

Policy Gradient Methods

(continued)

Lucas Janson and Sham Kakade

CS/Stat 184: Introduction to Reinforcement Learning
Fall 2022

Today

- Recap++
 - will clarify few points based on feedback.
- Today:
 1. Estimation of Stochastic Gradients
 2. Variance Reduction
 3. More Variance Reduction (baselines)

Recap++

(some new material and clarifications)

Recap Outline:

1. The Learning Setting
2. Objective: direct policy optimization.
3. General convergence: properties of SGD
4. Importance Sampling
& Deriving a Policy Gradient Expression

The Learning Setting:

We don't know the MDP, but we can obtain trajectories.

The Finite Horizon, Learning Setting. We can obtain trajectories as follows:

- We start at $s_0 \sim \mu_0$.
- We act for H steps and observe the trajectory $\tau = \{s_0, a_0, s_1, a_1, \dots, s_{H-1}, a_{H-1}\}$

The Infinite Horizon, Discounted Learning Setting. We can obtain trajectories as follows:

- We start at $s_0 \sim \mu_0$.
- We can obtain a “long trajectories” $\tau = \{s_0, a_0, s_1, a_1, \dots\}$
 - Suppose we can terminate the trajectory at will.
(and sufficient long trajectories will well approximate the discounted value function)

Note that with a simulator, we can sample trajectories as specified in the above.

Recap Outline:

1. The Learning Setting
2. Objective: direct policy optimization.
3. General convergence: properties of SGD
4. Importance Sampling
& Deriving a Policy Gradient Expression

Policy Optimization:

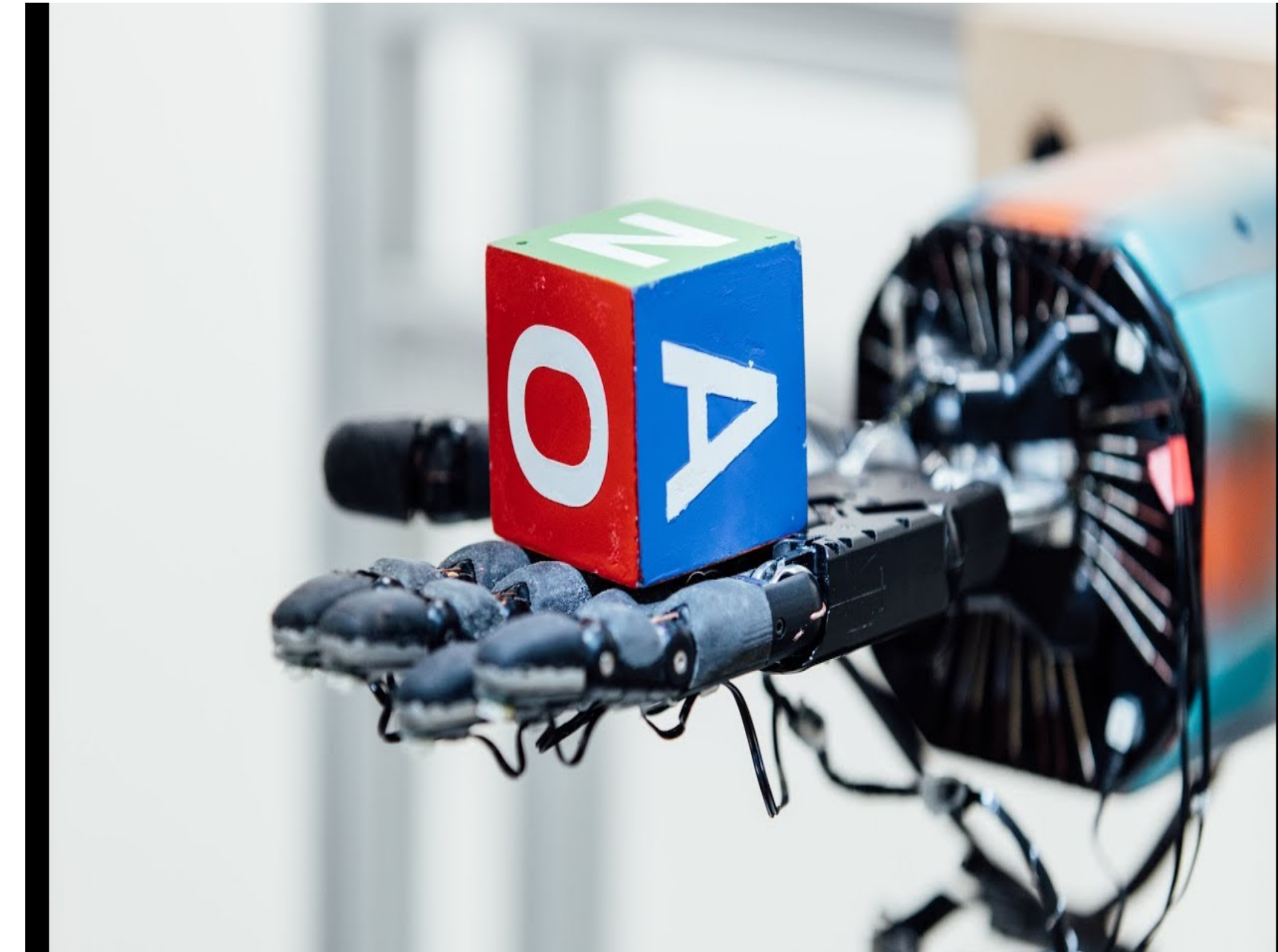
our goal is to do well on “large” problems



