

POLOCY GRADIENT
INTRODUCTION TO REINFORCEMENT LEARNING

Bogdan Ivanyuk-Skulskyi, Dmytro Kuzmenko

Department of Mathematics,
National University of Kyiv-Mohyla Academy

March 20, 2023

RL ALGORITHMS INVOLVE

- ▶ Optimization
- ▶ Delayed consequences
- ▶ Exploration
- ▶ Generalization
- ▶ And statistical and computational efficiency matters

LAST TIME: GENERALIZATION AND EFFICIENCY

Can use structure and additional knowledge to help constrain and speed reinforcement learning

POLICY-BASED REINFORCEMENT LEARNING

- ▶ In the last lecture we approximated the value or action-value function using parameters w ,

$$V_w(s) \approx V^\pi(s)$$

$$Q_w(s, a) \approx Q^\pi(s, a)$$

- ▶ A policy was generated directly from the value function (using ϵ -greedy)
- ▶ In this lecture we will directly parametrize the policy, and will typically use θ to show parameterization:

$$\pi_\theta(s, a) = \mathbb{P}[a|s; \theta]$$

- ▶ Goal is to find a policy π with the highest value function V^π
- ▶ We will focus again on model-free reinforcement learning

VALUE-BASED AND POLICY-BASED RL

- ▶ Value Based
 - learned Value Function
 - Implicit policy (e.g. -greedy)
- ▶ Policy Based
 - No Value Function
 - Learned Policy
- ▶ Actor-Critic
 - Learned Value Function
 - Learned Policy

TYPES OF POLICIES TO SEARCH OVER

- ▶ So far have focused on deterministic policies or ϵ -greedy policies
- ▶ Now we are thinking about direct policy search in RL, will focus heavily on stochastic policies

POLICY OBJECTIVE FUNCTIONS

- ▶ Goal: given a policy $\pi_\theta(s, a)$ with parameters θ , find best θ
- ▶ But how do we measure the quality for a policy π_θ ?
- ▶ In episodic environments can use policy value at start state $V(s_0, \theta)$
- ▶ For simplicity, today will mostly discuss the episodic case, but can easily extend to the continuing / infinite horizon case

POLICY OPTIMIZATION

- ▶ Policy based reinforcement learning is an optimization problem
- ▶ Find policy parameters θ that maximize $V(s_0, \theta)$

POLICY OPTIMIZATION

- ▶ Policy based reinforcement learning is an optimization problem
- ▶ Find policy parameters θ that maximize $V(s_0, \theta)$
- ▶ Can use gradient free optimization
 - Hill climbing
 - Simplex / amoeba / Nelder Mead
 - Genetic algorithms
 - Cross-Entropy method (CEM)
 - Covariance Matrix Adaptation (CMA)

HUMAN-IN-THE-LOOP EXOSKELETON OPTIMIZATION (ZHANG ET AL. SCIENCE 2017)

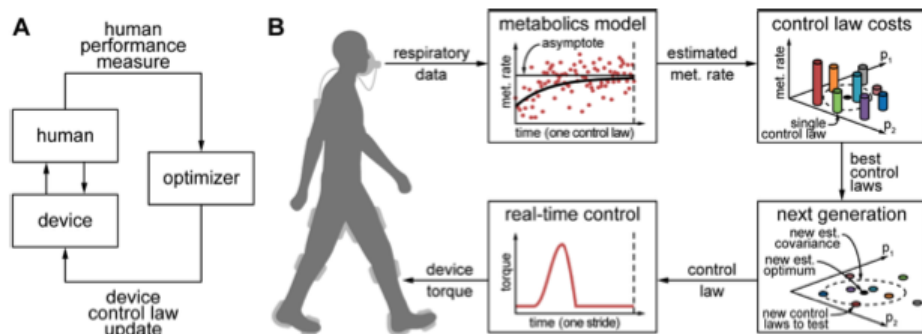


Figure: Zhang et al. Science 2017

Optimization was done using CMA-ES, variation of covariance matrix evaluation

GRADIENT FREE POLICY OPTIMIZATION

Can often work embarrassingly well: "discovered that evolution strategies (ES), an optimization technique that's been known for decades, rivals the performance of standard reinforcement learning (RL) techniques on modern RL benchmarks (e.g. Atari/MuJoCo)" (<https://blog.openai.com/evolution-strategies/>)

