



# CSC415: Introduction to Reinforcement Learning

## Lecture 6: Advanced Policy Gradients

Dr. Amey Pore

Winter 2026

February 11, 2026

Structure adapted from Chelsea Finn's course on Deep RL and Emma Brunskill's course on Introduction to RL.

# Course Logistics

- **Project Teams:**

- Still many unassigned teams on Quercus.
- Students not in a team will be assigned randomly soon.
- **Team Building:** put a note of the topic you want to work on the sticky note.

# Think pair wise

Which of the following are true about REINFORCE? In the following options, PG stands for policy gradient.

- ☐ A Adding a baseline term can help to reduce the variance of the PG updates
- ☐ B It will converge to a global optima
- ☐ C It can be initialized with a sub-optimal, deterministic policy and still converge to a local optima, given the appropriate step sizes
- ☐ D If we take one step of PG, it is possible that the resulting policy gets worse (in terms of achieved returns) than our initial policy

# Outline

- ➊ **Recall**
- ➋ Problems with Policy Gradient Methods
- ➌ Policy Performance Bounds
- ➍ Proximal Policy Optimization Algorithm
- ➎ Generalized Advantage Estimation
- ➏ Off-policy Actor Critic
- ➐ Monotonic Improvement Theory

# REINFORCE Algorithm

- Using policy gradient theorem
- Using return  $G_t$  as an unbiased estimate of  $Q^{\pi_\theta}(s_t, a_t)$
- Stochastic gradient ascent update:

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

---

```
1: Initialize policy parameters  $\theta$  arbitrarily
2: for each episode  $\{s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_\theta$  do
3:   for  $t = 1$  to  $T - 1$  do
4:      $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(s_t | a_t) G_t$ 
5:   end for
6: end for
7: return  $\theta$ 
```

---

Very high variance!

# Action-Value Actor-Critic

- Simple actor-critic algorithm based on action-value critic
  - Using linear value fn approx.  $Q_w(s, a) = \phi(s, a)^\top w$ 
    - **Critic** Updates  $w$  by linear TD(0)
    - **Actor** Updates  $\theta$  by policy gradient
- 

```

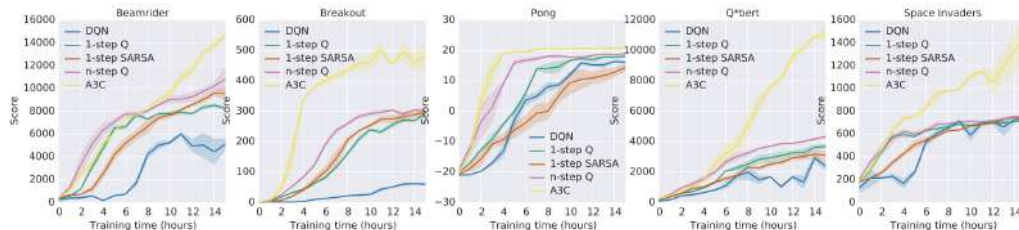
1: function Q-Actor Critic
2:   Initialize  $s, \theta$ 
3:   Sample action  $a \sim \pi_\theta(\cdot|s)$ 
4:   for each step do
5:     Sample reward  $r = R(s, a)$ ; sample transition  $s' \sim P(\cdot|s, a)$ 
6:     Sample action  $a' \sim \pi_\theta(\cdot|s')$ 
7:      $\delta \leftarrow r + \gamma Q_w(s', a') - Q_w(s, a)$ 
8:      $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(a|s) Q_w(s, a)$ 
9:      $w \leftarrow w + \beta \delta \phi(s, a)$ 
10:     $a \leftarrow a', s \leftarrow s'$ 
11:  end for
12: end function
  