[AlphaZero, Silver et.al, 17]



[OpenAI Five, 18]



[OpenAI,19]

Recap: Policy Parameterization

Recall that we consider parameterized policy $\pi_\theta(\cdot | s) \in \Delta(A), \forall s$

1. Softmax linear Policy (We will try this in HW2)

Feature vector $\phi(s, a) \in \mathbb{R}^d$, and
parameter $\theta \in \mathbb{R}^d$

$$\pi_\theta(a | s) = \frac{\exp(\theta^\top \phi(s, a))}{\sum_{a'} \exp(\theta^\top \phi(s, a'))}$$

2. Neural Policy:

Neural network
 $f_\theta : S \times A \mapsto \mathbb{R}$

$$\pi_\theta(a | s) = \frac{\exp(f_\theta(s, a))}{\sum_{a'} \exp(f_\theta(s, a'))}$$

Our Objective and Policy Gradient Ascent

We consider either discounted or finite horizon settings.

$$\begin{aligned} J(\theta) &:= E_{s_0 \sim \mu_0} [V^{\pi_\theta}(s_0)] \\ &= E \left[\sum_{h=0}^{\infty} \gamma^h r_h \mid \mu_0, \pi_\theta \right] \end{aligned}$$

$$\begin{aligned} J(\theta) &:= E_{s_0 \sim \mu_0} [V^{\pi_\theta}(s_0)] \\ &= E \left[\sum_{h=0}^{H-1} r_h \mid \mu_0, \pi_\theta \right] \end{aligned}$$

- **Objective:** try to find “good” parameters

$$\max_{\theta} J(\theta)$$

- **Approach:** stochastic gradient descent (or gradient descent)

$$\theta_{t+1} = \theta_t + \eta \widetilde{\nabla}_{\theta} J(\theta_t)$$

Recap Outline:

1. The Learning Setting
2. Objective: direct policy optimization.
3. General convergence: properties of SGD
4. Importance Sampling
& Deriving a Policy Gradient Expression

Stochastic Gradient Descent

Given an objective function $J(\theta) : \mathbb{R}^d \mapsto \mathbb{R}$, (e.g., $J(\theta) = \mathbb{E}_{x,y}(f_\theta(x) - y)^2$)

