

# Lecture 6: Policy Gradient II. Advanced policy gradient section slides from Joshua Achiam's slides, with minor modifications

Emma Brunskill

CS234 Reinforcement Learning.

Winter 2026

- Select all that are true about policy gradients:

- 1  $\nabla_{\theta} V(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q^{\pi_{\theta}}(s, a)]$
- 2  $\theta$  is always increased in the direction of  $\nabla_{\theta} \ln(\pi(S_t, A_t, \theta))$ .
- 3 State-action pairs with higher estimated  $Q$  values will increase in probability on average
- 4 Are guaranteed to converge to the global optima of the policy class
- 5 Not sure

- Select all that are true about policy gradients:
  - 1  $\nabla_{\theta} V(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q^{\pi_{\theta}}(s, a)]$
  - 2  $\theta$  is always increased in the direction of  $\nabla_{\theta} \ln(\pi(S_t, A_t, \theta))$ .
  - 3 State-action pairs with higher estimated Q values will increase in probability on average
  - 4 Are guaranteed to converge to the global optima of the policy class
  - 5 Not sure

- Last time: Policy Search
- This time: Policy search continued.

- Likelihood ratio / score function policy gradient
  - Baseline
  - Alternative targets
- Advanced policy gradient methods
  - Proximal policy optimization (PPO) (will implement in homework)

- Last time: time: DQN and REINFORCE
- This time: Policy Gradient and PPO
- Next time: Policy Search Cont.

## Policy Gradient Algorithms and Reducing Variance

- 1 Policy Gradient Algorithms and Reducing Variance
  - Baseline
  - Alternatives to MC Returns



- Leverages likelihood ratio / score function and temporal structure

$$\Delta\theta_t = \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) G_t$$

## **REINFORCE:**

Initialize policy parameters  $\theta$  arbitrarily

**for** each episode  $\{s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_{\theta}$  **do**

**for**  $t = 1$  to  $T - 1$  **do**

$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) G_t$

**endfor**

**endfor**

**return**  $\theta$

$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m R(\tau^{(i)}) \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{s}_t^{(i)})$$

- Unbiased but very noisy
- Fixes that can make it practical
  - Temporal structure
  - **Baseline**
  - Alternatives to using Monte Carlo returns  $R(\tau^{(i)})$  as targets

- Reduce variance by introducing a *baseline*  $b(s)$

$$\nabla_{\theta} \mathbb{E}_{\tau}[R] = \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t; \theta) \left( \sum_{t'=t}^{T-1} r_{t'} - b(s_t) \right) \right]$$

- For any choice of  $b$ , gradient estimator is unbiased.
- Near optimal choice is the expected return,

$$b(s_t) \approx \mathbb{E}[r_t + r_{t+1} + \dots + r_{T-1}]$$

- Interpretation: increase logprob of action  $a_t$  proportionally to how much returns  $\sum_{t'=t}^{T-1} r_{t'}$  are better than expected

## Baseline $b(s)$ Does Not Introduce Bias–Derivation

$$\begin{aligned} & \mathbb{E}_{\tau} [\nabla_{\theta} \log \pi(\mathbf{a}_t | \mathbf{s}_t; \theta) b(\mathbf{s}_t)] \\ &= \mathbb{E}_{\mathbf{s}_{0:t}, \mathbf{a}_{0:(t-1)}} \left[ \mathbb{E}_{\mathbf{s}_{(t+1):T}, \mathbf{a}_{t:(T-1)}} [\nabla_{\theta} \log \pi(\mathbf{a}_t | \mathbf{s}_t; \theta) b(\mathbf{s}_t)] \right] \end{aligned}$$

# Baseline $b(s)$ Does Not Introduce Bias–Derivation

$$\begin{aligned} & \mathbb{E}_\tau [\nabla_\theta \log \pi(a_t | s_t; \theta) b(s_t)] \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[ \mathbb{E}_{s_{(t+1):T}, a_{t:(T-1)}} [\nabla_\theta \log \pi(a_t | s_t; \theta) b(s_t)] \right] \text{ (break up expectation)} \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[ b(s_t) \mathbb{E}_{s_{(t+1):T}, a_{t:(T-1)}} [\nabla_\theta \log \pi(a_t | s_t; \theta)] \right] \text{ (pull baseline term out)} \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} [b(s_t) \mathbb{E}_{a_t} [\nabla_\theta \log \pi(a_t | s_t; \theta)]] \text{ (remove irrelevant variables)} \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[ b(s_t) \sum_a \pi_\theta(a_t | s_t) \frac{\nabla_\theta \pi(a_t | s_t; \theta)}{\pi_\theta(a_t | s_t)} \right] \text{ (likelihood ratio)} \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[ b(s_t) \sum_a \nabla_\theta \pi(a_t | s_t; \theta) \right] \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[ b(s_t) \nabla_\theta \sum_a \pi(a_t | s_t; \theta) \right] \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} [b(s_t) \nabla_\theta 1] \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} [b(s_t) \cdot 0] = 0 \end{aligned}$$

## Argument for Why Baseline $b(s)$ Can Reduce Variance

- Motivation was for introducing baseline  $b(s)$  was to reduce variance

$$\text{Var}[\nabla_{\theta} \mathbb{E}_{\tau}[R]] = \text{Var} \left[ \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t; \theta) (R_t(s_t) - b(s_t)) \right] \right] \quad (1)$$