```

# “Vanilla” Policy Gradient Algorithm

- 1: Initialize policy parameter  $\theta$ , baseline  $b$
- 2: **for** iteration = 1, 2,  $\dots$  **do**
- 3:     Collect a set of trajectories by executing the current policy
- 4:     At each timestep  $t$  in each trajectory  $\tau^i$ , compute:
- 5:         Return  $G_t^i = \sum_{t'=t}^{T-1} r_{t'}^i$
- 6:         Advantage estimate  $\hat{A}_t^i = G_t^i - b(s_t)$
- 7:         Re-fit the baseline, by minimizing  $\sum_i \sum_t \|b(s_t) - G_t^i\|^2$
- 8:         Update the policy, using a policy gradient estimate  $\hat{g}$ :
- 9:              $\hat{g} = \sum_i \sum_t \nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_t^i$
- 10:         (Plug  $\hat{g}$  into SGD or ADAM)
- 11: **end for**

- Other Baseline:  $\underbrace{A^{\pi}(s, a) = Q_w^{\pi}(s, a) - V_w^{\pi}(s)}_{\text{Advantage Actor Critic}}$

# Asynchronous Advantage Actor-Critic (A3C)



**Figure:** A3C: Multiple workers interact with their own environments and update a global network asynchronously.

Mnih, Badia, Mirza, Graves, Lillicrap, Harley, Silver, Kavukcuoglu. *Asynchronous Methods for Deep Reinforcement Learning*. ICML 2016.



# Outline

- 1 Recall
- 2 **Problems with Policy Gradient Methods**
- 3 Policy Performance Bounds
- 4 Proximal Policy Optimization Algorithm
- 5 Generalized Advantage Estimation
- 6 Off-policy Actor Critic
- 7 Monotonic Improvement Theory

# Policy Gradients Review

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  $\Pi = \{ \pi : \pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}, \sum_a \pi_{sa} = 1, \pi_{sa} \geq 0 \}$
  - Policy gradients take steps in parameter space
  - Step size is hard to get right as a result

# Sample Efficiency in Policy Gradients

- 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?

# 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?)

# Choosing a Step Size for Policy Gradients

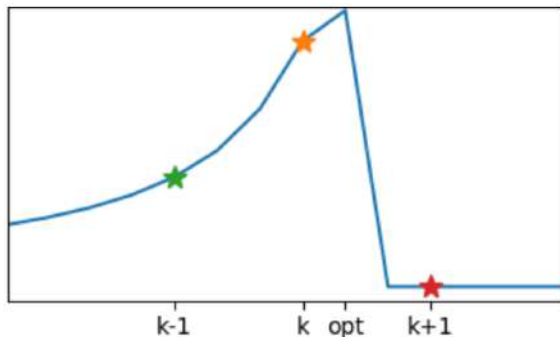
Policy gradient algorithms are stochastic gradient ascent:

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

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?

## Choosing a Step Size for Policy Gradients

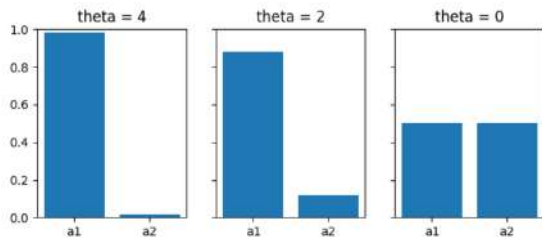


**Figure:** Policy parameters on x-axis and performance on y-axis. A bad step can lead to performance collapse, which may be hard to recover from.

# The Problem is More Than Step Size

Consider a family of policies with parametrization:

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



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

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

# Outline

- 1 Recall
- 2 Problems with Policy Gradient Methods
- 3 **Policy Performance Bounds**
- 4 Proximal Policy Optimization Algorithm
- 5 Generalized Advantage Estimation
- 6 Off-policy Actor Critic
- 7 Monotonic Improvement Theory



## 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:** For any policies  $\pi, \pi'$

$$J(\pi') - J(\pi) = \mathbb{E}_{\tau \sim \pi'} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi}(s_t, a_t) \right] \quad (1)$$

$$= \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d^{\pi'} \\ a \sim \pi'}} [A^{\pi}(s, a)] \quad (2)$$

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?