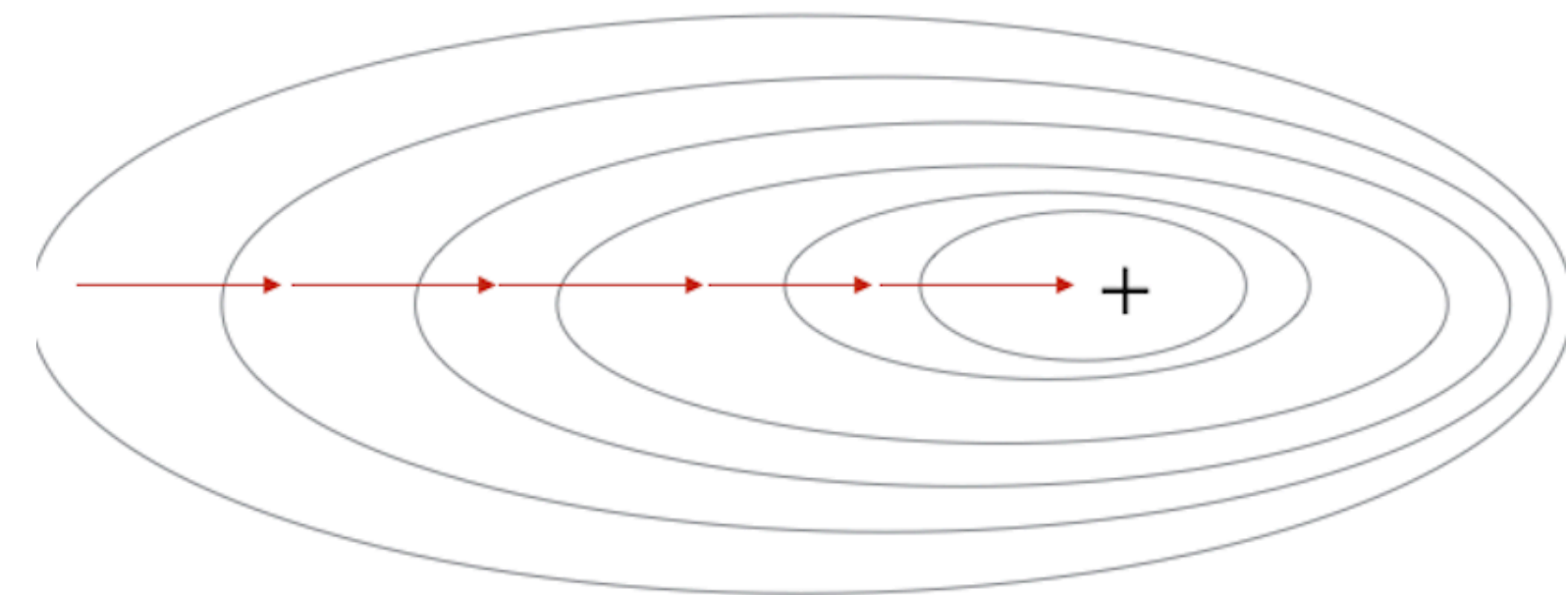
SGD minimizes the above objective function as follows:

Initialize θ_0 , for $t = 0, \dots$:

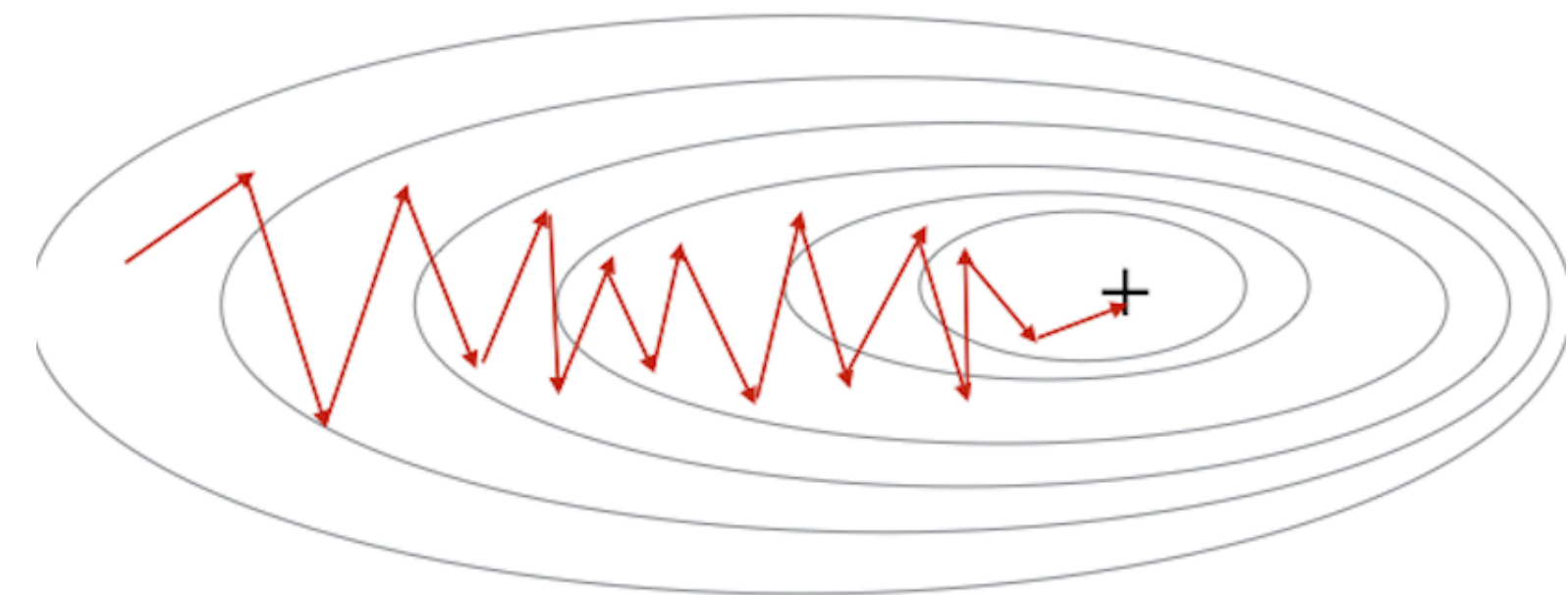
$$\theta_{t+1} = \theta_t - \eta_t \widetilde{\nabla}_\theta J(\theta_t)$$

