

PG Methods, Baselines, & fitted Value function methods

Lucas Janson and Sham Kakade

**CS/Stat 184: Introduction to Reinforcement Learning
Fall 2022**

Today

- Recap
- Today:
 1. Variance Reduction w/ Baselines
 2. Advantages and a better baseline
 3. An example: PG Example with (softmax) linear policies
 4. Fitted Value Functions:
 1. Direct approach
 2. An iterative approach

Recap

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: Policy Parameterization

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

1. Softmax linear Policy

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

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

Recap: the REINFORCE Algorithm

(finite horizon case)

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

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

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

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

Today:

Policy Gradients: Baselines
& Fitted Value Function Methods

Outline:

1. Variance Reduction w/ Baselines
2. Advantages and a better baseline
3. An example: PG Example with (softmax) linear policies
4. Fitted Value Functions:
 1. Direct approach
 2. An iterative approach

With a “baseline” function:

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

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

This is (basically) the method of control variates.

(M=1) PG 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:

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

(which also depends on θ)

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 $\tau \sim \rho_{\theta_t}$

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

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

Outline:

1. Variance Reduction w/ Baselines
2. Advantages and a better baseline
3. An example: PG Example with (softmax) linear policies
4. Fitted Value Functions:
 1. Direct approach
 2. An iterative approach

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:

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

- We have that:

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

- What do we know about $A_h^{\pi^\star}(s, a)$?
- For the discounted case, $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$

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 $b_h(s)$ to approximate $V_h^{\pi}(s)$.

Outline:

1. Variance Reduction w/ Baselines
2. Advantages and a better baseline
3. An example: PG Example with (softmax) linear policies
4. Fitted Value Functions:
 1. Direct approach
 2. An iterative approach

Policy Parameterizations

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

1. Softmax linear Policy

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

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

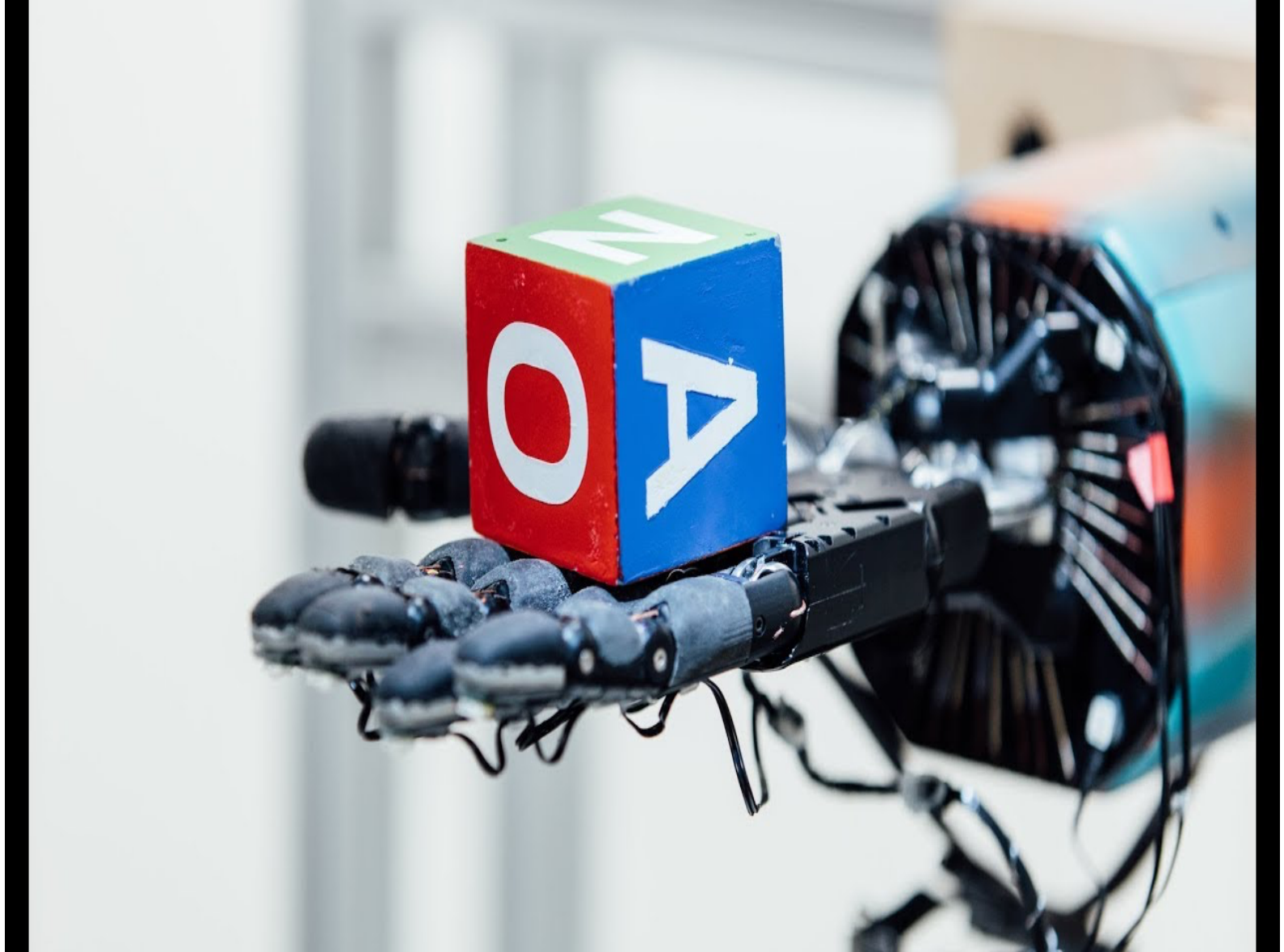
What is a “state” and a “feature vector”?