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

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  $\pi'$ ...

## Looking at it from Another Angle...

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:

$$J(\pi') - J(\pi) = \mathbb{E}_{\tau \sim \pi'} \left[ \sum_{t=0}^{\infty} \gamma^t A^\pi(s_t, a_t) \right] \quad (3)$$

$$= \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d^{\pi'} \\ a \sim \pi'}} [A^\pi(s, a)] \quad (4)$$

$$= \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d^{\pi'} \\ a \sim \pi}} \left[ \frac{\pi'(a|s)}{\pi(a|s)} A^\pi(s, a) \right] \quad (5)$$

Last step is an instance of importance sampling. **almost there! Only problem is  $s \sim d^{\pi'}$ .**

## A Useful Approximation

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

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

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>

$$J(\pi') \geq J(\pi) + L_{\pi}(\pi') - C \sqrt{\mathbb{E}_{s \sim d^{\pi}} [D_{\text{KL}}(\pi' \parallel \pi)[s]]}$$

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,

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

### Properties:

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

### What is KL-divergence between policies?

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

## A Useful Approximation

What did we gain from making that approximation?

$$J(\pi') - J(\pi) \approx L_{\pi}(\pi')$$

$$L_{\pi}(\pi') = \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d^{\pi} \\ a \sim \pi}} \left[ \frac{\pi'(a|s)}{\pi(a|s)} A^{\pi}(s, a) \right] \quad (6)$$

$$= \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \frac{\pi'(a_t|s_t)}{\pi(a_t|s_t)} A^{\pi}(s_t, a_t) \right] \quad (7)$$

- 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.

# Recommended Reading

- “Approximately Optimal Approximate Reinforcement Learning,” Kakade and Langford, 2002
- “Trust Region Policy Optimization,” Schulman et al. 2015
- “Constrained Policy Optimization,” Achiam et al. 2017



# Outline

- ➊ Recall
- ➋ Problems with Policy Gradient Methods
- ➌ Policy Performance Bounds
- ➍ **Proximal Policy Optimization Algorithm**
- ➎ Generalized Advantage Estimation
- ➏ Off-policy Actor Critic
- ➐ Monotonic Improvement Theory

# Proximal Policy Optimization

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

$$\theta_{k+1} = \arg \max_{\theta} L_{\theta_k}(\theta) - \beta_k \bar{D}_{\text{KL}}(\theta \| \theta_k) \quad (8)$$

$$\bar{D}_{\text{KL}}(\theta \| \theta_k) = \mathbb{E}_{s \sim d^{\pi_k}} [D_{\text{KL}}(\pi_{\theta_k}(\cdot | s), \pi_{\theta}(\cdot | s))] \quad (9)$$

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

# PPO with Adaptive KL Penalty

- 
- 1: **Input:** initial policy parameters  $\theta_0$ , initial KL penalty  $\beta_0$ , target KL-divergence  $\delta$
  - 2: **for**  $k = 0, 1, 2, \dots$  **do**
  - 3:     Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$
  - 4:     Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm
  - 5:     Compute policy update

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

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

- 6:     **if**  $\bar{D}_{\text{KL}}(\theta_{k+1} \| \theta_k) \geq 1.5\delta$  **then**
  - 7:          $\beta_{k+1} = 2\beta_k$
  - 8:     **else if**  $\bar{D}_{\text{KL}}(\theta_{k+1} \| \theta_k) \leq \delta/1.5$  **then**
  - 9:          $\beta_{k+1} = \beta_k/2$
  - 10:    **end if**
  - 11: **end for**
- 

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

# PPO with Adaptive KL Penalty: Multiple Gradient Steps

- 1: **Input:** initial policy parameters  $\theta_0$ , initial KL penalty  $\beta_0$ , target KL-divergence  $\delta$
- 2: **for**  $k = 0, 1, 2, \dots$  **do**
- 3:     Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$
- 4:     Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm
- 5:     Compute policy update