where $\mathbb{E} \left[\widetilde{\nabla}_\theta J(\theta_t) \right] = \nabla_\theta J(\theta_t)$

Gradient Descent



Stochastic Gradient Descent



SGD: Convergence to a Stationary Point for Nonconvex Functions

- Def of β -smooth: $\|\nabla_{\theta}J(\theta) - \nabla_{\theta}J(\theta_0)\|_2 \leq \beta\|\theta - \theta_0\|_2$

- **[Theorem]** (informal) Suppose we run SGD: $\theta_{t+1} = \theta_t - \eta \widetilde{\nabla}_{\theta}J(\theta_t)$, for T steps, where $\mathbb{E} \left[\widetilde{\nabla}_{\theta}J(\theta_t) \right] = \nabla_{\theta}J(\theta_t)$ with $\eta = O(1/\sqrt{T})$. Assume:
 - $J(\theta)$ is β -smooth.
 - $J(\theta)$ is bounded: $|J(\theta)| \leq M, \quad \forall \theta.$
 - $\widetilde{\nabla}_{\theta}J(\theta)$ has “bounded variance”: $\mathbb{E} \left[\|\nabla_{\theta}J(\theta_t) - \widetilde{\nabla}_{\theta}J(\theta_t)\|_2^2 \right] \leq \sigma^2,$

then, in T steps, SGD will find a θ such that:

$$\|\nabla_{\theta}J(\theta)\| \leq O \left((M\beta\sigma^2/T)^{1/4} + (M\beta/T)^{1/2} \right).$$

Formally, we have $\mathbb{E} \left[\min_{t \leq T} \|\nabla_{\theta}J(\theta_t)\|^2 \right] \leq O \left(\sqrt{M\beta\sigma^2/T} + M\beta/T \right)$

Recap Outline:

1. The Learning Setting
2. Objective: direct policy optimization.
3. General convergence: properties of SGD
4. Importance Sampling
& Deriving a Policy Gradient Expression

Importance Sampling (and the Likelihood Ratio Method)

For $J(\theta) = \mathbb{E}_{x \sim P_\theta} [f(x)]$, our goal is to accurately approximate $\nabla_\theta J(\theta) = \nabla_\theta \mathbb{E}_{x \sim P_\theta} f(x)$.
(We want to avoid computing the integral/sum.)

