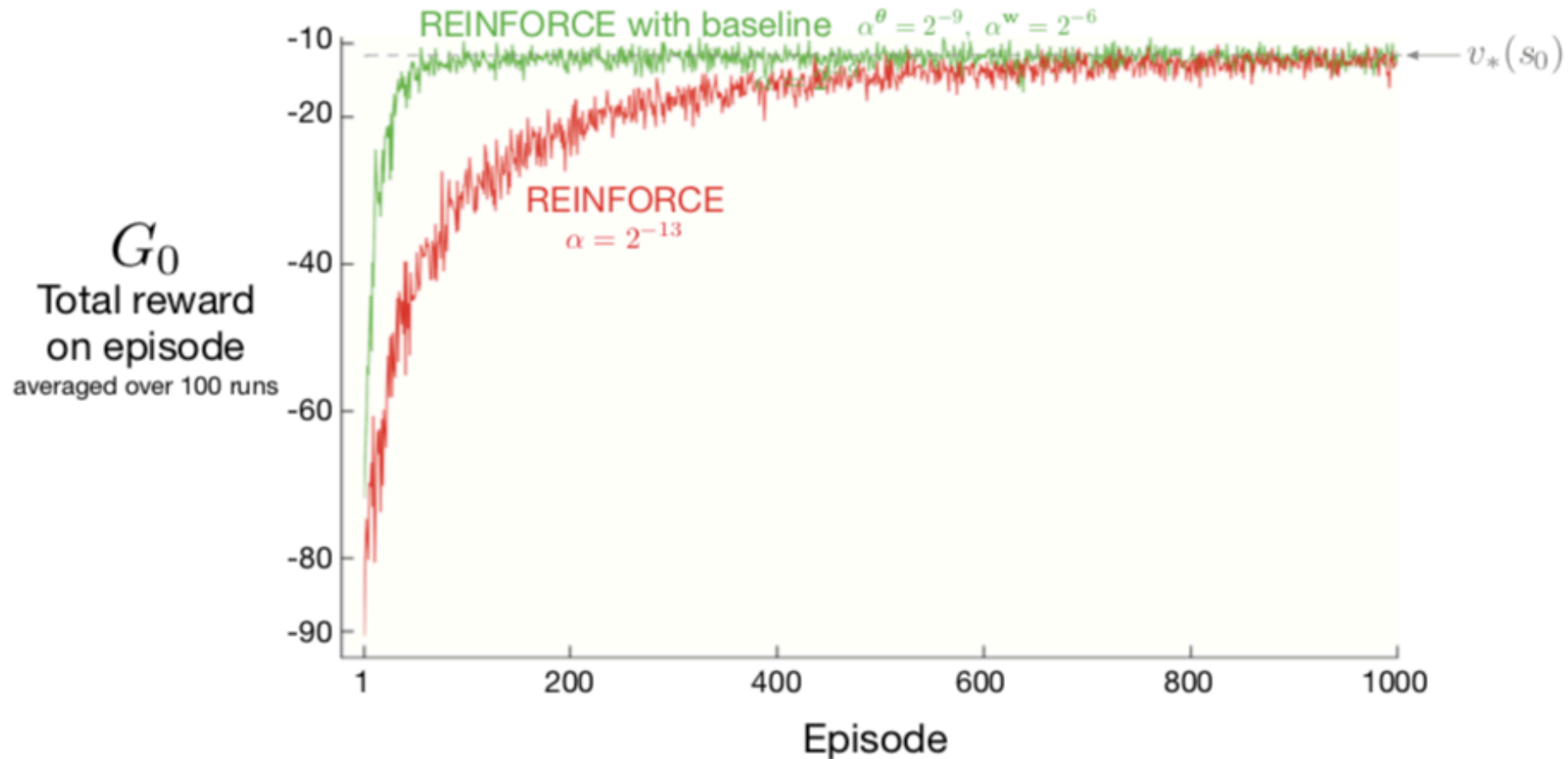
Claim: adding baseline does not change the value of the expected gradient