└AlphaZero. Silver



└OpenAI Five.



└OpenAI.

A state:

- **Tabular case:** an index in $[|S|] = \{1, \dots, |S|\}$
- **Real world:** a list/array of the relevant info about the world that makes the process Markovian.
(we sometimes append history info into the current state)
- Let's assume the current time h is contained in the state.
(e.g. you can always add the time into the “list” that specifies the state)

Softmax Policy Properties

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

1. Softmax linear Policy

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

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

Two properties (see HW):

- More probable actions have features which align with θ .
Precisely,

$$\pi_\theta(a | s) \geq \pi_\theta(a' | s) \text{ if and only if } \theta^\top \phi(s, a) \geq \theta^\top \phi(s, a')$$

- The gradient is:

$$\nabla_\theta \log(\pi_\theta(a | s)) = \phi(s, a) - \mathbb{E}_{a' \sim \pi_\theta(\cdot | s)}[\phi(s, a')]$$

PG for the (softmax) linear policies

- We have:

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \rho_{\theta}(\tau)} \left[\sum_{h=0}^{H-1} A_h^{\pi_{\theta}}(s_h, a_h) \left(\phi(s_h, a_h) - \mathbb{E}_{a' \sim \pi_{\theta}(\cdot | s_h)} [\phi(s_h, a')] \right) \right]$$

(also true Q_h instead of A_h)

- We can simplify this to:

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

- Why?

Outline:

1. Variance Reduction w/ Baselines
2. Advantages and a better baseline
3. An example: PG Example with (softmax) linear policies
4. Fitted Value Functions:
 1. Direct approach
 2. An iterative approach

(M=1) PG with a Learned Baseline:

1. Initialize θ_0 , parameters: η_1, η_2, \dots
2. For $t = 0, \dots$: **Now let's look at our baseline fitting step.**

1. Using N trajectories sampled under π_{θ_t} , try to learn a \tilde{b}_h

$$\tilde{b}(s) \approx V_h^{\pi_{\theta_t}}(s)$$

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

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

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

Note that regardless of our choice of $\tilde{b}_h(s)$, we still get unbiased gradient estimates.

Baseline/Value Function Parameterizations

Now let us consider parameterized classes of functions \mathcal{F} , where for each $f \in \mathcal{F}$, $f: S \rightarrow R$

1. Linear Functions

Feature vector $\psi(s) \in \mathbb{R}^k$, and
parameter $w \in \mathbb{R}^k$

$$f_w(s) = w^\top \psi(s)$$

2. Neural Policy:

Neural network $f_w: S \mapsto \mathbb{R}$

“Review”

- For a random variable $y \in R$, what is:

$$\arg \min_c E_{y \sim D}[(c - y)^2] = ??$$

- Now let us look at the “function” case where we have a distribution over (x, y) pairs

$$f^\star = \arg \min_{f \in \mathcal{F}} E_{(x,y) \sim D}[(f(x) - y)^2]$$

(where \mathcal{F} is the class of all possible functions)

What is $f^\star(x) = ??$

Let's look at our fitting step

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

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

1. Sample N trajectories under π_{θ_t} to make a dataset,

$$\widetilde{w} = \arg \min_w \sum_{\tau \in \text{Data}} \sum_{(s_h, a_h) \in \tau} \left(f_w(s_h) - R_h(\tau) \right)^2$$

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

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

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

Outline:

1. Variance Reduction w/ Baselines
2. Advantages and a better baseline
3. An example: PG Example with (softmax) linear policies
4. Fitted Value Functions:
 1. Direct approach
 2. An iterative approach

Is there an iterative version of Policy Evaluation?

(that is faster, but approximate?)

Algorithm (Iterative PE):

1. Initialization: $V^0 : \|V^0\|_\infty \in \left[0, \frac{1}{1-\gamma}\right]$
2. Iterate until convergence: $V^{t+1} \leftarrow R + \gamma P V^t$

This is a “fixed point” algorithm trying to enforce Bellman consistency:

$$\forall s, V^\pi(s) = r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi(s))} V^\pi(s')$$

Let's look at our fitting step

1. Initialize θ_0 , parameters: η_1, η_2, \dots
2. For $t = 0, \dots$:
 1. Using N trajectories sampled under π_{θ_t} , try to learn a \tilde{b}_h
 $\tilde{b}_h(s) \approx V_h^{\pi_{\theta_t}}(s)$
 2. Obtain a trajectory $\tau \sim \rho_{\theta_t}$
Set $\widetilde{\nabla}_{\theta} J(\theta_t) = \sum_{h=0}^{H-1} \nabla \ln \pi_{\theta_t}(a_h | s_h) (R_h(\tau) - \tilde{b}_h)$
3. Update: $\theta_{t+1} = \theta_t + \eta_t \widetilde{\nabla}_{\theta} J(\theta_t)$

Let's look at just our fitting step in the inner loop (where we want to fit the value of π_{θ_t})

1. Sample N trajectories under π_{θ_t} to make a dataset.
2. Initialize w_0
3. For $k = 0, \dots, K$:
 1. Update:

$$w_{k+1} = \arg \min_w \sum_{\tau \in \text{Data}} \sum_{(s_h, a_h) \in \tau} \left(f_w(s_h) - \left(r_h + f_{w_k}(s_{h+1}) \right) \right)^2$$

Temporal Difference Learning (TD) is a an online method to do the above.

Summary so far:

1. Variance Reduction w/ Baselines & Advantages.
2. An example: PG Example with (softmax) linear policies
3. Fitted Value Functions:
 1. Direct approach
 2. An iterative approach & TD

Next up: Why not just directly use a “fitted” approach for Value Iteration or Policy Iteration?

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