# Argument for Why Baseline $b(s)$ Can Reduce Variance

- Motivation was for introducing baseline  $b(s)$  was to reduce variance

$$\text{Var}[\nabla_{\theta} \mathbb{E}_{\tau}[R]] = \text{Var} \left[ \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t; \theta) (R_t(s_t) - b(s_t)) \right] \right] \quad (2)$$

$$\approx \sum_{t=0}^{T-1} \mathbb{E}_{\tau} \text{Var} \left[ \nabla_{\theta} \log \pi(a_t | s_t; \theta) (R_t(s_t) - b(s_t)) \right] \quad (3)$$

- Focus on the variance of one term.

$$\begin{aligned} \text{Var} \left[ \nabla_{\theta} \log \pi(a_t | s_t; \theta) (R_t(s_t) - b(s_t)) \right] &= E \left[ \left[ \nabla_{\theta} \log \pi(a_t | s_t; \theta) (R_t(s_t) - b(s_t)) \right]^2 \right] \\ &\quad - \left[ E \left[ \nabla_{\theta} \log \pi(a_t | s_t; \theta) (R_t(s_t) - b(s_t)) \right] \right]^2 \end{aligned}$$

- Choosing a baseline to minimize variance
- Recall the baseline  $b(s)$  does not impact the expectation. Therefore sufficient to consider

$$\begin{aligned} \arg \max_b \text{Var} \left[ \nabla_{\theta} \log \pi(a_t | s_t; \theta) (G_t(s_t) - b(s_t)) \right] &= \arg \min_b E \left[ \left[ \nabla_{\theta} \log \pi(a_t | s_t; \theta) (G_t(s_t) - b(s_t)) \right]^2 \right] \quad (4) \\ &= \arg \min_b E_{s \sim d} \pi \left[ E_{a \sim \pi(\cdot | s), G | s, a} \left[ \nabla_{\theta} \log \pi(a_t | s; \theta) (G_t(s) - b(s))^2 \right] \right] \end{aligned}$$

- This is a weighted least squares problem. Taking the derivative and setting to zero yields

$$b(s) = \frac{E_{a \sim \pi(\cdot | s), G | s, a} \left[ \nabla_{\theta} \log \pi(a_t | s; \theta) (G_t(s))^2 \right]}{E_{a \sim \pi(\cdot | s), G | s, a} \left[ \nabla_{\theta} \log \pi(a_t | s; \theta) \right]^2} \approx E_{a \sim \pi(\cdot | s), G | s, a} [G_t(s)] \quad (5)$$

# "Vanilla" Policy Gradient Algorithm

Initialize policy parameter  $\theta$ , baseline  $b$

**for** iteration= $1, 2, \dots$  **do**

Collect a set of trajectories by executing the current policy

At each timestep  $t$  in each trajectory  $\tau^i$ , compute

Return  $G_t^i = \sum_{t'=t}^{T-1} r_{t'}^i$ , and

Advantage estimate  $\hat{A}_t^i = G_t^i - b(s_t^i)$ .

Re-fit the baseline, by minimizing  $\sum_i \sum_t |b(s_t^i) - G_t^i|^2$ ,

Update the policy, using a policy gradient estimate  $\hat{g}$ ,

Which is a sum of terms  $\nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_t$ .

(Plug  $\hat{g}$  into SGD or ADAM)

**endfor**



Initialize policy parameter  $\theta$ , baseline  $b$

**for** iteration=1, 2, ... **do**

Collect a set of trajectories by executing the current policy

At each timestep  $t$  in each trajectory  $\tau^i$ , compute

Return  $G_t^i = \sum_{t'=t}^{T-1} r_{t'}^i$ , and

Advantage estimate  $\hat{A}_t^i = G_t^i - b(s_t^i)$ .

Re-fit the baseline, by minimizing  $\sum_i \sum_t |b(s_t^i) - G_t^i|^2$ ,

Update the policy, using a policy gradient estimate  $\hat{g}$ ,

Which is a sum of terms  $\nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_t$ .

(Plug  $\hat{g}$  into SGD or ADAM)

**endfor**

- Recall Q-function / state-action-value function:

$$Q^\pi(s, a) = \mathbb{E}_\pi \left[ r_0 + \gamma r_1 + \gamma^2 r_2 \cdots \mid s_0 = s, a_0 = a \right]$$

- State-value function can serve as a great baseline

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi \left[ r_0 + \gamma r_1 + \gamma^2 r_2 \cdots \mid s_0 = s \right] \\ &= \mathbb{E}_{a \sim \pi} [Q^\pi(s, a)] \end{aligned}$$

## 1 Policy Gradient Algorithms and Reducing Variance

- Baseline
- Alternatives to MC Returns

- Policy gradient:

$$\nabla_{\theta} \mathbb{E}[R] \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t, s_t) (G_t^{(i)} - b(s_t))$$

- Fixes that improve simplest estimator
  - Temporal structure (shown in above equation)
  - Baseline (shown in above equation)
  - **Alternatives to using Monte Carlo returns  $G_t^i$  as estimate of expected discounted sum of returns for the policy parameterized by  $\theta$ ?**