- Often, we are in setting where:
 - P_θ is “easy” to compute.
 - We have a distribution ρ , that is easy to sample from and where $\max_x P_\theta(x)/\rho(x) < \infty$
(sometimes we use P_θ itself as ρ)

$$\nabla_\theta J(\theta) = \nabla_\theta \mathbb{E}_{x \sim P_\theta} f(x) = \nabla_\theta \mathbb{E}_{x \sim \rho} \frac{P_\theta(x)}{\rho(x)} f(x) = \mathbb{E}_{x \sim \rho} \frac{\nabla_\theta P_\theta(x)}{\rho(x)} f(x) \approx \frac{1}{N} \sum_{i=1}^N \frac{\nabla_\theta P_\theta(x_i)}{\rho(x_i)} f(x_i)$$

To compute gradient at θ_0 : $\nabla_\theta J(\theta_0)$ (in short of $\nabla_\theta J(\theta) |_{\theta=\theta_0}$)

By setting the sampling distribution $\rho = P_{\theta_0}$

$$\nabla_\theta J(\theta_0) = \mathbb{E}_{x \sim P_{\theta_0}} \left[\nabla_\theta \ln P_{\theta_0}(x) \cdot f(x) \right]$$

Recap: the REINFORCE Algorithm (discounted case)

We derived the most basic PG formulation:

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\left(\sum_{h=0}^{\infty} \nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) \right) R(\tau) \right], \text{ where } R(\tau) = \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h)$$

Increase the likelihood of sampling an trajectory with high total reward

Recap: the REINFORCE Algorithm

(finite horizon case)

$$\begin{aligned}\tau &= \{s_0, a_0, s_1, a_1, \dots, s_{H-1}, a_{H-1}\} \\ \rho_\theta(\tau) &= \mu(s_0)\pi_\theta(a_0 | s_0)P(s_1 | s_0, a_0)\pi_\theta(a_1 | s_1)\dots \\ J(\theta) &= \mathbb{E}_{\tau \sim \rho_\theta(\tau)} \underbrace{\left[\sum_{h=0}^{H-1} r(s_h, a_h) \right]}_{R(\tau)}\end{aligned}$$

$$\nabla_\theta J(\theta) := \mathbb{E}_{\tau \sim \rho_\theta(\tau)} \left[\left(\sum_{h=0}^{H-1} \nabla_\theta \ln \pi_\theta(a_h | s_h) \right) R(\tau) \right]$$

Derivation of Policy Gradient: REINFORCE for finite H

$$\tau = \{s_0, a_0, s_1, a_1, \dots\}$$

$$\rho_{\theta}(\tau) = \mu(s_0)\pi_{\theta}(a_0 | s_0)P(s_1 | s_0, a_0)\pi_{\theta}(a_1 | s_1)\dots P(s_{H-1} | s_{H-2}, a_{H-2})\pi_{\theta}(a_{H-1} | s_{H-1})$$

$$J(\theta) = \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \underbrace{\left[\sum_{h=0}^{H-1} r(s_h, a_h) \right]}_{R(\tau)}$$

$$\nabla_{\theta} J(\theta_0) = \mathbb{E}_{\tau \sim \rho_{\theta_0}(\tau)} \left[\nabla_{\theta} \ln \rho_{\theta_0}(\tau) R(\tau) \right]$$

$$= \mathbb{E}_{\tau \sim \rho_{\theta_0}(\tau)} \left[\nabla_{\theta} \left(\ln \rho(s_0) + \ln \pi_{\theta_0}(a_0 | s_0) + \ln P(s_1 | s_0, a_0) + \dots \right) R(\tau) \right]$$

$$= \mathbb{E}_{\tau \sim \rho_{\theta_0}(\tau)} \left[\nabla_{\theta} \left(\ln \pi_{\theta_0}(a_0 | s_0) + \ln \pi_{\theta_0}(a_1 | s_1) \dots \right) R(\tau) \right] = \mathbb{E}_{\tau \sim \rho_{\theta_0}(\tau)} \left[\left(\sum_{h=0}^{H-1} \nabla_{\theta} \ln \pi_{\theta_0}(a_h | s_h) \right) R(\tau) \right]$$

Today:

Policy Gradient Descent

Outline:

1. Estimation of Stochastic Gradients
2. Variance Reduction
3. More Variance Reduction (baselines)

Obtaining an Unbiased Gradient Estimate at θ_0

$$\nabla_{\theta} J(\theta) := \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\left(\sum_{h=0}^{H-1} \nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) \right) R(\tau) \right]$$