GRADIENT FREE POLICY OPTIMIZATION

- ▶ Often a great simple baseline to try
- ▶ Benefits
 - Can work with any policy parametrizations, including non-differentiable
 - Frequently very easy to parallelize
- ▶ Limitations
 - Often less sample efficient because it ignores temporal structure

POLICY OPTIMIZATION

- ▶ Policy based reinforcement learning is an optimization problem
- ▶ Find policy parameters θ that maximize $V(s_0, \theta)$
- ▶ Can use gradient free optimization:
- ▶ Greater efficiency often possible using gradient
 - Gradient descent
 - Conjugate gradient
 - Quasi-newton
- ▶ We focus on gradient descent, many extensions possible
- ▶ And on methods that exploit sequential structure

POLICY GRADIENT

- ▶ Define $V(\theta) = V(s_0, \theta)$ to make explicit the dependence of the value on the policy parameters [but don't confuse with value function approximation, where parameterized value function]
- ▶ Assume episodic MDPs (easy to extend to related objectives, like average reward)

POLICY GRADIENT

- ▶ Define $V^{\pi_{\theta}} = V(s_0, \theta)$ to make explicit the dependence of the value on the policy parameters
- ▶ Assume episodic MDPs
- ▶ Policy gradient algorithms search for a local maximum in $V(s_0, \theta)$ by ascending the gradient of the policy, w.r.t parameters

$$\Delta\theta = \alpha \nabla_{\theta} V(s_0, \theta)$$

- ▶ Where $\nabla_{\theta} V(s_0, \theta)$ is the policy gradient

$$\nabla_{\theta} V(s_0, \theta) = \begin{pmatrix} \frac{\partial V(s_0, \theta)}{\partial \theta_1} \\ \frac{\partial V(s_0, \theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial V(s_0, \theta)}{\partial \theta_n} \end{pmatrix}$$

- ▶ and α is a step-size parameter

SUMMARY OF BENEFITS OF POLICY-BASED RL

Advantages:

- ▶ Better convergence properties
- ▶ Effective in high-dimensional or continuous action spaces
- ▶ Can learn stochastic policies

Disadvantages:

- ▶ Typically converge to a local rather than global optimum
- ▶ Evaluating a policy is typically inefficient and high variance

Shortly will see some ideas to help with this last limitation

TABLE OF CONTENTS

- ▶ **Differentiable Policies**
- ▶ Temporal Structure
- ▶ Baseline
- ▶ Alternatives to MC Returns

COMPUTING THE GRADIENT ANALYTICALLY

- ▶ We now compute the policy gradient analytically
- ▶ Assume policy π_θ is differentiable whenever it is non-zero
- ▶ Assume we can calculate gradient $\nabla_\theta \pi_\theta(s, a)$ analytically
- ▶ What kinds of policy classes can we do this for?

DIFFERENTIABLE POLICY CLASSES

- ▶ Many choices of differentiable policy classes including:
 - Softmax
 - Gaussian
 - Neural networks

VALUE OF A PARAMETERIZED POLICY

- ▶ Now assume policy π_θ is differentiable whenever it is non-zero and we know the gradient $\nabla_\theta \pi_\theta(s, a)$
- ▶ Recall policy value is $V(s_0, \theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T R(s_t, a_t); \pi_\theta, s_0 \right]$ where the expectation is taken over the states actions visited by π_θ
- ▶ We can re-express this in multiple ways
 - $V(s_0, \theta) = \sum_a \pi_\theta(a|s_0) Q(s_0, a, \theta)$

VALUE OF A PARAMETERIZED POLICY

- ▶ Now assume policy π_θ is differentiable whenever it is non-zero and we know the gradient $\nabla_\theta \pi_\theta(s, a)$
- ▶ Recall policy value is $V(s_0, \theta) = \sum_a \pi_\theta(a|s_0)Q(s_0, a, \theta)$ where the expectation is taken over the states actions visited by π_θ
- ▶ We can re-express this in multiple ways
 - $V(s_0, \theta) = \sum_a \pi_\theta(a|s_0)Q(s_0, a, \theta)$
 - $V(s_0, \theta) = \sum_\tau P(\tau; \theta)R(\tau)$
 - ▶ where $\tau = (s_0, a_0, r_0, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T)$ is a state-action trajectory,
 - ▶ $P(\tau; \theta)$ is used to denote the probability over trajectories when executing policy $\pi(\theta)$ starting in state s_0 , and
 - ▶ $R(\tau) = \sum_{i=0}^T R(s_i, a_i)$ the sum of rewards for a trajectory τ
- ▶ To start will focus on this latter definition. See Chp 13.1-13.3 of SB for a nice discussion starting with the other definition