- $G_t^i$  is an estimation of the value function at  $s_t$  from a single roll out
- Unbiased but high variance
- Reduce variance by introducing bias using bootstrapping and function approximation
  - Just like we saw for TD vs MC, and value function approximation

- Estimate of  $V/Q$  is done by a **critic**
- **Actor-critic** methods maintain an explicit representation of policy and the value function, and update both
- A3C (Mnih et al. ICML 2016) is a very popular actor-critic method

- Recall:

$$\nabla_{\theta} \mathbb{E}_{\tau} [R] = \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t; \theta) \left( \sum_{t'=t}^{T-1} r_{t'} - b(s_t) \right) \right]$$

$$\nabla_{\theta} \mathbb{E}_{\tau} [R] \approx \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t; \theta) (Q(s_t, a_t; \mathbf{w}) - b(s_t)) \right]$$

- Letting the baseline be an estimate of the value  $V$ , we can represent the gradient in terms of the state-action advantage function

$$\nabla_{\theta} \mathbb{E}_{\tau} [R] \approx \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t; \theta) \hat{A}^{\pi}(s_t, a_t) \right]$$

- where the advantage function  $A^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s)$

## Advanced Policy Gradients



## Theory:

- 1 Problems with Policy Gradient Methods
- 2 Policy Performance Bounds
- 3 Monotonic Improvement Theory (next time)

## Algorithms:

- 1 Proximal Policy Optimization

## The Problems with Policy Gradients

Policy gradient algorithms try to solve the optimization problem

$$\max_{\theta} J(\pi_{\theta}) \doteq \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$

by taking stochastic gradient ascent on the policy parameters  $\theta$ , using the *policy gradient*

$$g = \nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A^{\pi_{\theta}}(s_t, a_t) \right].$$

Limitations of policy gradients:

- Sample efficiency is poor
- Distance in parameter space  $\neq$  distance in policy space!
  - What is policy space? For tabular case, set of matrices

$$\Pi = \left\{ \pi : \pi \in \mathbb{R}^{|S| \times |A|}, \sum_a \pi_{sa} = 1, \pi_{sa} \geq 0 \right\}$$

- Policy gradients take steps in parameter space
- Step size is hard to get right as a result

- Sample efficiency for vanilla policy gradient methods is poor
- Discard each batch of data immediately after **just one gradient step**
- Why? PG is an **on-policy expectation**.
- Two main approaches to obtaining an unbiased estimate of the policy gradient
  - Collect sample trajectories from policy, then form sample estimate. (More stable)
  - Use trajectories from other policies (Less stable)
- Opportunity: use old data to take **multiple gradient steps** before using the resulting new policy to gather more data
- Challenge: even if this is possible to use old data to estimate multiple gradients, how many steps should be taken?

Policy gradient algorithms are stochastic gradient ascent:

$$\theta_{k+1} = \theta_k + \alpha_k \hat{g}_k$$

with step  $\Delta_k = \alpha_k \hat{g}_k$ .

- If the step is too large, **performance collapse** is possible (Why?)

## Choosing a Step Size for Policy Gradients

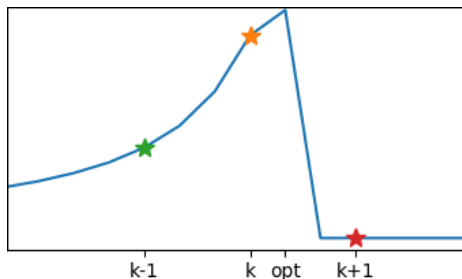
Policy gradient algorithms are stochastic gradient ascent:

$$\theta_{k+1} = \theta_k + \alpha_k \hat{g}_k$$

with step  $\Delta_k = \alpha_k \hat{g}_k$ .

- If the step is too large, **performance collapse** is possible (Why?)
- If the step is too small, progress is unacceptably slow
- “Right” step size changes based on  $\theta$

Automatic learning rate adjustment like advantage normalization, or Adam-style optimizers, can help. But does this solve the problem?



**Figure:** Policy parameters on x-axis and performance on y-axis. A bad step can lead to performance collapse, which may be hard to recover from.

# The Problem is More Than Step Size

Consider a family of policies with parametrization:

$$\pi_{\theta}(a) = \begin{cases} \sigma(\theta) & a = 1 \\ 1 - \sigma(\theta) & a = 2 \end{cases}$$

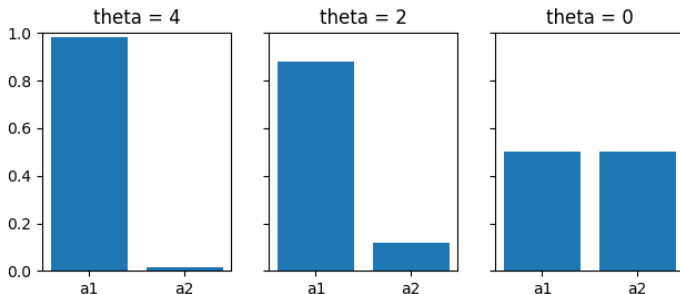


Figure: Small changes in the policy parameters can unexpectedly lead to **big** changes in the policy.

Big question: how do we come up with an update rule that doesn't ever change the policy more than we meant to?