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

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

- 6:     **if**  $\bar{D}_{\text{KL}}(\theta_{k+1} \| \theta_k) \geq 1.5\delta$  **then**
- 7:          $\beta_{k+1} = 2\beta_k$
- 8:     **else if**  $\bar{D}_{\text{KL}}(\theta_{k+1} \| \theta_k) \leq \delta/1.5$  **then**
- 9:          $\beta_{k+1} = \beta_k/2$
- 10:    **end if**
- 11: **end for**

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

# Proximal Policy Optimization

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  

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

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

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

# Think Pair wise: Proximal Policy Optimization

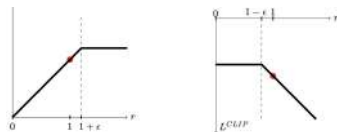
Clipped Objective function: let  $r_t(\theta) = \pi_\theta(a_t|s_t)/\pi_{\theta_k}(a_t|s_t)$ . Then

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

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

Consider the figure. Select all that are true.  $\epsilon \in (0, 1)$ .

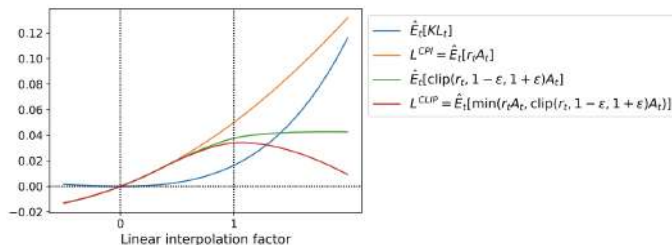
- The left graph shows the  $L^{\text{CLIP}}$  objective when the advantage function  $A > 0$  and the right graph shows when  $A < 0$
- The right graph shows the  $L^{\text{CLIP}}$  objective when the advantage function  $A > 0$  and the left graph shows when  $A < 0$
- It depends on the value of  $\epsilon$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



**Figure:** 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:** Figure from Schulman et al., 2017: Various objectives as a function of interpolation factor  $\alpha$  between  $\theta_{k+1}$  and  $\theta_k$  after one update of PPO-Clip

# Proximal Policy Optimization with Clipped Objective

- 1: **Input:** initial policy parameters  $\theta_0$ , clipping threshold  $\epsilon$
- 2: **for**  $k = 0, 1, 2, \dots$  **do**
- 3:     Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$
- 4:     Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm
- 5:     Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} L_{\theta_k}^{\text{CLIP}}(\theta)$$

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

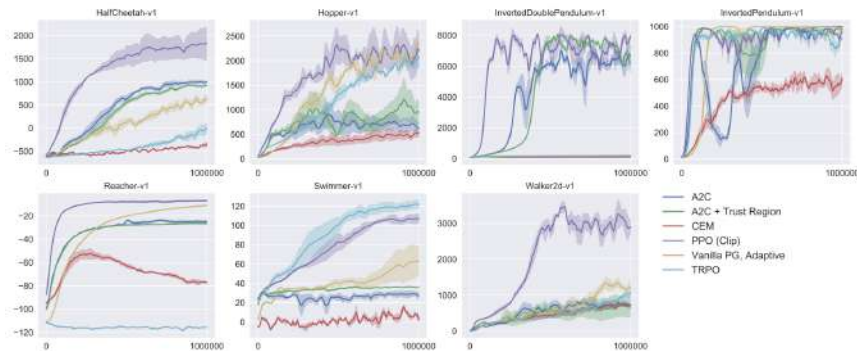
$$L_{\theta_k}^{\text{CLIP}}(\theta) = \mathbb{E}_{\tau \sim \pi_k} \left[ \sum_{t=0}^T \min \left( r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon) \hat{A}_t^{\pi_k} \right) \right]$$

6: **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.