LIKELIHOOD RATIO POLICIES

- ▶ Denote a state-action trajectory as $\tau = (s_0, a_0, r_0, \dots, s_{T1}, a_{T1}, r_{T1}, s_T)$
- ▶ Use $R(\tau) = \sum_{t=0}^T R(s_t, a_t)$ to be the sum of rewards for a trajectory τ
- ▶ Policy value is

$$V(\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T R(s_t, a_t); \pi_\theta \right] = \sum_{\tau} P(\tau; \theta) R(\tau)$$

- ▶ where $P(\tau; \theta)$ is used to denote the probability over trajectories when executing policy $\pi(\theta)$
- ▶ In this new notation, our goal is to find the policy parameters θ :

$$\operatorname{argmax}_{\theta} V(\theta) = \operatorname{argmax}_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

LIKELIHOOD RATIO POLICY GRADIENT

- Goal is to find the policy parameters θ :

$$\operatorname{argmax}_{\theta} V(\theta) = \operatorname{argmax}_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

- Take the gradient with respect to θ :

$$\begin{aligned} \nabla_{\theta} V(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) R(\tau) \frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)} \\ &= \sum_{\tau} P(\tau; \theta) R(\tau) \nabla_{\theta} \log P(\tau; \theta) \end{aligned}$$

LIKELIHOOD RATIO POLICY GRADIENT

- Goal is to find the policy parameters θ :

$$\operatorname{argmax}_{\theta} V(\theta) = \operatorname{argmax}_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

- Take the gradient with respect to θ :

$$\nabla_{\theta} V(\theta) = \sum_{\tau} P(\tau; \theta) R(\tau) \nabla_{\theta} \log P(\tau; \theta)$$

- Approximate with empirical estimate for m sample trajectories under policy π_{θ} :

$$\nabla_{\theta} V(\theta) \approx \hat{g} = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \nabla_{\theta} \log P(\tau^{(i)}; \theta)$$

DECOMPOSING THE TRAJECTORIES INTO STATES AND ACTIONS

- Approximate with empirical estimate for m sample paths under policy :

$$\nabla_{\theta} V(\theta) \approx \hat{g} = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \nabla_{\theta} \log P(\tau^{(i)}; \theta)$$

$$\begin{aligned} \nabla_{\theta} \log P(\tau^{(i)}; \theta) &= \nabla_{\theta} \log \left[\mu(s_0) \prod_{t=0}^{T-1} \pi_{\theta}(a_t | s_t) P(s_{t+1} | s_t, a_t) \right] \\ &= \nabla_{\theta} \left[\log \mu(s_0) + \sum_{t=0}^{T-1} \log \pi_{\theta}(a_t | s_t) + \log P(s_{t+1} | s_t, a_t) \right] \\ &= \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \quad (\text{score function}) \end{aligned}$$

SCORE FUNCTION

- ▶ A score function is the derivative of the log of a parameterized probability / likelihood
- ▶ Example: let $\pi(s, \theta)$ be the probability of state s under parameter
- ▶ Then the score function would be

$$\nabla_{\theta} \log \pi(s; \theta)$$

- ▶ For many policy classes, it is not hard to compute the score function

SOFTMAX POLICY

- ▶ Weight actions using linear combination of features $\phi(s, a)^T \theta$
- ▶ Probability of action is proportional to exponentiated weight

$$\pi_{\theta}(s, a) = e^{\phi(s, a)^T \theta} / \left(\sum_a e^{\phi(s, a)^T \theta} \right)$$