## Policy Performance Bounds



## Relative Performance of Two Policies

In a policy optimization algorithm, we want an update step that

- uses rollouts collected from the most recent policy as efficiently as possible,
- and takes steps that respect **distance in policy space** as opposed to distance in parameter space.

To figure out the right update rule, we need to exploit relationships between the performance of two policies.

**Performance difference lemma:** In CS234 HW2 we ask you to prove that for any policies  $\pi, \pi'$

$$J(\pi') - J(\pi) = \mathbb{E}_{\tau \sim \pi'} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi}(s_t, a_t) \right] \quad (6)$$

$$= \frac{1}{1-\gamma} \mathbb{E}_{\substack{s \sim d^{\pi'} \\ a \sim \pi'}} [A^{\pi}(s, a)] \quad (7)$$

where

$$d^{\pi}(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi)$$

# What is it good for?

Can we use this for policy improvement, where  $\pi'$  represents the new policy and  $\pi$  represents the old one?

$$\begin{aligned}\max_{\pi'} J(\pi') &= \max_{\pi'} J(\pi') - J(\pi) \\ &= \max_{\pi'} \mathbb{E}_{\tau \sim \pi'} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi}(s_t, a_t) \right]\end{aligned}$$

This is suggestive, but not useful yet.

Nice feature of this optimization problem: defines the performance of  $\pi'$  in terms of the advantages from  $\pi$ !

But, problematic feature: still requires trajectories sampled from  $\pi'$ ...

In terms of the **discounted future state distribution**  $d^\pi$ , defined by

$$d^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi),$$

we can rewrite the relative policy performance identity:

$$J(\pi') - J(\pi) = \mathbb{E}_{\tau \sim \pi'} \left[ \sum_{t=0}^{\infty} \gamma^t A^\pi(s_t, a_t) \right]$$

## Note: Instance of Importance Sampling

In terms of the **discounted future state distribution**  $d^\pi$ , defined by

$$d^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi),$$

we can rewrite the relative policy performance identity:

$$\begin{aligned} J(\pi') - J(\pi) &= \mathbb{E}_{\tau \sim \pi'} \left[ \sum_{t=0}^{\infty} \gamma^t A^\pi(s_t, a_t) \right] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d^{\pi'} \\ a \sim \pi'}} [A^\pi(s, a)] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d^{\pi'} \\ a \sim \pi}} \left[ \frac{\pi'(a|s)}{\pi(a|s)} A^\pi(s, a) \right] \end{aligned}$$

Last step is an instance of **importance sampling** (more on this next time)

In terms of the **discounted future state distribution**  $d^\pi$ , defined by

$$d^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi),$$

we can rewrite the relative policy performance identity:

$$\begin{aligned} J(\pi') - J(\pi) &= \mathbb{E}_{\tau \sim \pi'} \left[ \sum_{t=0}^{\infty} \gamma^t A^\pi(s_t, a_t) \right] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d^{\pi'} \\ a \sim \pi'}} [A^\pi(s, a)] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d^{\pi'} \\ a \sim \pi}} \left[ \frac{\pi'(a|s)}{\pi(a|s)} A^\pi(s, a) \right] \end{aligned}$$

...almost there! Only problem is  $s \sim d^{\pi'}$ .

What if we just said  $d^{\pi'} \approx d^{\pi}$  and didn't worry about it?

$$\begin{aligned} J(\pi') - J(\pi) &\approx \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d^{\pi'} \\ a \sim \pi}} \left[ \frac{\pi'(a|s)}{\pi(a|s)} A^{\pi}(s, a) \right] \\ &\doteq \mathcal{L}_{\pi}(\pi') \end{aligned}$$

Turns out: this approximation is pretty good when  $\pi'$  and  $\pi$  are close! But why, and how close do they have to be?

**Relative policy performance bounds:**<sup>1</sup>

$$|J(\pi') - (J(\pi) + \mathcal{L}_{\pi}(\pi'))| \leq C \sqrt{\mathbb{E}_{s \sim d^{\pi}} [D_{KL}(\pi' || \pi)[s]]} \quad (8)$$

If policies are close in KL-divergence—the approximation is good!

---

<sup>1</sup>Achiam, Held, Tamar, Abbeel, 2017

# What is KL-divergence?

For probability distributions  $P$  and  $Q$  over a discrete random variable,

$$D_{KL}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

Properties:

- $D_{KL}(P||P) = 0$
- $D_{KL}(P||Q) \geq 0$
- $D_{KL}(P||Q) \neq D_{KL}(Q||P)$  — Non-symmetric!

What is KL-divergence between policies?

$$D_{KL}(\pi' || \pi)[s] = \sum_{a \in \mathcal{A}} \pi'(a|s) \log \frac{\pi'(a|s)}{\pi(a|s)}$$

What did we gain from making that approximation?

$$J(\pi') - J(\pi) \approx \mathcal{L}_\pi(\pi')$$

$$\begin{aligned}\mathcal{L}_\pi(\pi') &= \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d^\pi \\ a \sim \pi}} \left[ \frac{\pi'(a|s)}{\pi(a|s)} A^\pi(s, a) \right] \\ &= \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \frac{\pi'(a_t|s_t)}{\pi(a_t|s_t)} A^\pi(s_t, a_t) \right]\end{aligned}$$

- This is something we can optimize using trajectories sampled from the old policy  $\pi$ !
- Similar to using importance sampling, but because weights only depend on current timestep (and not preceding history), they don't vanish or explode.