- Wildly popular, and key component of ChatGPT

# Recommended Reading

## PPO

- “Proximal Policy Optimization Algorithms,” Schulman et al. 2017  
<https://arxiv.org/pdf/1707.06347.pdf>
- OpenAI blog post on PPO, 2017  
<https://blog.openai.com/openai-baselines-ppo/>

# PPO: Algorithm and Code Implementation Details

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

# Course Logistics

- **Mid-term Exam:** will be distributed on Feb 25th (in class).
- **Project Teams:**
  - Still many unassigned teams on Quercus.
  - Students not in a team will be assigned randomly soon.
  - **Team Building:** put a note of the topic you want to work on the sticky note.
- **A1 Code Implementations Clarifications:**
  - Goal: Reproduce the exact results reported in the paper.
  - Form a hypothesis on how the behaviour changes with variations in specific parameters.
- **Project Proposal (Per Team):**
  - Due date pushed to **March 27th**.
  - Expectation: Reading papers, proposing a new idea, and providing some preliminary results.
- **Logistics for A2:**
  - Still checking if peer review is possible for A2.
  - If not, we will have another quiz. (10%)

# Break: 10 minutes

Paste your ideas on the sticky note and find your team.

# Outline

- 1 Recall
- 2 Problems with Policy Gradient Methods
- 3 Policy Performance Bounds
- 4 Proximal Policy Optimization Algorithm
- 5 **Generalized Advantage Estimation**
- 6 Off-policy Actor Critic
- 7 Monotonic Improvement Theory

# Recall Proximal Policy Optimization

PPO is a family of methods that approximately enforce KL constraint

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

**How do we estimate the advantage function inside the policy update?**

# Recall N-step estimators

$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} A_{ti} \nabla_{\theta} \log \pi_{\theta}(a_{ti} | s_{ti})$$

- Recall the N-step advantage estimators

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

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

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

- Define  $\delta_t^V = r_t + \gamma V(s_{t+1}) - V(s_t)$ . Then

$$\hat{A}_t^{(1)} = \delta_t^V$$

$$\hat{A}_t^{(2)} = \delta_t^V + \gamma \delta_{t+1}^V$$

$$\hat{A}_t^{(k)} = \sum_{l=0}^{k-1} \gamma^l \delta_{t+l}^V$$



# Generalized Advantage Estimator (GAE)

$$\hat{A}_t^{(k)} = \sum_{l=0}^{k-1} \gamma^l r_{t+l} + \gamma^k V(s_{t+k}) - V(s_t) \quad (1)$$

- GAE is an exponentially-weighted average of  $k$ -step estimators

$$\begin{aligned} \hat{A}_t^{GAE(\gamma, \lambda)} &= (1 - \lambda)(\hat{A}_t^{(1)} + \lambda \hat{A}_t^{(2)} + \lambda^2 \hat{A}_t^{(3)} + \dots) \\ &= (1 - \lambda)(\delta_t^V + \lambda(\delta_t^V + \gamma \delta_{t+1}^V) + \lambda^2(\delta_t^V + \gamma \delta_{t+1}^V + \gamma^2 \delta_{t+2}^V) + \dots) \\ &= (1 - \lambda)(\delta_t^V(1 + \lambda + \lambda^2 + \dots) + \gamma \delta_{t+1}^V(\lambda + \lambda^2 + \dots) \\ &\quad + \gamma^2 \delta_{t+2}^V(\lambda^2 + \lambda^3 + \dots) + \dots) \\ &= (1 - \lambda)(\delta_t^V \frac{1}{1 - \lambda} + \gamma \lambda \delta_{t+1}^V \frac{1}{1 - \lambda} + \gamma^2 \lambda^2 \delta_{t+2}^V \frac{1}{1 - \lambda} + \dots) \\ &= \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}^V \end{aligned}$$

# Generalized Advantage Estimator (GAE) in PPO

- GAE is an exponentially-weighted average of  $k$ -step estimators

$$\hat{A}_t^{(k)} = \sum_{l=0}^{k-1} \gamma^l r_{t+l} + \gamma^k V(s_{t+k}) - V(s_t)$$