1. Obtain a trajectory $\tau \sim \rho_{\theta_0}$
(which we can do in our learning setting)
2. Set:

$$\widetilde{\nabla}_{\theta} J(\theta_0) := \sum_{h=0}^{H-1} \nabla \ln \pi_{\theta_0}(a_h | s_h) R(\tau)$$

We have: $\mathbb{E}[\widetilde{\nabla}_{\theta} J(\theta_0)] = \nabla_{\theta} J(\pi_{\theta_0})$

PG with REINFORCE:

1. Initialize θ_0 , parameters: η_1, η_2, \dots

2. For $t = 0, \dots$:

1. Obtain a trajectory $\tau \sim \rho_{\theta_t}$

$$\text{Set } \widetilde{\nabla}_{\theta} J(\theta_t) = \sum_{h=0}^{H-1} \nabla \ln \pi_{\theta_t}(a_h | s_h) R(\tau)$$

2. Update: $\theta_{t+1} = \theta_t + \eta_t \widetilde{\nabla}_{\theta} J(\theta_t)$

The (mini-batch) PG procedure with REINFORCE

(reducing variance using batch sizes of M)

1. Initialize θ_0 , parameters: η_1, η_2, \dots

2. For $t = 0, \dots$:

1. Init $g = 0$ and do M times:

Obtain a trajectory $\tau \sim \rho_{\theta_t}$

Set $g = g + \sum_{h=0}^{H-1} \nabla \ln \pi_{\theta_t}(a_h | s_h) R(\tau)$

Set $\widetilde{\nabla}_{\theta} J(\theta_t) := \frac{1}{M} g$

2. Update: $\theta_{t+1} = \theta_t + \eta_t \widetilde{\nabla}_{\theta} J(\theta_t)$

Outline:

1. Estimation of Stochastic Gradients
2. Variance Reduction
3. More Variance Reduction: Baselines and Advantages

A improved PG formulation, for sampling (finite horizon setting)

$$\begin{aligned}\nabla J(\theta) &= \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\left(\sum_{h=0}^{H-1} \nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) \right) R(\tau) \right] \\ &= \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\sum_{h=0}^{H-1} \left(\nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) \sum_{t=h}^{H-1} r_t \right) \right] \\ &= \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\sum_{h=0}^{H-1} \nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) Q_h^{\pi_{\theta}}(s_h, a_h) \right]\end{aligned}$$

Intuition: Change action distribution at h only affects rewards later on...)

HW: You will show these simplified version are also valid PG expressions

Proof sketch

Let $f(s_0, a_0, \dots, s_{h-1}, a_{h-1}, s_h)$ be an arbitrary function.

$$\begin{aligned} & \mathbb{E}_{a_h \sim \pi_\theta(\cdot | s_h)} \left[\nabla_\theta \ln \pi_\theta(a_h | s_h) f(s_0, a_0, \dots, s_{h-1}, a_{h-1}, s_h) \mid s_0, a_0, \dots, s_{h-1}, a_{h-1}, s_h \right] \\ &= f(s_0, a_0, \dots, s_{h-1}, a_{h-1}, s_h) \mathbb{E}_{a_h \sim \pi_\theta(\cdot | s_h)} \left[\nabla_\theta \ln \pi_\theta(a_h | s_h) \mid s_0, a_0, \dots, s_{h-1}, a_{h-1}, s_h \right] = ?? \end{aligned}$$

An improved PG procedure:

1. Initialize θ_0 , parameters: η_1, η_2, \dots

2. For $t = 0, \dots$:

1. Obtain a trajectory $\tau \sim \rho_{\theta_t}$

$$\text{Set } \widetilde{\nabla}_{\theta} J(\theta_t) = \sum_{h=0}^{H-1} \left(\nabla \ln \pi_{\theta_t}(a_h | s_h) \sum_{t=h}^{H-1} r_t \right)$$

2. Update: $\theta_{t+1} = \theta_t + \eta_t \widetilde{\nabla}_{\theta} J(\theta_t)$

A improved PG formulation, for sampling (for the discounted setting)

$$\begin{aligned}\nabla J(\theta) &= \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\left(\sum_{h=0}^{\infty} \nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) \right) R(\tau) \right] \\ &= \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\sum_{h=0}^{\infty} \left(\nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) \sum_{t=h}^{\infty} \gamma^t r_t \right) \right] \\ &= \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\sum_{h=0}^{\infty} \nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) \gamma^h Q^{\pi_{\theta}}(s_h, a_h) \right]\end{aligned}$$