- “Approximately Optimal Approximate Reinforcement Learning,” Kakade and Langford, 2002 <sup>2</sup>
- “Trust Region Policy Optimization,” Schulman et al. 2015 <sup>3</sup>
- “Constrained Policy Optimization,” Achiam et al. 2017 <sup>4</sup>

---

<sup>2</sup><https://people.eecs.berkeley.edu/~pabbeel/cs287-fa09/readings/KakadeLangford-icml2002.pdf>

<sup>3</sup><https://arxiv.org/pdf/1502.05477.pdf>

<sup>4</sup><https://arxiv.org/pdf/1705.10528.pdf>

# Algorithms

Proximal Policy Optimization (PPO) is a family of methods that approximately penalize policies from changing too much between steps. Two variants:

- Adaptive KL Penalty
  - Policy update solves unconstrained optimization problem

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}(\theta) - \beta_k \bar{D}_{KL}(\theta || \theta_k) \quad (9)$$

$$\bar{D}_{KL}(\theta || \theta_k) = E_{s \sim d^{\pi_k}} D_{KL}(\theta_k(\cdot | s), \pi_{\theta}(\cdot | s)) \quad (10)$$

- Penalty coefficient  $\beta_k$  changes between iterations to approximately enforce KL-divergence constraint

---

## Algorithm PPO with Adaptive KL Penalty

---

Input: initial policy parameters  $\theta_0$ , initial KL penalty  $\beta_0$ , target KL-divergence  $\delta$

**for**  $k = 0, 1, 2, \dots$  **do**

Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}(\theta) - \beta_k \bar{D}_{KL}(\theta || \theta_k)$$

by taking  $K$  steps of minibatch SGD (via Adam)

**if**  $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \geq 1.5\delta$  **then**

$$\beta_{k+1} = 2\beta_k$$

**else if**  $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \leq \delta/1.5$  **then**

$$\beta_{k+1} = \beta_k/2$$

**end if**

**end for**

---

- Initial KL penalty not that important—it adapts quickly
- Some iterations may violate KL constraint, but most don't

---

## Algorithm PPO with Adaptive KL Penalty

---

Input: initial policy parameters  $\theta_0$ , initial KL penalty  $\beta_0$ , target KL-divergence  $\delta$

**for**  $k = 0, 1, 2, \dots$  **do**

Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}(\theta) - \beta_k \bar{D}_{KL}(\theta || \theta_k)$$

**by taking**  $K$  **steps of minibatch SGD (via Adam)**

**if**  $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \geq 1.5\delta$  **then**

$$\beta_{k+1} = 2\beta_k$$

**else if**  $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \leq \delta/1.5$  **then**

$$\beta_{k+1} = \beta_k/2$$

**end if**

**end for**

---

- Initial KL penalty not that important—it adapts quickly
- Some iterations may violate KL constraint, but most don't

Proximal Policy Optimization (PPO) is a family of methods that approximately enforce KL constraint **without computing natural gradients**. Two variants:

- Adaptive KL Penalty
  - Policy update solves unconstrained optimization problem

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}(\theta) - \beta_k \bar{D}_{KL}(\theta || \theta_k)$$

- Penalty coefficient  $\beta_k$  changes between iterations to approximately enforce KL-divergence constraint
- Clipped Objective
  - New objective function: let  $r_t(\theta) = \pi_{\theta}(a_t | s_t) / \pi_{\theta_k}(a_t | s_t)$ . Then

$$\mathcal{L}_{\theta_k}^{CLIP}(\theta) = \mathbb{E}_{\tau \sim \pi_k} \left[ \sum_{t=0}^T \left[ \min(r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^{\pi_k}) \right] \right]$$

where  $\epsilon$  is a hyperparameter (maybe  $\epsilon = 0.2$ )

- Policy update is  $\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}^{CLIP}(\theta)$

## L6 Check Your Understanding: Proximal Policy Optimization

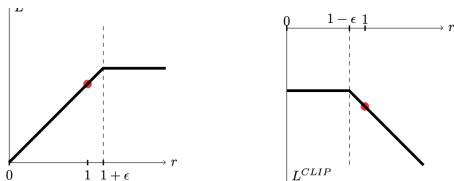
- Clipped Objective function: let  $r_t(\theta) = \pi_\theta(a_t|s_t)/\pi_{\theta_k}(a_t|s_t)$ . Then

$$\mathcal{L}_{\theta_k}^{CLIP}(\theta) = \mathbb{E}_{\tau \sim \pi_k} \left[ \sum_{t=0}^T \left[ \min(r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^{\pi_k}) \right] \right]$$

- where  $\epsilon$  is a hyperparameter (maybe  $\epsilon = 0.2$ )
- Policy update is  $\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}^{CLIP}(\theta)$ .

Consider the figure<sup>5</sup>. Select all that are true.  $\epsilon \in (0, 1)$ .

