

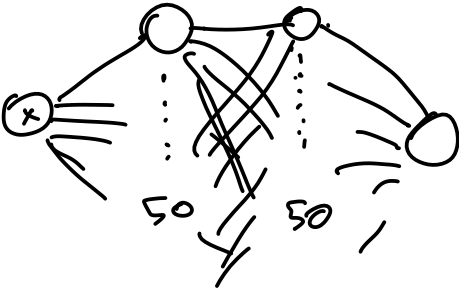
$$\theta \leftarrow \theta + \alpha \frac{\partial Q}{\partial \theta} (Q_{\theta}(s, a) - y)$$

$$\mathcal{L}(s, a, \theta) = (Q_{\theta}(s, a) - y)^2 \quad \uparrow \quad r + \gamma \max_{a'} Q_{\theta'}(s', a')$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = 2 (Q_{\theta}(s, a) - y) \frac{\partial Q}{\partial \theta}$$

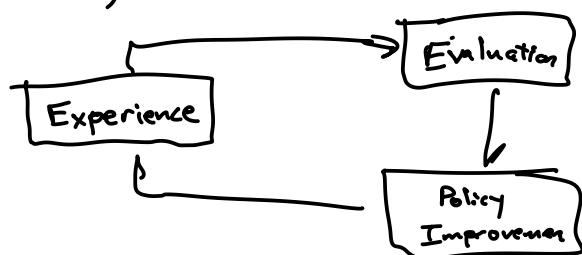
$$y = w_3 \sigma(w_2 \sigma(w_1 x + b_1) + b_2) + b_3$$

$\begin{matrix} 1 \times 50 & 50 \times 50 & 50 \times 50 & 50 & 55 & 1 \end{matrix}$



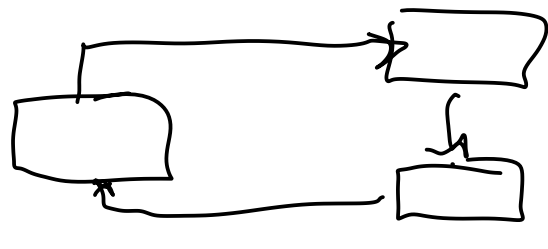
Policy Gradients

Q Learning



$$\theta \leftarrow \theta + \alpha \frac{\partial Q}{\partial \theta} (Q_\theta(s,a) - r - \gamma \max_{a'} Q_\theta(s',a'))$$

Policy Grad



$$\pi(a|s) \begin{cases} 1-\epsilon & \text{if } a = \arg \max Q_\theta(s,a) \\ \frac{\epsilon}{|A|-1} & \text{o.w.} \end{cases}$$

$$\hat{Q} = \sum_t r_t$$

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

$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

$$\tau = (s_1, a_1, s_2, a_2, \dots, s_T, a_T)$$

$$J(\theta) = E_{\tau \sim p(\tau|\pi_\theta)} [r(\tau)]$$

$$J(\theta) = \int p(\tau|\pi_\theta) r(\tau) d\tau$$

$$\nabla_\theta J(\theta) = \int \nabla_\theta p(\tau|\pi_\theta) r(\tau) d\tau$$

$$= \int \frac{\nabla_\theta p(\tau|\pi_\theta)}{p(\tau|\pi_\theta)} p(\tau|\pi_\theta) r(\tau) d\tau$$

$$= \int \nabla_\theta \log p(\tau|\pi_\theta) p(\tau|\pi_\theta) r(\tau) d\tau$$

$$= E_{\tau \sim p(\tau|\pi_\theta)} [\nabla_\theta \log p(\tau|\pi_\theta) r(\tau)]$$

$$\frac{d}{dx} \log(f(x)) = \frac{df(x)}{dx} \frac{1}{f(x)}$$

$$\pi: S \rightarrow A$$

$$\pi: S \times A \rightarrow \mathbb{R}$$

prob
a in s

$$p(\tau|\pi_\theta) = p(s_1) \prod_t \pi_\theta(a_t|s_t) T(s_{t+1}|s_t, a_t)$$

$$\log p = \log p(s_1) + \sum_t \log \pi(a_t|s_t) + \log T(s_{t+1}|s_t, a_t) \quad \left\{ \log(xy) = \log x + \log y \right.$$

$$\nabla_\theta \log p = \sum_t \nabla_\theta \log \pi_\theta(a_t|s_t)$$

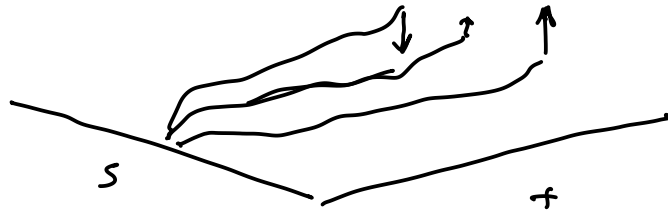
$$\nabla_\theta J(\theta) = E_{\tau \sim p(\tau|\pi_\theta)} \left[\left(\sum_t \nabla_\theta \log \pi_\theta \right) \left(\sum_t R(s_t, a_t) \right) \right]$$

