

Policy Gradient

Last Time

- Bandits

Exploration
vs Exploitation

ϵ -greedy

UCB

$$\mu_a + c \sqrt{\frac{\log N}{N_a}}$$

Guiding Questions

Guiding Questions

- What is Policy Optimization?
- What is Policy Gradient?

Guiding Questions

- What is Policy Optimization?
- What is Policy Gradient?
- What tricks are needed for it to work effectively?

Map

Map

Challenges in RL

- Exploration and Exploitation
- Credit Assignment
- Generalization

Map

Challenges in RL

- Exploration and Exploitation
- Credit Assignment 
- Generalization

Policy Optimization

Policy Optimization

$$\underset{\pi}{\text{maximize}} \underset{s \sim b}{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s, a_t = \pi(s_t) \right]$$

Policy Optimization

$$\underset{\pi}{\text{maximize}} \underset{s \sim b}{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s, a_t = \pi(s_t) \right]$$

$$\underset{\pi}{\text{maximize}} U(\pi) = \underset{s \sim b}{E} [U^{\pi}(s)]$$

Policy Optimization

$$\underset{\pi}{\text{maximize}} \underset{s \sim b}{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s, a_t = \pi(s_t) \right]$$

$$\underset{\pi}{\text{maximize}} U(\pi) = \underset{s \sim b}{E} [U^{\pi}(s)]$$

Two approximations:

Policy Optimization

$$\underset{\pi}{\text{maximize}} E_{s \sim b} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s, a_t = \pi(s_t) \right]$$

$$\underset{\pi}{\text{maximize}} U(\pi) = E_{s \sim b} [U^{\pi}(s)]$$

$$a = \pi(s)$$

Two approximations:

1. Parameterized stochastic policies

$$\underset{\theta}{\text{maximize}} \quad \underline{U(\pi_{\theta})} = \underline{U(\theta)}$$

$$\underline{a \sim \pi_{\theta}(a \mid s)}$$

Policy Optimization

$$\underset{\pi}{\text{maximize}} \underset{s \sim b}{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s, a_t = \pi(s_t) \right]$$

$$\underset{\pi}{\text{maximize}} U(\pi) = \underset{s \sim b}{E} [U^{\pi}(s)]$$

Two approximations:

1. Parameterized stochastic policies

$$\underset{\theta}{\text{maximize}} \quad U(\pi_{\theta}) = U(\theta) \quad a \sim \pi_{\theta}(a \mid s)$$

2. Monte Carlo Utility

$$U(\pi) \approx \underbrace{\frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})}_{\text{return}}$$

trajectory:

$$\tau = \underbrace{(s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_d, a_d, r_d)}$$

Policy Optimization

$$\underset{\pi}{\text{maximize}} \underset{s \sim b}{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s, a_t = \pi(s_t) \right]$$

$$\underset{\pi}{\text{maximize}} U(\pi) = \underset{s \sim b}{E} [U^{\pi}(s)]$$

Two approximations:

1. Parameterized stochastic policies $\underset{\theta}{\text{maximize}} \quad U(\pi_{\theta}) = U(\theta) \quad a \sim \pi_{\theta}(a \mid s)$
2. Monte Carlo Utility $U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$ trajectory:
 $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_d, a_d, r_d)$

Two classes of optimization algorithms:

Policy Optimization

$$\underset{\pi}{\text{maximize}} \underset{s \sim b}{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid \underline{s_0 = s}, a_t = \pi(s_t) \right]$$

$$\underset{\pi}{\text{maximize}} U(\pi) = \underset{s \sim b}{E} [U^\pi(s)]$$

Two approximations:

1. Parameterized stochastic policies $\underset{\theta}{\text{maximize}} \quad U(\pi_\theta) = U(\theta) \quad a \sim \pi_\theta(a \mid s)$
2. Monte Carlo Utility $U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$ trajectory:
 $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_d, a_d, r_d)$

Two classes of optimization algorithms:

1. Zeroth order (use only $U(\theta)$)
2. First order (use $U(\theta)$ and $\nabla_\theta U(\theta)$)

1. Zeroth-Order Optimization

1. Zeroth-Order Optimization

Common zeroth-order approaches:

1. Genetic Algorithms
2. Pattern Search
3. Cross-Entropy

1. Zeroth-Order Optimization

Common zeroth-order approaches:

1. Genetic Algorithms
2. Pattern Search
3. Cross-Entropy

Cross Entropy:

Initialize d

loop:

population \leftarrow sample(d)

elite $\leftarrow m$ with highest $U(\theta)$

$d \leftarrow$ fit(elite)

1. Zeroth-Order Optimization

Common zeroth-order approaches:

1. Genetic Algorithms
2. Pattern Search
3. Cross-Entropy

Cross Entropy:

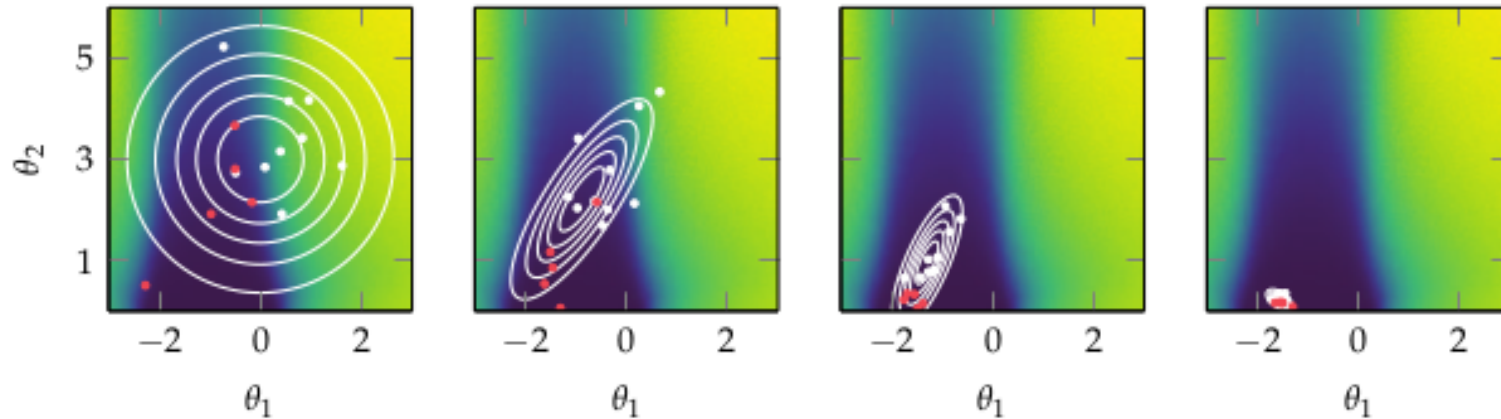
Initialize d

loop:

population \leftarrow sample(d)

elite $\leftarrow m$ with highest $U(\theta)$

$d \leftarrow \text{fit}(\text{elite})$



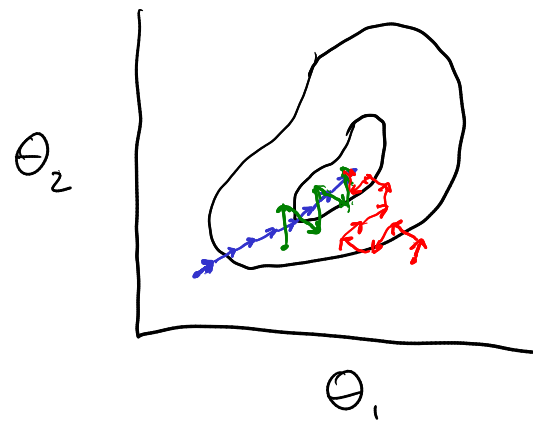
2. First Order Optimization

$$\nabla_{\theta} U(\theta) = \left[\frac{\partial}{\partial \theta_1} U \Big|_{\theta}, \frac{\partial}{\partial \theta_2} U \Big|_{\theta}, \dots, \frac{\partial}{\partial \theta_n} U \Big|_{\theta} \right]$$

loop

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} U(\theta)$$

$$\nabla_{\theta} U(\theta) = E[\nabla_{\theta} U(\theta)]$$



Roughly:

Convergence to local optimum $\iff \sum_{k=1}^{\infty} \alpha^{(k)} = \infty, \sum_{k=1}^{\infty} (\alpha^{(k)})^2 < \infty$