- 1 The left graph shows the  $L^{CLIP}$  objective when the advantage function  $A > 0$  and the right graph shows when  $A < 0$
- 2 The right graph shows the  $L^{CLIP}$  objective when the advantage function  $A > 0$  and the left graph shows when  $A < 0$
- 3 It depends on the value of  $\epsilon$
- 4 Not sure



<sup>5</sup>Schulman, Wolski, Dhariwal, Radford, Klimov, 2017

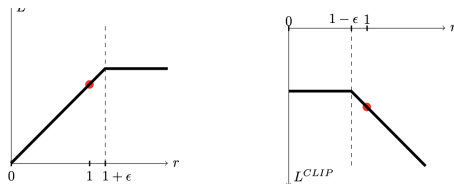
## L6: Check Your Understanding Proximal Policy Optimization Solutions

- Clipped Objective function: let  $r_t(\theta) = \pi_\theta(a_t|s_t)/\pi_{\theta_k}(a_t|s_t)$ . Then

$$\mathcal{L}_{\theta_k}^{CLIP}(\theta) = \mathbb{E}_{\tau \sim \pi_k} \left[ \sum_{t=0}^T \left[ \min(r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^{\pi_k}) \right] \right]$$

- where  $\epsilon$  is a hyperparameter (maybe  $\epsilon = 0.2$ )
- Policy update is  $\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}^{CLIP}(\theta)$ .

Consider the figure<sup>6</sup>. Select all that are true.  $\epsilon \in (0, 1)$ .

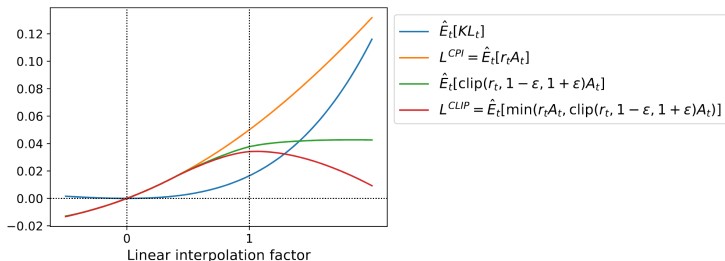


<sup>6</sup>Schulman, Wolski, Dhariwal, Radford, Klimov, 2017



# Proximal Policy Optimization with Clipped Objective

But *how* does clipping keep policy close? By making objective as pessimistic as possible about performance far away from  $\theta_k$ :



**Figure:** Various objectives as a function of interpolation factor  $\alpha$  between  $\theta_{k+1}$  and  $\theta_k$  after one update of PPO-Clip <sup>7</sup>

<sup>7</sup>Schulman, Wolski, Dhariwal, Radford, Klimov, 2017

---

## Algorithm PPO with Clipped Objective

---

Input: initial policy parameters  $\theta_0$ , clipping threshold  $\epsilon$

**for**  $k = 0, 1, 2, \dots$  **do**

Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}^{CLIP}(\theta)$$

by taking  $K$  steps of minibatch SGD (via Adam), where

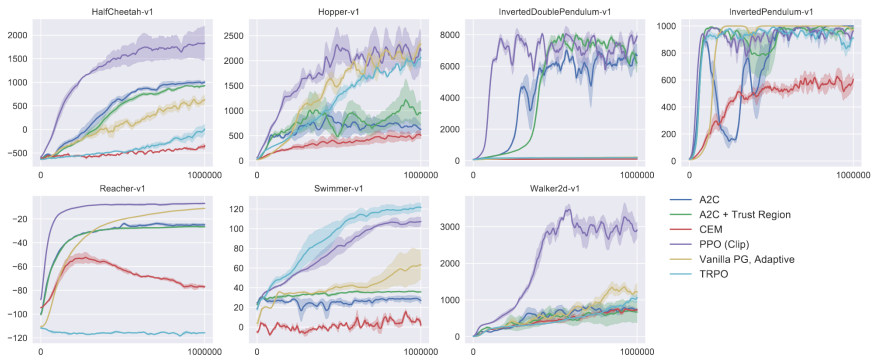
$$\mathcal{L}_{\theta_k}^{CLIP}(\theta) = \mathbb{E}_{\tau \sim \pi_k} \left[ \sum_{t=0}^T \left[ \min(r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^{\pi_k}) \right] \right]$$

**end for**

---

- Clipping prevents policy from having incentive to go far away from  $\theta_{k+1}$
- Clipping seems to work at least as well as PPO with KL penalty, but is simpler to implement

# Empirical Performance of PPO



**Figure:** Performance comparison between PPO with clipped objective and various other deep RL methods on a slate of MuJoCo tasks. <sup>8</sup>

- Wildly popular, and key component of ChatGPT

<sup>8</sup>Schulman, Wolski, Dhariwal, Radford, Klimov, 2017

### PPO

- “Proximal Policy Optimization Algorithms,” Schulman et al. 2017 <sup>9</sup>
- OpenAI blog post on PPO, 2017 <sup>10</sup>

---

<sup>9</sup><https://arxiv.org/pdf/1707.06347.pdf>

<sup>10</sup><https://blog.openai.com/openai-baselines-ppo/>

- Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph, and Aleksander Madry. Implementation Matters in Deep RL: A Case Study on PPO and TRPO. ICLR 2020  
<https://openreview.net/forum?id=r1etN1rtPB>
- Reward scaling, learning rate annealing, etc. can make a significant difference

- Likelihood ratio / score function policy gradient
  - Baseline
  - Alternative targets
- Advanced policy gradient methods
  - Proximal policy optimization (PPO) algorithm (will implement in homework)