Outline:

1. Estimation of Stochastic Gradients
2. Variance Reduction
3. More Variance Reduction: Baselines and Advantages

With a “baseline” function:

For any function $b_h(s)$, we have:

$$\begin{aligned}\nabla J(\theta) &= \mathbb{E}_{\tau \sim \rho_\theta(\tau)} \left[\sum_{h=0}^{H-1} \nabla_\theta \ln \pi_\theta(a_h | s_h) \left(\sum_{t=h}^{H-1} r_t - b_h(s_h) \right) \right] \\ &= \mathbb{E}_{\tau \sim \rho_\theta(\tau)} \left[\sum_{h=0}^{H-1} \nabla_\theta \ln \pi_\theta(a_h | s_h) \left(Q_h^{\pi_\theta}(s_h, a_h) - b_h(s_h) \right) \right]\end{aligned}$$

(M=1) SGD with a Naive (constant) Baseline:

- On a trajectory τ , define:

$$R_h(\tau) = \sum_{t=h}^{H-1} r_t.$$

- Let try to use a constant (time-dependent) baseline:

$$b_h^\theta = \mathbb{E}_{\tau \sim \rho_\theta(\tau)} E [R_h(\tau)]$$

1. Initialize θ_0 , parameters: η_1, η_2, \dots
2. For $t = 0, \dots$:
 1. Using N trajectories sampled under π_{θ_t} , set

$$\tilde{b}_h = \frac{1}{N} \sum_{i=1}^N R_h(\tau_i)$$

2. Obtain a trajectory $\tau \sim \rho_{\theta_t}$

$$\text{Set } \widetilde{\nabla}_\theta J(\theta_t) = \sum_{h=0}^{H-1} \nabla \ln \pi_{\theta_t}(a_h | s_h) \left(R_h(\tau) - \tilde{b}_h \right)$$

3. Update: $\theta_{t+1} = \theta_t + \eta_t \widetilde{\nabla}_\theta J(\theta_t)$

The Advantage Function (finite horizon)

$$V_h^\pi(s) = \mathbb{E} \left[\sum_{\tau=h}^{H-1} r(s_\tau, a_\tau) \mid s_h = s \right] \quad Q_h^\pi(s, a) = \mathbb{E} \left[\sum_{\tau=h}^{H-1} r(s_\tau, a_\tau) \mid (s_h, a_h) = (s, a) \right]$$

- The Advantage function is defined as:

$$A_h^\pi(s, a) = Q_h^\pi(s, a) - V_h^\pi(s, a)$$

- We have that:

$$E_{a \sim \pi(\cdot | s)} [A_h(s, a) \mid s, h] = \sum_a \pi(a \mid s) A_h(s, a) = ??$$

- For the discounted case, $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s, a)$

The Advantage-based PG:

$$\begin{aligned}\nabla J(\theta) &= \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\sum_{h=0}^{H-1} \nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) \left(Q_h^{\pi_{\theta}}(s_h, a_h) - b_h(s_h) \right) \right] \\ &= \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\sum_{h=0}^{H-1} \nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) A_h^{\pi_{\theta}}(s_h, a_h) \right]\end{aligned}$$

- The second step follows by choosing $b_h(s) = V_h^{\pi}(s)$.
- In practice, the most common approach is to use $b_h(s)$ as an estimate of $V_h^{\pi}(s)$.

Summary so far:

Variance reduction with:

- Improvement over REINFORCE
- baseline functions (and the “advantage” formulation)

$$\begin{aligned}\nabla J(\theta) &= \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\sum_{h=0}^{H-1} \nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) \left(\sum_{t=h}^{H-1} r_t - b_h(s_h) \right) \right] \\ &= \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\sum_{h=0}^{H-1} \nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) \left(Q_h^{\pi_{\theta}}(s_h, a_h) - b_h(s_h) \right) \right] \\ &= \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\sum_{h=0}^{H-1} \nabla_{\theta} \ln \pi_{\theta}(a_h | s_h) A_h^{\pi_{\theta}}(s_h, a_h) \right]\end{aligned}$$

1-minute feedback form: <https://bit.ly/3RHtlxy>