Thm 3

Patterns, Pred
and Actions
Hardt + Recht

$$\rho_0 = \|x_0 - x^*\|, \quad \|\nabla f\| \leq B$$

Suppose we run SGD on convex function $f(x)$ with minimum f^* for N steps with step size α

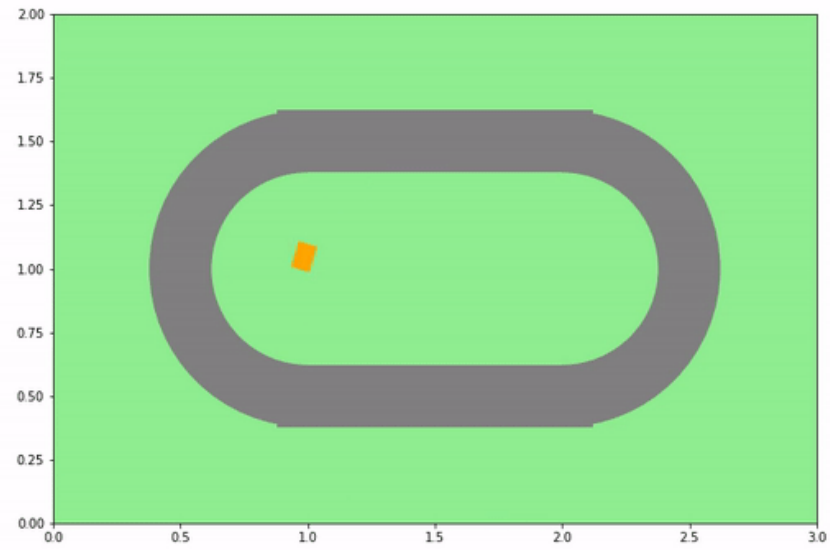
$$\alpha_{\text{opt}} = \frac{\rho_0}{B\sqrt{N}}, \quad \frac{\alpha}{\alpha_{\text{opt}}} = \theta$$

$$\text{then } E[f(x_N) - f^*] \leq \left(\frac{1}{2}\theta + \frac{1}{2}\theta^{-1} \right) \frac{B\rho_0}{\sqrt{N}}$$

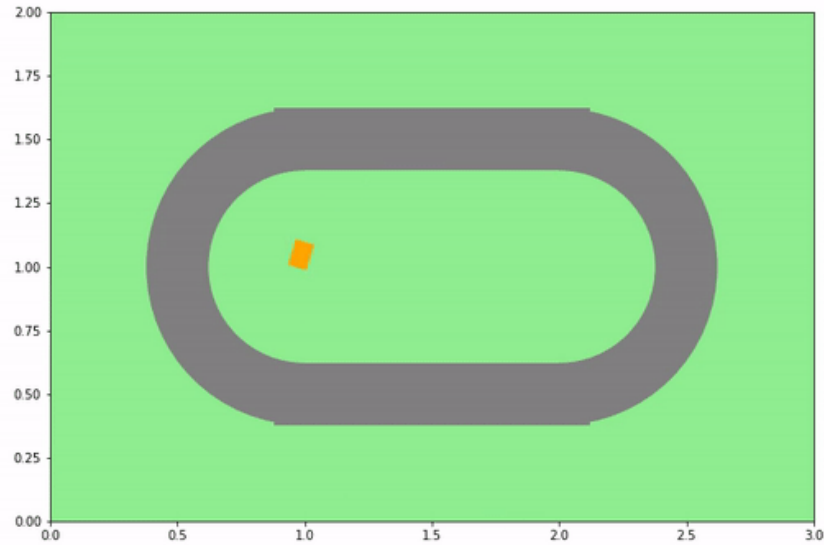
- Definition of Gradient
- Gradient Ascent
- Stochastic Gradient Ascent

Tricks

Tricks



Tricks



For policy gradient, 3 tricks

- Likelihood Ratio/Log Derivative
- Reward to go
- Baseline Subtraction

Log Derivative

$$U(\theta) = E[R(\tau)] \\ = \int p_{\theta}(\tau) R(\tau) d\tau$$

$$\nabla_{\theta} U(\theta) = \nabla_{\theta} \int p_{\theta}(\tau) R(\tau) d\tau \\ = \int \nabla_{\theta} p_{\theta}(\tau) R(\tau) d\tau$$

$$= \int \underbrace{p_{\theta}(\tau)}_{\text{cancel}} \underbrace{\nabla_{\theta} \log p_{\theta}(\tau)}_{\text{cancel}} R(\tau) d\tau$$

$$\nabla_{\theta} U(\theta) = E \left[\underbrace{\nabla_{\theta} \log p_{\theta}(\tau)}_{\nabla_{\theta} U(\theta)} R(\tau) \right]$$

$$\nabla_{\theta} \log p_{\theta}(\tau) = \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} \\ \therefore \nabla_{\theta} p_{\theta}(\tau) = p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)$$