- Last time: Policy Search
- This time: Policy search continued.
- Next time: Proximal Policy Optimization (PPO) cont (theory and additional discussion)

- SLIDES FOR NEXT CLASS (LIKELY)



# Monotonic Improvement Theory

From the bound on the previous slide, we get

$$J(\pi') - J(\pi) \geq \mathcal{L}_\pi(\pi') - C \sqrt{\mathbb{E}_{s \sim d^\pi} [D_{KL}(\pi' || \pi)[s]]}.$$

- If we maximize the RHS with respect to  $\pi'$ , we are **guaranteed to improve over  $\pi$** .
  - This is a *majorize-maximize* algorithm w.r.t. the true objective, the LHS.
- **And**  $\mathcal{L}_\pi(\pi')$  and the KL-divergence term *can both be estimated with samples from  $\pi$* !

# Monotonic Improvement Theory

Proof of improvement guarantee: Suppose  $\pi_{k+1}$  and  $\pi_k$  are related by

$$\pi_{k+1} = \arg \max_{\pi'} \mathcal{L}_{\pi_k}(\pi') - C \sqrt{\mathbb{E}_{s \sim d^{\pi_k}} [D_{KL}(\pi' || \pi_k)[s]]}.$$

Proof of improvement guarantee: Suppose  $\pi_{k+1}$  and  $\pi_k$  are related by

$$\pi_{k+1} = \arg \max_{\pi'} \mathcal{L}_{\pi_k}(\pi') - C \sqrt{\mathbb{E}_{s \sim d^{\pi_k}} [D_{KL}(\pi' || \pi_k)[s]]}.$$

- $\pi_k$  is a feasible point, and the objective at  $\pi_k$  is equal to 0.
  - $\mathcal{L}_{\pi_k}(\pi_k) \propto \mathbb{E}_{s, a \sim d^{\pi_k}, \pi_k} [A^{\pi_k}(s, a)] = 0$
  - $D_{KL}(\pi_k || \pi_k)[s] = 0$
- $\implies$  optimal value  $\geq 0$
- $\implies$  by the performance bound,  $J(\pi_{k+1}) - J(\pi_k) \geq 0$

This proof works even if we restrict the domain of optimization to an arbitrary class of parametrized policies  $\Pi_\theta$ , as long as  $\pi_k \in \Pi_\theta$ .

$$\pi_{k+1} = \arg \max_{\pi'} \mathcal{L}_{\pi_k}(\pi') - C \sqrt{\mathbb{E}_{s \sim d^{\pi_k}} [D_{KL}(\pi' || \pi_k)[s]]}. \quad (11)$$

Problem:

- $C$  provided by theory is quite high when  $\gamma$  is near 1
- $\implies$  steps from (11) are too small.

Potential Solution:

- Tune the KL penalty
- Use KL constraint (called **trust region**).

## Importance Sampling for Off Policy, Policy Gradient

# Importance Sampling

Importance sampling is a technique for estimating expectations using samples drawn from a different distribution.

$$\mathbb{E}_{x \sim P} [f(x)] =$$

Importance sampling is a technique for estimating expectations using samples drawn from a different distribution.

$$\mathbb{E}_{x \sim P} [f(x)] = \mathbb{E}_{x \sim Q} \left[ \frac{P(x)}{Q(x)} f(x) \right] \approx \frac{1}{|D|} \sum_{x \in D} \frac{P(x)}{Q(x)} f(x), \quad D \sim Q$$

The ratio  $P(x)/Q(x)$  is the **importance sampling weight** for  $x$ .



Importance sampling is a technique for estimating expectations using samples drawn from a different distribution.

$$\mathbb{E}_{x \sim P} [f(x)] = \mathbb{E}_{x \sim Q} \left[ \frac{P(x)}{Q(x)} f(x) \right] \approx \frac{1}{|D|} \sum_{x \in D} \frac{P(x)}{Q(x)} f(x), \quad D \sim Q$$

The ratio  $P(x)/Q(x)$  is the **importance sampling weight** for  $x$ .

What is the variance of an importance sampling estimator?

$$\begin{aligned} \text{var}(\hat{\mu}_Q) &= \frac{1}{N} \text{var} \left( \frac{P(x)}{Q(x)} f(x) \right) \\ &= \frac{1}{N} \left( \mathbb{E}_{x \sim Q} \left[ \left( \frac{P(x)}{Q(x)} f(x) \right)^2 \right] - \mathbb{E}_{x \sim Q} \left[ \frac{P(x)}{Q(x)} f(x) \right]^2 \right) \\ &= \frac{1}{N} \left( \mathbb{E}_{x \sim P} \left[ \frac{P(x)}{Q(x)} f(x)^2 \right] - \mathbb{E}_{x \sim P} [f(x)]^2 \right) \end{aligned}$$

The term in red is problematic—if  $P(x)/Q(x)$  is large in the wrong places, the variance of the estimator explodes.