MC simulations

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_t \nabla_\theta \log \pi_\theta(a_{t,i}|s_{t,i}) \sum_t r_{t,i}$$

simulation

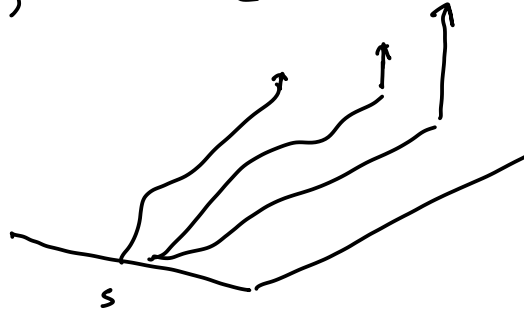
\int
 $p(s)$ $p(s')$



REINFORCETM

1. sample $\{\tau_i\}$ (run π_θ)
 2. $\nabla_\theta J(\theta) = \frac{1}{N} \sum_i \sum_\tau \nabla_\theta \log \pi_\theta r(\tau)$
 3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$
- } Evaluation
 } improvement

Won't work: High Variance



1. Causality

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_i \sum_{\tau} \left(\nabla_{\theta} \log \pi_{\theta} \sum_{t'=\tau}^T r_{i,t'} \right)$$

Always do this

$\uparrow Q(s_{i,\tau}, a_{i,\tau})$

2. Baselines

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_i \sum_{\tau} \left(\nabla_{\theta} \log \pi_{\theta} \left[r(\tau_i) - \underline{b} \right] \right)$$

$$b = \frac{1}{N} \sum r(\tau)$$

$$\begin{aligned} \mathbb{E}_{\tau} [\nabla_{\theta} \log \pi_{\theta} b] &= \int \pi_{\theta} \nabla_{\theta} \log \pi_{\theta} b d\tau \\ &= \int \nabla_{\theta} p(\tau | \pi_{\theta}) b d\tau \\ &= b \nabla_{\theta} \left(\int p(\tau | \pi_{\theta}) d\tau \right) \\ &= 0 \end{aligned}$$

Adding a baseline does not bias gradient

$$\frac{dVar}{db} = 0$$

$$b^* = \frac{\mathbb{E}[g(\tau)^2 r(\tau)]}{\mathbb{E}[g(\tau)^2]}$$

In practice use $b = \frac{1}{N} \sum r(\tau)$

REINFORCE +

it works :)

REINFORCE is on policy :)

$\pi \neq \pi_{\theta}$

Importance Sampling

$$\begin{aligned} \mathbb{E}_{x \sim p(x)} [f(x)] &= \int p(x) f(x) dx \\ &= \int \frac{q(x)}{p(x)} p(x) f(x) dx \\ &= \int q(x) \frac{p(x)}{q(x)} f(x) dx \\ &= \mathbb{E}_{x \sim q(x)} \left[\frac{p(x)}{q(x)} f(x) \right] \end{aligned}$$

$$J(\theta) = E_{\tau \sim p(\tau|\bar{\pi})} \left[\frac{p(\tau|\pi_\theta)}{p(\tau|\bar{\pi})} r(\tau) \right]$$

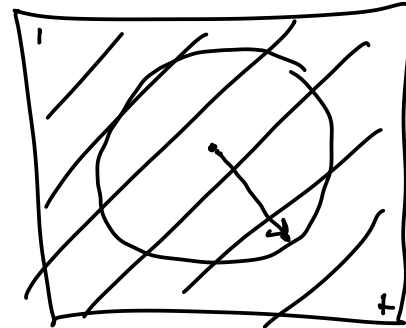
$$\nabla_{\theta'} J(\theta') = E \left[\underbrace{\frac{1}{T} \sum_{t=1}^T \frac{\pi_{\theta'}(a_t | s_t)}{\bar{\pi}(a_t | s_t)}}_{\uparrow} \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left(r(\tau) \right) \right]$$

gets really small

~~Gradient Ascent~~ Policy Iteration

$$\theta' = \underset{\theta'}{\operatorname{argmax}} \nabla_{\theta} J(\theta)^T (\theta' - \theta)$$

$$\text{s.t. } \|\theta - \theta'\|^2 \leq \epsilon$$



Actually want something like

$$\theta' = \underset{\theta'}{\operatorname{argmax}} \nabla_{\theta} J(\theta)^T (\theta' - \theta)$$

$$\text{s.t. } D_{KL}(\pi_{\theta'}(a|s) || \pi_{\theta}(a|s)) \leq \epsilon$$

$$D_{KL}(p || q) =$$

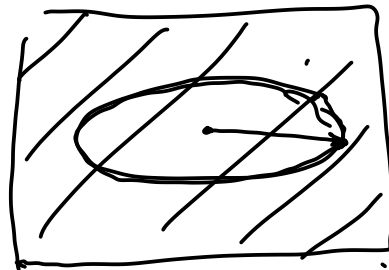
$$E \left[\log \frac{p(x)}{q(x)} \right]$$

Taylor Expansion

$$D_{KL}(\) \approx \frac{1}{2} (\theta' - \theta)^T F (\theta' - \theta)$$

$$\theta' = \theta + \alpha F^{-1} \nabla_{\theta} J(\theta)$$

Fisher info matrix
= "Natural Gradient"



Trust Region Policy Opt TRPO

$$\alpha = \sqrt{\frac{2\epsilon}{\nabla_{\theta} J(\theta)^T F \nabla_{\theta} J(\theta)}}$$

Proximal Policy Optimization PPO

Use regularization to stay close to old policy
so that we can use importance sampling