Trajectory Probability Gradient

$$\nabla_{\theta} \log p_{\theta}(\tau)$$

$$\log(ab) = \log a + \log b$$

$$p_{\theta}(\tau) = b(s_0) \prod_{k=0}^d T(s_{k+1} | s_k, a_k) \pi_{\theta}(a_k | s_k)$$

$$\log p_{\theta}(\tau) = \log b(s_0) + \sum_{k=0}^d \log T(s_{k+1} | s_k, a_k) + \sum_{k=0}^d \pi_{\theta}(a_k | s_k)$$

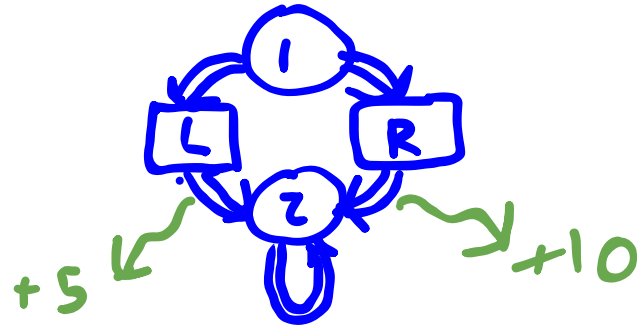
$$\nabla_{\theta} \log p_{\theta}(\tau) = \sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k | s_k)$$

$$\nabla_{\theta} U(\theta) = E \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) R(\tau) \right]$$

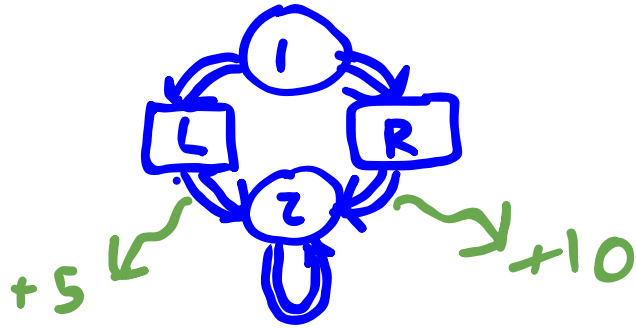
$\nabla_{\theta} U(\theta)$

$A = \{L, R\}$

Example



$$A = \{L, R\}$$



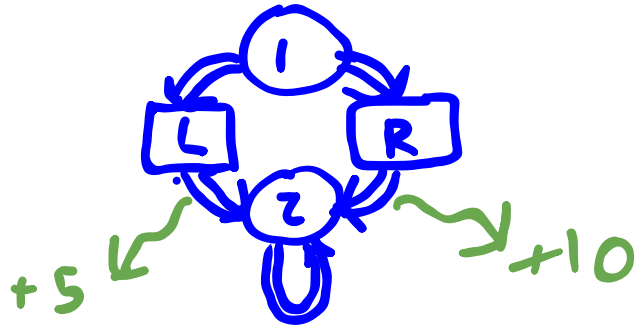
Example

θ

$$\pi_{\theta}(a = L \mid s = 1) = \text{clamp}(\theta, 0, 1)$$

$$\pi_{\theta}(a = R \mid s = 1) = \text{clamp}(1 - \theta, 0, 1)$$

$$A = \{L, R\}$$



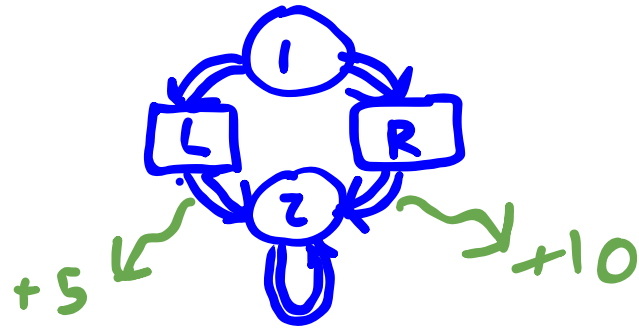
Example

$$\pi_{\theta}(a = L \mid s = 1) = \text{clamp}(\theta, 0, 1)$$

$$\pi_{\theta}(a = R \mid s = 1) = \text{clamp}(1 - \theta, 0, 1)$$

$$\nabla U(\theta) = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) R(\tau) \right]$$