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

$$\begin{aligned} \hat{A}_t^{GAE(\gamma, \lambda)} &= (1 - \lambda)(\hat{A}_t^{(1)} + \lambda \hat{A}_t^{(2)} + \lambda^2 \hat{A}_t^{(3)} + \dots) \\ &= \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}^V \end{aligned}$$

- PPO uses a truncated version of a GAE

$$\hat{A}_t = \sum_{l=0}^{T-t-1} (\gamma \lambda)^l \delta_{t+l}^V$$

- Benefits: Only have to run policy in environment for  $T$  timesteps before updating,

# Final Proximal Policy Optimization

- 1: **Input:** initial policy parameters  $\theta_0$ , clipping threshold  $\epsilon$
- 2: **for**  $k = 0, 1, 2, \dots$  **do**
- 3:     Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$
- 4:     Estimate advantages  $\hat{A}_t^{\pi_k}$  using Generalized Advantage Estimation (GAE)
- 5:     Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} L_{\theta_k}^{\text{CLIP}}(\theta)$$

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

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

- 6: **end for**

**Some example hyperparameters:**

~2000 timesteps in batch of data

~10 epochs when updating policy

( $M \approx 300$  gradient steps with batch size 64)

# Outline

- 1 Recall
- 2 Problems with Policy Gradient Methods
- 3 Policy Performance Bounds
- 4 Proximal Policy Optimization Algorithm
- 5 Generalized Advantage Estimation
- 6 **Off-policy Actor Critic**
- 7 Monotonic Improvement Theory

# Off-Policy Actor-Critic Methods

So far:

- use one batch of policy data for one gradient step (fully on-policy)
- use one batch of policy data for multiple gradient steps (starting to be off-policy)

**Can we be even more off-policy?**

**Idea:** Can we reuse data from previous batches, i.e. all of the past trial-and-error data?

## Two key ideas

- 1 Maintain a **replay buffer** of all past data
- 2 Adjust equations to remove on-policy assumptions

# DQN Pseudocode

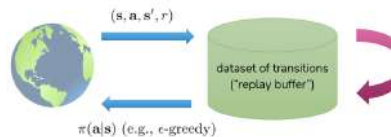
---

```

1: Input:  $E, \alpha, s, a, r, s' \sim \pi$ ; Initialize  $\mathcal{D} = \emptyset, \mathbf{w} = 0$ 
2: Set other state  $\mathbf{w}_0$ 
3: for episode = 1, ...,  $E$  do do
4:   Initialize  $s_1$ 
5:   for  $t = 1, \dots, T$  do do
6:     Observe reward  $r_t$  and next state  $s_{t+1}$ 
7:     Store transition  $(s_t, a_t, r_t, s_{t+1})$  in replay buffer  $\mathcal{D}$ 
8:     for  $i = 1, \dots, K$  do do
9:       Sample random minibatch of transitions  $(s, a, r, s')$  from  $\mathcal{D}$ 
10:      if  $s_{t+1}$  is terminal at step  $t + 1$  then then
11:        Set  $y_t = r_t$ 
12:      else
13:        Set  $y_t = r_t + \gamma \max_{a'} \hat{Q}(s_{t+1}, a'; \mathbf{w}^-)$ 
14:      end if
15:      Perform gradient descent step on  $(y_t - \hat{Q}(s_t, a_t; \mathbf{w}))^2$  w.r.t.  $\mathbf{w}$ 
16:    end for
17:    Every  $C$  steps:  $\mathbf{w}^- = \mathbf{w}$ 
18:  end for
19: end for

```

---



# Fixing the Value Function

**Solution:** Fit  $Q_w(s, a)$  instead of  $V_v(s)$ .

The datapoints we have:  $(s, a, r, s') + \text{future reward}$ .

- If we fit  $V_v(s)$ , we assume the next action comes from  $\pi$  (but  $a$  in buffer is from old policy).
- If we fit  $Q_w(s, a)$ , we pass the action as input!
- It is okay if  $a$  is different from what the current policy would have done.