- ▶ The score function is

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

GAUSSIAN POLICY

- ▶ In continuous action spaces, a Gaussian policy is natural
- ▶ Mean is a linear combination of state features $\mu(s) = \phi(s)^T \theta$
- ▶ Variance may be fixed σ^2 , or can also be parametrised
- ▶ Policy is Gaussian $a \sim N(\mu(s), \sigma^2)$
- ▶ The score function is

$$\nabla_{\theta} \log \pi_{\theta}(s, a) = \frac{(a - \mu(s))\phi(s)}{\sigma^2}$$

LIKELIHOOD RATIO / SCORE FUNCTION POLICY GRADIENT

- Goal is to find the policy parameters :

$$\operatorname{argmax}_{\theta} V(\theta) = \operatorname{argmax}_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

- Approximate with empirical estimate for m sample paths under policy π_{θ} using score function:

$$\begin{aligned} \nabla_{\theta} V(\theta) &\approx \hat{g} = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \nabla_{\theta} \log P(\tau^{(i)}; \theta) \\ &= (1/m) \sum_{i=1}^m R(\tau^{(i)}) \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \end{aligned}$$

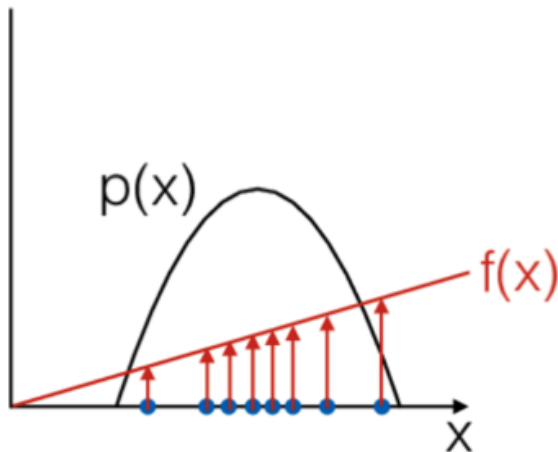
- Do not need to know dynamics model

SCORE FUNCTION GRADIENT ESTIMATOR: INTUITION

- ▶ Consider generic form of $R(\tau^{(i)})\nabla_{\theta} \log P(\tau^{(i)}; \theta)$:
 $\hat{g}_i = f(x_i)\nabla_{\theta} \log p(x_i|\theta)$
- ▶ $f(x)$ measures how good the sample x is.
- ▶ Moving in the direction \hat{g}^i pushes up the logprob of the sample, in proportion to how good it is
- ▶ Valid even if $f(x)$ is discontinuous, and unknown, or sample space (containing x) is a discrete set

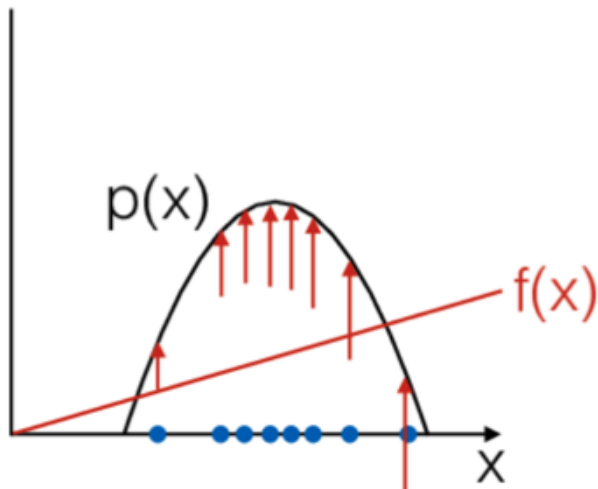
SCORE FUNCTION GRADIENT ESTIMATOR: INTUITION

$$\hat{g}_i = f(x_i) \nabla_{\theta} \log p(x_i | \theta)$$



SCORE FUNCTION GRADIENT ESTIMATOR: INTUITION

$$\hat{g}_i = f(x_i) \nabla_{\theta} \log p(x_i | \theta)$$



POLICY GRADIENT THEOREM

The policy gradient theorem generalizes the likelihood ratio approach

Theorem 1

For any differentiable policy $\pi_\theta(s, a)$, for any of the policy objective function $J = J_1$, (episodic reward), J_{avR} (average reward per time step), or $\frac{1}{J_{avV}}$ (average value), the policy gradient is