$$A = \{L, R\}$$



Example

$$\pi_{\theta}(a = L \mid s = 1) = \text{clamp}(\theta, 0, 1) = \min(1, \max(0, \theta))$$

$$\pi_{\theta}(a = R \mid s = 1) = \text{clamp}(1 - \theta, 0, 1)$$

$$a) \quad \nabla_{\theta} \log \pi_{\theta}(a|s) = \frac{\partial}{\partial \theta} \log \text{clamp}(\theta, 0, 1) \Big|_{\theta=0.2} = \frac{1}{\theta} = \frac{1}{0.2}$$

$$\widehat{\nabla_{\theta} U(\theta)} = \frac{1}{0.2} 5 = \underline{\underline{25}}$$

$$b) \quad \nabla_{\theta} \log \pi_{\theta}(a=R|s=1) = \frac{\partial}{\partial \theta} \log \text{clamp}(1-\theta, 0, 1) \Big|_{\theta=0.2} = \frac{1}{1-\theta} (-1) = -\frac{1}{0.8}$$

$$\widehat{\nabla_{\theta} U(\theta)} = -\frac{1}{0.8} R(\tau) = \underline{\underline{-12.5}}$$

$$\nabla U(\theta) = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) R(\tau) \right]$$

$$\tau_{(a)} (s_0, a_0=L, r_0=5, s_1=2)$$

$$\mathbb{E}_{a \sim \pi_{\theta}} [\nabla_{\theta} U(\theta)] < 0$$

Given $\theta = 0.2$ calculate $\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) R(\tau)$ for two cases, (a) where $a_0 = L$ and (b) where $a_0 = R$

Policy Gradient

.

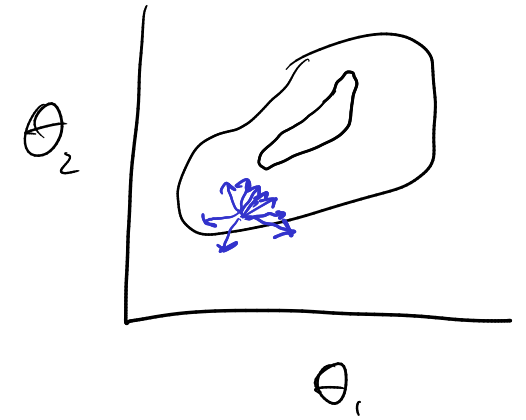
Policy Gradient

loop

$\tau \leftarrow \text{simulate}(\pi_\theta)$

$\theta \leftarrow \theta + \alpha \sum_{k=0}^d \nabla_\theta \log \pi_\theta(a_k | s_k) R(\tau)$

On Policy



Policy Gradient

loop

$\tau \leftarrow \text{simulate}(\pi_\theta)$

$\theta \leftarrow \theta + \alpha \sum_{k=0}^d \underbrace{\nabla_\theta \log \pi_\theta(a_k | s_k)}_{\text{On Policy!}} R(\tau)$

On Policy!

Causality

Causality

$$\nabla U(\theta) = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) R(\tau) \right]$$

Causality

$$\begin{aligned}\nabla U(\theta) &= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) R(\tau) \right] \\ &= \mathbb{E} \left[\left(\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \right) \left(\sum_{k=0}^d \gamma^k r_k \right) \right]\end{aligned}$$

Causality

$$\begin{aligned}\nabla U(\theta) &= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) R(\tau) \right] \\ &= \mathbb{E} \left[\underbrace{\left(\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \right)}_{f_k} \left(\sum_{k=0}^d \gamma^k r_k \right) \right]\end{aligned}$$