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{x_{\tau}, u_{\tau}, \dots} \left[ \left( \sum_{t \geq \tau} r(x_t, u_t) - b(x_{\tau}) \right) \nabla_{\theta} \log \pi(u_t | x_t, \theta) \right] \\ &= \mathbb{E} \left[ \sum_{t \geq \tau} r(x_t, u_t) \nabla_{\theta} \log \pi(u_t | x_t, \theta) \right] - \mathbb{E} [b(x_{\tau}) \nabla_{\theta} \log \pi(u_t | x_t, \theta)] \\ \mathbb{E} [b(x_{\tau}) \nabla_{\theta} \log \pi(u_t | x_t, \theta)] &= \mathbb{E}_{x_{\tau}} [b(x_{\tau}) \mathbb{E}_{u_{\tau}} [\nabla_{\theta} \log \pi(u_t | x_t, \theta)]] \\ \mathbb{E}_{u_{\tau}} [\nabla_{\theta} \log \pi(u_t | x_t, \theta)] &= \nabla_{\theta} \mathbb{E}_{u_{\tau}} [1] = 0\end{aligned}$$

Any state-dependent function, indep. of action, works.

# Example

# Performance improvement on gridworld



# Actor-critic

Particularly good choice: value function

Actor-critic: use both **actor** (policy) and **critic** (value function).

Loop forever:

Generate episode  $x_0, u_0, r_0, x_1, u_1, r_1 \dots$  with  $\pi_\theta$

Loop for all  $t = 0, \dots, N - 1$ :

$$G \leftarrow \sum_{k=t+1}^N r_k$$

$$\delta_w \leftarrow G - V_w(x_t)$$

$$w \leftarrow w + \alpha_w \delta_w \nabla_w V_w(x_t)$$

$$\theta \leftarrow \theta + \alpha_\theta \delta_w \nabla_\theta \log \pi_\theta(u_t | x_t)$$

# Policy gradient theorem with Q function

- Previously, have used  $J(\theta) = \mathbb{E}_{x_\tau, u_\tau, \dots} [\sum_{t \geq \tau} r(x_t, u_t)]$
- Note that

$$J(\theta) = \mathbb{E}_{u_\tau \sim \pi(\cdot | x_\tau)} [Q^\pi(x_\tau, u_\tau)]$$

Yields policy gradient

$$\nabla_\theta J(\theta) = \mathbb{E}_{u_\tau \sim \pi(\cdot | x_\tau)} [Q^\pi(x_\tau, u_\tau) \nabla_\theta \log \pi(u_\tau | x_\tau)]$$

Note that  $Q^\pi(x_\tau, u_\tau) = \mathbb{E}_{u_\tau \sim \pi(\cdot | x_\tau), x_{\tau+1}} [r(x_\tau, u_\tau) + V^\pi(x_{\tau+1})]$

# Advantage policy gradient

- Combining the Q function policy gradient and the value baselines, we have

$$\nabla_{\theta} J(\theta) = \mathbb{E}[\delta^{\pi} \nabla_{\theta} \log \pi(u_{\tau} | x_{\tau})]$$

For  $\delta^{\pi} = (r_t + V^{\pi}(x_{t+1}) - V^{\pi}(x_t))$ . This is the TD error for policy evaluation!

- Note that  $\mathbb{E}_{\pi}[\delta^{\pi} | x, u] = Q^{\pi}(x, u) - V^{\pi}(x) = A^{\pi}(x, u)$ .
  - This is called the **advantage**.

# Advantage actor-critic

Loop forever:

Generate episode  $x_0, u_0, r_0, x_1, u_1, r_1 \dots$  with  $\pi_\theta$

Loop for all  $t = 0, \dots, N - 1$ :

$$\delta_w \leftarrow r_t + V_w(x_{t+1}) - V_w(x_t)$$

$$w \leftarrow w + \alpha_w \delta_w \nabla_w V_w(x_t)$$

$$\theta \leftarrow \theta + \alpha_\theta \delta_w \nabla_\theta \log \pi_\theta(u_t | x_t)$$



# Alternative estimators

- Many possible estimators for the advantage

- Multistep TD error:

$$\delta \leftarrow r_t + r_{t+1} + \dots + r_{t+\tau} + V_w(x_{t+\tau+1}) - V_w(x_t)$$

As  $\tau$  gets larger, this gets closer to Monte Carlo with value baseline.