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q^{\pi_\theta}(s, a)]$$

Chapter 13.2 in SB has a nice derivation of the policy gradient theorem for episodic tasks and discrete states

TABLE OF CONTENTS

- ▶ Differentiable Policies
- ▶ **Temporal Structure**
- ▶ Baseline
- ▶ Alternatives to MC Returns

LIKELIHOOD RATIO / SCORE FUNCTION POLICY GRADIENT

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

- ▶ Unbiased but very noisy
- ▶ Fixes that can make it practical
 - Temporal structure
 - Baseline

POLICY GRADIENT: USE TEMPORAL STRUCTURE

- Previously:

$$\nabla_{\theta} \mathbb{E}_{\tau}[R] = \mathbb{E}_{\tau} \left[\left(\sum_{t=0}^{T-1} r_t \right) \left(\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \right]$$

- We can repeat the same argument to derive the gradient estimator for a single reward term $r_{t'}$.

$$\nabla_{\theta} \mathbb{E}[r_{t'}] = \mathbb{E} \left[r_{t'} \sum_{t=0}^{t'} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- To see this, recall $V(s_0, \theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^T R(s_t, a_t); \pi_{\theta}, s_0 \right]$ expectation is taken over the states actions visited by π_{θ}

POLICY GRADIENT: USE TEMPORAL STRUCTURE

- Previously:

$$\nabla_{\theta} \mathbb{E}_{\tau}[R] = \mathbb{E}_{\tau} \left[\left(\sum_{t=0}^{T-1} r_t \right) \left(\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \right]$$

- We can repeat the same argument to derive the gradient estimator for a single reward term $r_{t'}$.

$$\nabla_{\theta} \mathbb{E}[r_{t'}] = \mathbb{E} \left[r_{t'} \sum_{t=0}^{t'} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- Summing this formula over t , we obtain

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

POLICY GRADIENT: USE TEMPORAL STRUCTURE

Recall for a particular trajectory $\tau^{(i)}$, $\sum_{t'=t}^{T-1} r_{t'}^{(i)}$ is the return $G_t^{(i)}$

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

Monte-Carlo Policy Gradient (REINFORCE)

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 θ 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 θ

LIKELIHOOD RATIO / SCORE FUNCTION POLICY GRADIENT

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

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

DESIRED PROPERTIES OF A POLICY GRADIENT RL ALGORITHM

Goal: Converge as quickly as possible to a local optima

- ▶ Incurring reward / cost as execute policy, so want to minimize number of iterations / time steps until reach a good policy

TABLE OF CONTENTS

- ▶ Differentiable Policies
- ▶ Temporal Structure
- ▶ **Baseline**
- ▶ Alternatives to MC Returns

POLICY GRADIENT: INTRODUCE BASELINE

- ▶ 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

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

"VANILLA" POLICY GRADIENT ALGORITHM

Initialize policy parameter θ , 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 τ^i , compute
 Advantage estimate \hat{A}_{it}^n
 Update the policy, using a policy gradient estimate \hat{g} ,
 Which is a sum of terms $\nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_{it}^n$.
 (Plug \hat{g} into SGD or ADAM)
endfor

OTHER CHOICES FOR BASELINE?

Initialize policy parameter θ , baseline b

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

Collect a set of trajectories by executing the current policy

At each timestep t in each trajectory τ^i , compute

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

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

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

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

CHOOSING THE BASELINE: VALUE FUNCTIONS

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

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

- ▶ State-value function can serve as a great baseline

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

TABLE OF CONTENTS

- ▶ Differentiable Policies
- ▶ Temporal Structure
- ▶ Baseline
- ▶ Alternatives to MC Returns

LIKELIHOOD RATIO / SCORE FUNCTION POLICY GRADIENT

- Policy gradient:

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

- 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 θ ?

CHOOSING THE TARGET

- ▶ 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 in we saw for TD vs MC, and value function approximation

ACTOR-CRITIC METHODS

- ▶ 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

POLICY GRADIENT FORMULAS WITH VALUE FUNCTIONS

- Recall:

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

- 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 $A^{\pi}(s_t, a_t) = Q^{\pi}(s, a) - V^{\pi}(s)$

CHOOSING THE TARGET: N-STEP ESTIMATORS

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

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

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

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

...

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

- If subtract baselines from the above, get advantage estimators

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

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