Causality

$$\begin{aligned}\nabla U(\theta) &= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) R(\tau) \right] \\ &= \mathbb{E} \left[\left(\sum_{k=0}^d \underbrace{\nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k)}_{f_k} \right) \left(\sum_{k=0}^d \gamma^k r_k \right) \right] \\ &= \mathbb{E} \left[(f_0 + \dots + f_d) (\gamma^0 r_0 + \dots \gamma^d r_d) \right]\end{aligned}$$

Causality

$$\begin{aligned}
 \nabla U(\theta) &= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) R(\tau) \right] \\
 &= \mathbb{E} \left[\underbrace{\left(\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \right)}_{f_k} \left(\sum_{k=0}^d \gamma^k r_k \right) \right] \\
 &= \mathbb{E} \left[(f_0 + \dots + f_d) (\gamma^0 r_0 + \dots \gamma^d r_d) \right] \\
 &= \mathbb{E} \left[\begin{array}{l} f_0 \gamma^0 r_0 + f_0 \gamma^1 r_1 + f_0 \gamma^2 r_2 + \dots + f_0 \gamma^d r_d \\ + f_1 \gamma^0 r_0 + f_1 \gamma^1 r_1 + f_1 \gamma^2 r_2 + \dots + f_1 \gamma^d r_d \\ \vdots \\ + f_d \gamma^0 r_0 + f_d \gamma^1 r_1 + f_d \gamma^2 r_2 + \dots + f_d \gamma^d r_d \end{array} \right]
 \end{aligned}$$

Causality

$$\begin{aligned}
 \nabla U(\theta) &= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) R(\tau) \right] \\
 &= \mathbb{E} \left[\underbrace{\left(\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) \right)}_{f_k} \left(\sum_{k=0}^d \gamma^k r_k \right) \right] \\
 &= \mathbb{E} \left[(f_0 + \dots + f_d) (\gamma^0 r_0 + \dots \gamma^d r_d) \right] \\
 &= \mathbb{E} \left[\begin{array}{l} f_0 \gamma^0 r_0 + f_0 \gamma^1 r_1 + f_0 \gamma^2 r_2 + \dots + f_0 \gamma^d r_d \\ \cancel{+ f_1 \gamma^0 r_0} + f_1 \gamma^1 r_1 + f_1 \gamma^2 r_2 + \dots + f_1 \gamma^d r_d \\ \vdots \\ \cancel{+ f_d \gamma^0 r_0} + \cancel{f_d \gamma^1 r_1} + \cancel{f_d \gamma^2 r_2} + \dots + f_d \gamma^d r_d \end{array} \right]
 \end{aligned}$$

Causality

$$\begin{aligned}\nabla U(\theta) &= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) R(\tau) \right] \\ &= \mathbb{E} \left[\underbrace{\left(\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \right)}_{f_k} \left(\sum_{k=0}^d \gamma^k r_k \right) \right]\end{aligned}$$

$$= \mathbb{E} \left[(f_0 + \dots + f_d) (\gamma^0 r_0 + \dots + \gamma^d r_d) \right]$$

$$= \mathbb{E} \left[\begin{array}{l} f_0 \gamma^0 r_0 + f_0 \gamma^1 r_1 + f_0 \gamma^2 r_2 + \dots + f_0 \gamma^d r_d \\ + \cancel{f_1 \gamma^0 r_0} + f_1 \gamma^1 r_1 + f_1 \gamma^2 r_2 + \dots + f_1 \gamma^d r_d \\ \vdots \\ + \cancel{f_d \gamma^0 r_0} + \cancel{f_d \gamma^1 r_1} + \cancel{f_d \gamma^2 r_2} + \dots + f_d \gamma^d r_d \end{array} \right]$$

$$= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \left(\sum_{l=k}^d \gamma^l r_l \right) \right]$$

Causality

$$\begin{aligned}\nabla U(\theta) &= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) R(\tau) \right] \\ &= \mathbb{E} \left[\underbrace{\left(\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \right)}_{f_k} \left(\sum_{k=0}^d \gamma^k r_k \right) \right]\end{aligned}$$

$$= \mathbb{E} \left[(f_0 + \dots + f_d) (\gamma^0 r_0 + \dots + \gamma^d r_d) \right]$$

$$= \mathbb{E} \left[\begin{array}{l} f_0 \gamma^0 r_0 + f_0 \gamma^1 r_1 + f_0 \gamma^2 r_2 + \dots + f_0 \gamma^d r_d \\ + \cancel{f_1 \gamma^0 r_0} + f_1 \gamma^1 r_1 + f_1 \gamma^2 r_2 + \dots + f_1 \gamma^d r_d \\ \vdots \\ + \cancel{f_d \gamma^0 r_0} + \cancel{f_d \gamma^1 r_1} + \cancel{f_d \gamma^2 r_2} + \dots + f_d \gamma^d r_d \end{array} \right]$$

$$= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \left(\sum_{l=k}^d \gamma^l r_l \right) \right] = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k r_{k,\text{to-go}} \right]$$