# Trust region policy optimization (TRPO)

[Schulman et al., ICML 2015]

- Main idea : instead of choosing step size, use trust region

$$\begin{aligned} & \max E\left[\frac{\pi_{\theta}(u_t|x_t)}{\pi_{\theta_{old}}(u_t|x_t)} \hat{A}_t\right] \\ & s. t. E_{x \sim \rho_{old}} \left[ D_{KL} \left( \pi_{\theta_{old}}(\cdot | x) || \pi_{\theta}(\cdot | x) \right) \right] \leq \delta \end{aligned}$$

- Can show that this leads to monotonic improvement in the ideal case.
- Simpler, more popular version: proximal policy optimization (PPO).
  - Replaces TRPO CG solve with simple adaptive KL penalty.

# Deterministic policy gradient (DPG)

[Silver et al., ICML 2014]

- Instead of using stochastic policy with value estimation baseline:
  - Maintain estimate of Q function via minimizing TD error
  - Optimize deterministic policy via

$$\max_{\theta} E_x[Q(s, \pi_{\theta}(x))]$$

- Policy simply *amortizes* optimization of the Q function.
- Can be used off policy, relatively unstable in practice.

# Maximization bias

- Even though state-action value estimates are unbiased, may still have biased value estimates
- Example [Mannor et al., 2007]:

# Double Q-learning

- Several possible solutions; in general, want to avoid using *max of estimates as estimate of max*.
- Double Q-learning [van Hasselt, NeurIPS 2010]: use two independent estimates  $Q_1, Q_2$ 
  - $u^* = \operatorname{argmax}_u Q_1(x, u)$
  - Use value estimate  $Q_2(x, u^*)$
- Alternative approach: maintain two independent critics, always use min [Fujimoto et al, ICML 2018]

# Criticism of model-free methods

- Despite recent progress (including much not discussed here), questions about whether model-free methods are doing more than random search in parameter space.

---

## Why did TD-Gammon Work?

---

**Jordan B. Pollack & Alan D. Blair**  
Computer Science Department  
Brandeis University  
Waltham, MA 02254  
{pollack,blair}@cs.brandeis.edu

### Abstract

Although TD-Gammon is one of the major successes in machine learning, it has not led to similar impressive breakthroughs in temporal difference learning for other applications or even other games. We were able to replicate some of the success of TD-Gammon, developing a competitive evaluation function on a 4000 parameter feed-forward neural network, without using back-propagation, reinforcement or temporal difference learning methods. Instead we apply simple hill-climbing in a relative fitness environment. These results and further analysis suggest that the surprising success of Tesauro's program had more to do with the co-evolutionary structure of the learning task and the dynamics of the backgammon game itself.

---

## Simple random search of static linear policies is competitive for reinforcement learning

---

**Horia Mania**      **Aurelia Guy**      **Benjamin Recht**  
hmania@berkeley.edu      lia@berkeley.edu      brecht@berkeley.edu

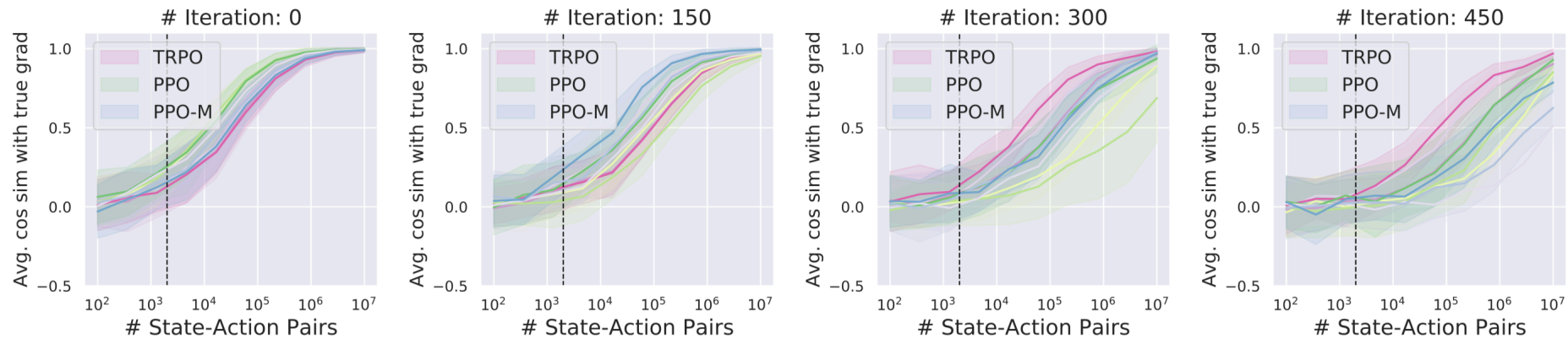
Department of Electrical Engineering and Computer Science  
University of California, Berkeley