Here, we compress the notation  $\pi_\theta$  down to  $\theta$  in some places for compactness.

$$\begin{aligned}g &= \nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \theta} \left[ \sum_{t=0}^{\infty} \gamma^t \nabla_\theta \log \pi_\theta(a_t | s_t) A^\theta(s_t, a_t) \right] \\&= \sum_{\tau} \sum_{t=0}^{\infty} \gamma^t P(\tau_t | \theta) \nabla_\theta \log \pi_\theta(a_t | s_t) A^\theta(s_t, a_t) \\&= \mathbb{E}_{\tau \sim \theta'} \left[ \sum_{t=0}^{\infty} \frac{P(\tau_t | \theta)}{P(\tau_t | \theta')} \gamma^t \nabla_\theta \log \pi_\theta(a_t | s_t) A^\theta(s_t, a_t) \right]\end{aligned}$$

Here, we compress the notation  $\pi_\theta$  down to  $\theta$  in some places for compactness.

$$\begin{aligned}g &= \nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \theta} \left[ \sum_{t=0}^{\infty} \gamma^t \nabla_\theta \log \pi_\theta(a_t | s_t) A^\theta(s_t, a_t) \right] \\&= \sum_{\tau} \sum_{t=0}^{\infty} \gamma^t P(\tau_t | \theta) \nabla_\theta \log \pi_\theta(a_t | s_t) A^\theta(s_t, a_t) \\&= \mathbb{E}_{\tau \sim \theta'} \left[ \sum_{t=0}^{\infty} \frac{P(\tau_t | \theta)}{P(\tau_t | \theta')} \gamma^t \nabla_\theta \log \pi_\theta(a_t | s_t) A^\theta(s_t, a_t) \right]\end{aligned}$$

$$\frac{P(\tau_t | \theta)}{P(\tau_t | \theta')} =$$

# Importance Sampling for Policy Gradients

Here, we compress the notation  $\pi_\theta$  down to  $\theta$  in some places for compactness.

$$\begin{aligned} g &= \nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \theta} \left[ \sum_{t=0}^{\infty} \gamma^t \nabla_\theta \log \pi_\theta(a_t | s_t) A^\theta(s_t, a_t) \right] \\ &= \sum_{\tau} \sum_{t=0}^{\infty} \gamma^t P(\tau_t | \theta) \nabla_\theta \log \pi_\theta(a_t | s_t) A^\theta(s_t, a_t) \\ &= \mathbb{E}_{\tau \sim \theta'} \left[ \sum_{t=0}^{\infty} \frac{P(\tau_t | \theta)}{P(\tau_t | \theta')} \gamma^t \nabla_\theta \log \pi_\theta(a_t | s_t) A^\theta(s_t, a_t) \right] \end{aligned}$$

Challenge? **Exploding or vanishing importance sampling weights.**

$$\frac{P(\tau_t | \theta)}{P(\tau_t | \theta')} = \frac{\mu(s_0) \prod_{t'=0}^t P(s_{t'+1} | s_{t'}, a_{t'}) \pi_\theta(a_{t'} | s_{t'})}{\mu(s_0) \prod_{t'=0}^t P(s_{t'+1} | s_{t'}, a_{t'}) \pi_{\theta'}(a_{t'} | s_{t'})} = \prod_{t'=0}^t \frac{\pi_\theta(a_{t'} | s_{t'})}{\pi_{\theta'}(a_{t'} | s_{t'})}$$

Even for policies only slightly different from each other, **many small differences multiply to become a big difference.**

Big question: how can we make efficient use of the data we already have from the old policy, while avoiding the challenges posed by importance sampling?

## Theory:

- 1 Problems with Policy Gradient Methods
- 2 Policy Performance Bounds
- 3 Monotonic Improvement Theory

## Proximal Policy Optimization:

- 1 Approximately constraints policy steps
- 2 Relatively simple to implement
- 3 Good empirical success and very widely used

$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} R_t^i \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

- Note that critic can select any blend between TD and MC estimators for the target to substitute for the true state-action value function.

$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} R_t^i \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

- Note that critic can select any blend between TD and MC estimators for the target to substitute for the true state-action value function.

$$\hat{R}_t^{(1)} = r_t + \gamma V(s_{t+1})$$

$$\hat{R}_t^{(2)} = r_t + \gamma r_{t+1} + \gamma^2 V(s_{t+2}) \quad \dots$$

$$\hat{R}_t^{(\text{inf})} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

- If subtract baselines from the above, get advantage estimators

$$\hat{A}_t^{(1)} = r_t + \gamma V(s_{t+1}) - V(s_t)$$

$$\hat{A}_t^{(\text{inf})} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots - V(s_t)$$

$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} R_t^i \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

- If subtract baselines from the above, get advantage estimators

$$\hat{A}_t^{(1)} = r_t + \gamma V(s_{t+1}) - V(s_t)$$

$$\hat{A}_t^{(\text{inf})} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+1} + \dots - V(s_t)$$

- Select all that are true
- $\hat{A}_t^{(1)}$  has low variance & low bias.
- $\hat{A}_t^{(1)}$  has high variance & low bias.
- $\hat{A}_t^{(\infty)}$  low variance and high bias.
- $\hat{A}_t^{(\infty)}$  high variance and low bias.
- Not sure



$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} R_t^i \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

- If subtract baselines from the above, get advantage estimators

$$\hat{A}_t^{(1)} = r_t + \gamma V(s_{t+1}) - V(s_t)$$

$$\hat{A}_t^{(\text{inf})} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+1} + \dots - V(s_t)$$