Causality

$$\begin{aligned}\nabla U(\theta) &= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) R(\tau) \right] \\ &= \mathbb{E} \left[\underbrace{\left(\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) \right)}_{f_k} \left(\sum_{k=0}^d \gamma^k r_k \right) \right]\end{aligned}$$

$$= \mathbb{E} \left[(f_0 + \dots + f_d) (\gamma^0 r_0 + \dots + \gamma^d r_d) \right]$$

$$= \mathbb{E} \left[\begin{array}{l} f_0 \gamma^0 r_0 + f_0 \gamma^1 r_1 + f_0 \gamma^2 r_2 + \dots + f_0 \gamma^d r_d \\ \cancel{+ f_1 \gamma^0 r_0} + f_1 \gamma^1 r_1 + f_1 \gamma^2 r_2 + \dots + f_1 \gamma^d r_d \\ \vdots \\ \cancel{+ f_d \gamma^0 r_0} + \cancel{f_d \gamma^1 r_1} + \cancel{f_d \gamma^2 r_2} + \dots + f_d \gamma^d r_d \end{array} \right]$$

$$= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) \left(\sum_{l=k}^d \gamma^l r_l \right) \right] = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) \gamma^k \underline{r_{k,\text{to-go}}} \right]$$

$$Q^{\pi}(s_k, a_k)$$

Baseline Subtraction

Baseline Subtraction

$$\nabla U(\theta) = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k r_{k,\text{to-go}} \right]$$

Baseline Subtraction

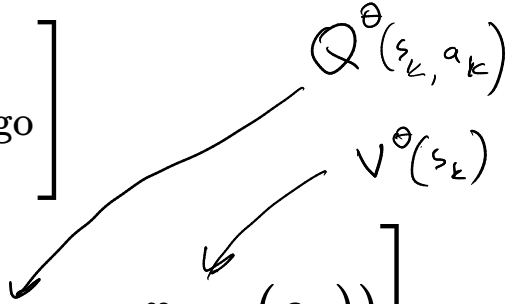
$$\begin{aligned}\nabla U(\theta) &= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k r_{k,\text{to-go}} \right] \\ \nabla U(\theta) &= \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k)) \right]\end{aligned}$$


Diagram illustrating the relationship between the two equations. The first equation shows the gradient of the expected return $\nabla U(\theta)$ as the expectation of the sum of the gradients of the log-probability of the action $\nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k)$ multiplied by the discounted return $\gamma^k r_{k,\text{to-go}}$. The second equation shows the gradient of the expected return $\nabla U(\theta)$ as the expectation of the sum of the gradients of the log-probability of the action $\nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k)$ multiplied by the discounted advantage $\gamma^k (r_{k,\text{to-go}} - r_{\text{base}}(s_k))$. The handwritten terms $Q^{\theta}(s_k, a_k)$ and $V^{\theta}(s_k)$ are shown with arrows pointing to $r_{k,\text{to-go}}$ and $r_{\text{base}}(s_k)$ respectively, indicating that $r_{k,\text{to-go}} = Q^{\theta}(s_k, a_k)$ and $r_{\text{base}}(s_k) = V^{\theta}(s_k)$.