### Abstract

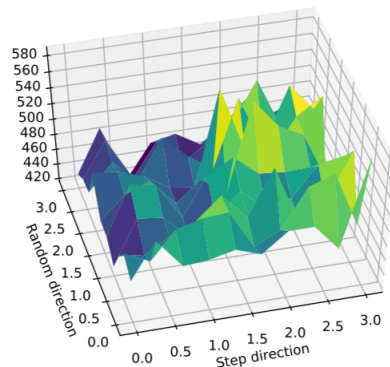
Model-free reinforcement learning aims to offer off-the-shelf solutions for controlling dynamical systems without requiring models of the system dynamics. We introduce a model-free random search algorithm for training static, linear policies for continuous control problems. Common evaluation methodology shows that our method matches state-of-the-art sample efficiency on the benchmark MuJoCo locomotion tasks. Nonetheless, more rigorous evaluation reveals that the assessment of performance on these benchmarks is optimistic. We evaluate the performance of our method over hundreds of random seeds and many different hyperparameter configurations for each benchmark task. This extensive evaluation is possible because of the small computational footprint of our method. Our simulations highlight a high variability in performance in these benchmark tasks, indicating that commonly used estimations of sample efficiency do not adequately evaluate the performance of RL algorithms. Our results stress the need for new baselines, benchmarks and evaluation methodology for RL algorithms.

# Are Deep Policy Gradient Algorithms Truly Policy Gradient Algorithms?

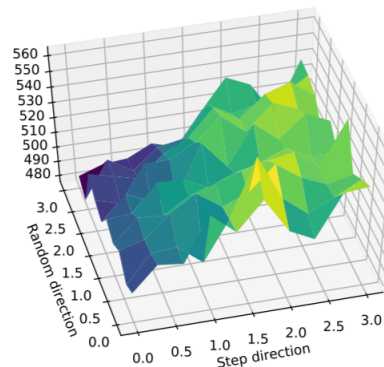
Andrew Ilyas<sup>\*1</sup>, Logan Engstrom<sup>\*1</sup>, Shibani Santurkar<sup>1</sup>, Dimitris Tsipras<sup>1</sup>,  
Firdaus Janoos<sup>2</sup>, Larry Rudolph<sup>1,2</sup>, and Aleksander Mądry<sup>1</sup>



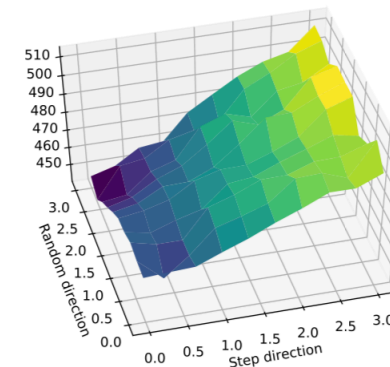
2,000 state-action pairs  
(19 trajectories)



20,000 state-action pairs  
(198 trajectories)



100,000 state-action pairs  
(1068 trajectories)



# Why model-free?

- Advantages
  - Very few assumptions
  - Many state of the art methods reach better performance than model-based methods
- Weaknesses
  - Extremely high sample complexity



# Next time

- Combining policy optimization with model learning