## Last Time

Bandit

- — -ε greedy
- = softmax
- ⌐ Thompson Sampling
- – Interval
- ⟹ - UCB
- – Optimal Dynamic

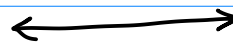Relationship to MCTS

$Q$

---

## This Time

Policy Gradient

| Model Based | ⟷ | Model Free |
|---|---|---|
| estimate $T, R$ | | directly |
| solve with $T, R$ | | optimize $\pi$ or $Q$ |
| | | w/o $T, R$ |

$$\nabla_x f(x) = \left[ \frac{\partial f}{\partial x_i}(x), \ldots \frac{\partial f}{\partial x_n}(x) \right]$$

## Gradient Ascent        optimize $U(\theta)$

loop
$$\theta' \leftarrow \theta + \alpha \nabla U(\theta)$$
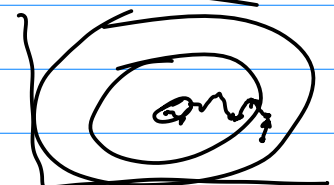$$\theta \leftarrow \theta'$$

step size
decaying
ADAM

Stochastic
Gradient Descent

Probabilistic/Stochastic
Parameterized Policies

$$\pi_\theta(a|s)$$

$\theta = |s| \times |A|$ matrix

$$\pi_\theta(a|s) = \frac{\theta[s,a]}{\sum_a \theta[s,a]}$$

$$(S, A, T, R, \gamma, \underset{\equiv}{\rho_0})$$

initial state . ist

$$\tau = \left( s^{(1)}, a^{(0)}, r^{(0)}, \dots s^{(d)}, a^{(d)}, r^{(d)} \right)$$

step

Episode / Trajectory

Advantage $\qquad A(s,a) = Q_{s,a} - V(s)$

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

$$\nabla U(\theta) =$$

$e^{(2)} = [0, 1, 00]$

$f(x)$

Finite Differencing

$[1, 00, 0]$

$$\nabla U(\theta) \approx \left[ \frac{U(\theta + \delta e^{(1)}) - U(\theta)}{\delta}, \quad \dots \quad \frac{U(\theta + \delta e^{(n)}) - U(\theta)}{\delta} \right]$$

$$U(\theta) \approx \frac{1}{m} \sum_{i=1}^{m} R(\tau_i)$$

Leverage $\qquad \nabla \pi_\theta$

$\frac{\partial}{\partial x} \log x = \frac{1}{x}$

Likelihood Ratio Trick

$$\nabla_\theta \log p_\theta(\tau) = \nabla_\theta p_\theta(\tau) / p_\theta(\tau)$$

$$\therefore \quad \nabla_\theta p_\theta(\tau) = p_\theta(\tau) \nabla_\theta \log p_\theta(\tau)$$

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

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

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

$$= \int p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) R(\tau) d\tau$$

$$= E\left[ \nabla_\theta \log p_\theta(\tau) R(\tau) \right]$$

$$E[f(x)] = \int p(x) f(x) dx$$

$$p_\theta(\tau) = p\left(s^{(1)}\right) \overset{d}{\underset{k=1}{\prod}} T\left(s^{(k+1)} \mid s^{(k)}, a^{(k)}\right) \pi_\theta\left(a^{(k)} \mid s^{(k)}\right) \quad \swarrow^{p_\theta}$$

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

$$\log p_\theta(\tau) = \log p(s^{(1)}) + \overset{d}{\underset{k=1}{\sum}} \log T\left(s^{k+1} \mid s^k, a^k\right) + \overset{d}{\underset{k=1}{\sum}} \log \pi_\theta\left(a^k \mid s^k\right)$$

$$\nabla \log p_\theta(\tau) = \overset{d}{\underset{k=1}{\sum}} \nabla_\theta \log \pi_\theta(a^k \mid s^k)$$

## Policy Gradient

$\theta \leftarrow \text{rand}()$

loop

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

$\quad \theta \leftarrow \theta + \alpha \overset{d}{\underset{k=1}{\sum}} \nabla_\theta \log \pi_\theta\left(a^{(k)} \mid s^{(k)}\right) R(\tau)$

Variance

$$\overline{\text{Break} \cancel{\text{at Room}}}$$

## Causality

$$\nabla U(\theta) = E\left[ \left( \overset{d}{\underset{k=1}{\sum}} \nabla_\theta \log \pi_\theta\left(a^{(k)} \mid s^{(k)}\right) \right) \left( \overset{d}{\underset{k=1}{\sum}} r^{(k)} \gamma^{k-1} \right) \right]$$

$\underbrace{\qquad}_{f^k} \qquad \underbrace{\qquad}_{R(\tau)}$

$$= E\left[ \left( f^1 + f^2 + \dots + f^d \right) \left( r^{(1)} + r^{(2)} \gamma \dots r^d \gamma^{d-1} \right) \right]$$

$$= E\left[ \begin{array}{l} f^1 r^1 + f^1 r^2 \gamma + \dots \quad f^{(1)} r^d \gamma^{d-1} \\ + f^2 r^1 + f^2 r^2 \gamma \quad\quad f^2 r^d \gamma^{d-1} \\ \vdots \\ + f^d r^1 + \quad\quad\quad\quad\quad + f^d r^d \gamma^{d-1} \end{array} \right]$$

$$\nabla U(\theta) = E\left[ \overset{d}{\underset{k=1}{\sum}} \left( \nabla_\theta \log \pi_\theta\left(a^k \mid s^k\right) \left( \overset{d}{\underset{\ell=k}{\sum}} r^d \gamma^{\ell-k} \right) \right) \right]$$

$$= E\left[ \overset{d}{\underset{k=1}{\sum}} \nabla_\theta \log \pi_\theta\left(a^k \mid s^k\right) \gamma^{k-1} r^{(k)}_{\text{to go}} \right]$$

$\hat{Q}(s^k, a^k)$

## Baseline Subtraction

$$\nabla U(\theta) \simeq E\left[ \sum_{k=1}^{d} \nabla_\theta \log \pi_\theta(a^k/s^k) \, \gamma^{k-1} \left( r_{to-go}^k - r_{base}(s^k) \right) \right]$$

Does not bias grad est
(proof in book)

$\longrightarrow$ Good: $r_{base}(s^k) = \hat{V}(s^k) = \frac{1}{m}\sum_{i=1}^{k} r_{togo,i}$

$\hookleftarrow$ previous simulations

Optimal: $r_{base,i} = \dfrac{E\left[ \ell_i(a,s,k)^2 \, r_{togo} \right]}{E\left[ \ell_i(a,s,k)^2 \right]}$

reduces variance

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

$r_{togo}^k - r_{base}(s^k)$

$\longrightarrow$
$$\nabla U(\theta) = E\left[ \sum_{k=1}^{d} \nabla_\theta \log \pi_\theta(a^k/s^k) \, \gamma^{k-1} \, \hat{A}.(s^k, a^k) \right]$$

---

## Recap

Policy Gradient

Running a bunch of simulations
increasing $\pi_\theta(a/s)$ for a that resulted in high
reward

$\nabla U(\theta) \quad \nabla \pi_\theta$

- Likelihood Ratio
- Causality
- Baseline Subtraction

---

$R(s,a)$

$s' \leftarrow$ step!$(env, a)$
$r \leftarrow$ act!$(env, a)$