Baseline Subtraction

$$\nabla U(\theta) = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k r_{k,\text{to-go}} \right]$$

$$\nabla U(\theta) = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_{k,\text{to-go}} - \underline{r_{\text{base}}(s_k)}) \right]$$

does not bias
(proof in book)

Baseline Subtraction

$$\nabla U(\theta) = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k r_{k,\text{to-go}} \right]$$

$$\nabla U(\theta) = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k \mid s_k) \gamma^k (r_{k,\text{to-go}} - \underline{r_{\text{base}}(s_k)}) \right]$$

does not bias
(proof in book)

$$r_{\text{base},i} = \frac{\mathbb{E}_{a,s,r_{\text{to-go}},k} [\ell_i(a,s,k)^2 r_{\text{to-go}}]}{\mathbb{E}_{a,s,k} [\ell_i(a,s,k)^2]}$$

Baseline Subtraction

$$\nabla U(\theta) = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) \gamma^k r_{k,\text{to-go}} \right]$$

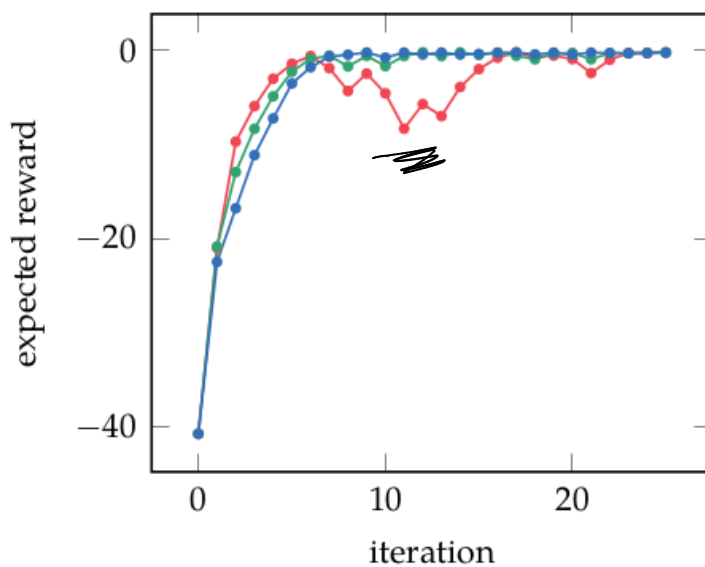
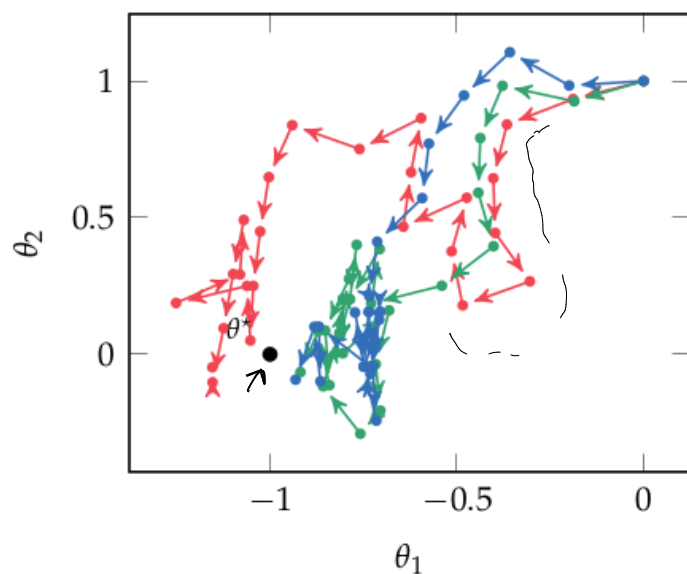
$$\nabla U(\theta) = \mathbb{E} \left[\sum_{k=0}^d \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) \gamma^k (r_{k,\text{to-go}} - \underline{r_{\text{base}}(s_k)}) \right]$$

does not bias
(proof in book)

$$r_{\text{base},i} = \frac{\mathbb{E}_{a,s,r_{\text{to-go}},k} [\ell_i(a,s,k)^2 r_{\text{to-go}}]}{\mathbb{E}_{a,s,k} [\ell_i(a,s,k)^2]}$$

$$\ell_i(a,s,k) = \gamma^{k-1} \frac{\partial}{\partial \theta_i} \log \pi_{\theta}(a | s)$$

In practice $\hat{V}^{\theta}(s_k)$



- likelihood ratio
- reward-to-go
- baseline subtraction

Figure 11.3. Several policy gradient methods used to optimize policies for the simple regulator problem from the same initial parameterization. Each gradient evaluation ran six rollouts to depth 10. The magnitude of the gradient was limited to 1, and step updates were applied with step size 0.2. The optimal policy parameterization is shown in black.

Guiding Questions

Policy Optimization

- What is Policy Gradient?
- What tricks are needed for it to work effectively?