**Q-Learning Update Strategy:**

$$Q_w(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a), a' \sim \pi(\cdot | s')} [Q_w(s', a')]$$

We use samples  $(s, a, r, s')$  from the replay buffer, but sample the *next action*  $a'$  from the *current policy*.

# Off Policy Actor-Critic

- 
- 1: **Input:** initial policy parameters  $\theta$ , Q-function parameters  $w$ , Replay Buffer  $\mathcal{D}$
  - 2: **for**  $k = 0, 1, 2, \dots$  **do**
  - 3:     Collect set of partial trajectories on policy  $\pi_k$  and add to  $\mathcal{D}$
  - 4:     Sample batch of transitions  $B \sim \mathcal{D}$
  - 5:     Update critic  $Q_w$  using Q-learning update strategy
  - 6:     Compute policy update
    - Original Policy Gradient:  $\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} [\sum_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t]$
    - Adapted for SAC (Q-function):

$$\nabla_{\theta} J(\theta) \approx \frac{1}{|B|} \sum_{s \in B} \nabla_{\theta} \log \pi_{\theta}(a|s) Q_w(s, a)$$

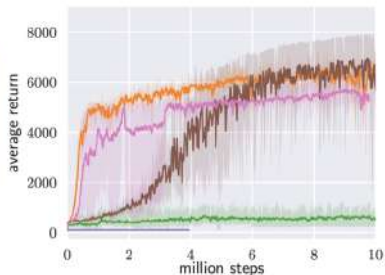
- 7: **end for**
- 

Haarnoja, Zhou, Abbeel, Levine. *Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor*. 2018.

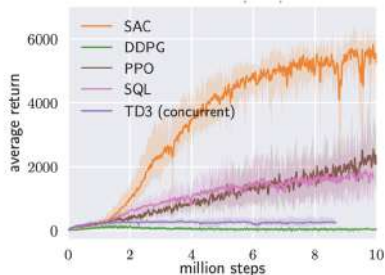
<https://sites.google.com/view/sac-and-applications/>



# Soft actor critic



(e) Humanoid-v1



(f) Humanoid (rllab)

<- more off-policy  
(i.e. with replay buffer)

<- less off-policy  
(i.e. no replay buffer)

- + Off-policy with replay buffer (e.g. soft actor-critic) can be far more data efficient
- They can also generally be a lot harder to tune hyperparameters, less stable (than e.g. PPO)

# When to use one online RL algorithm vs. another?

- **PPO & variants**

- When you care about stability, ease-of-use
- When you don't care about data efficiency

- **DQN & variants**

- When you have discrete actions or low-dimensional continuous actions

- **SAC & variants**

- When you care most about data efficiency
- When you are okay with tuning hyperparameters, less stability

# Outline

- 1 Recall
- 2 Problems with Policy Gradient Methods
- 3 Policy Performance Bounds
- 4 Proximal Policy Optimization Algorithm
- 5 Generalized Advantage Estimation
- 6 Off-policy Actor Critic
- 7 **Monotonic Improvement Theory**

# Monotonic Improvement Theory

From the bound on the previous slide, we get

$$J(\pi') - J(\pi) \geq L_{\pi}(\pi') - C\sqrt{\mathbb{E}_{s \sim d^{\pi}}[D_{\text{KL}}(\pi' \parallel \pi)[s]]}$$

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  $L_{\pi}(\pi')$  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

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

- $\pi_k$  is a feasible point, and the objective at  $\pi_k$  is equal to 0.
  - $L_{\pi_k}(\pi_k) \propto \mathbb{E}_{s, a \sim d^{\pi_k}, \pi_k} [A^{\pi_k}(s, a)] = 0$
  - $D_{\text{KL}}(\pi_k \parallel \pi_k)[s] = 0$
- $\implies$  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  $\pi_k \in \Pi_\theta$ .

# Approximate Monotonic Improvement

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

## Problem:

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

## Potential Solution:

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