

# 1) PG with value functions

claim:  $s, a \sim d_{\mu_0}^{\pi_\theta}$

$$g = \frac{1}{1-\gamma} \nabla_{\theta} \log(\pi_{\theta}(a|s)) \left[ \overbrace{Q^{\pi_{\theta}}(s,a) - b(s)}^{\substack{\text{action} \\ \text{independent} \\ \text{baseline}}} \right]$$

is an unbiased estimate of  $\nabla J(\theta)$

Proof:  $\nabla J(\theta) = \nabla_{\theta} \mathbb{E}_{s_0 \sim \mu_0} [V^{\pi_{\theta}}(s_0)]$

skipping steps (lecture notes)

$$= \mathbb{E}_{\substack{s_0 \sim \mu_0 \\ a_0 \sim \pi_{\theta}(s_0)}} \left[ \nabla_{\theta} \log \pi_{\theta}(a_0|s_0) Q^{\pi_{\theta}}(s_0, a_0) \right] + \gamma \mathbb{E}_{s_1 \sim P_1^{\pi_{\theta}}} \left[ \nabla_{\theta} V^{\pi_{\theta}}(s_1) \right]$$

iteration

marginalized over states

$$\nabla J(\theta) = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s_t, a_t \sim P_t^{\pi_{\theta}}} \left[ \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) Q^{\pi_{\theta}}(s_t, a_t) \right]$$

$$= \frac{1}{1-\gamma} \mathbb{E}_{s, a \sim d_{\mu_0}^{\pi_{\theta}}} \left[ \nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s, a) \right]$$

## "Actor Critic" Algorithm

1)  $s, a \sim d_{\mu_0}^{\pi_{\theta_t}}$

2)  $g = \frac{1}{1-\gamma} \nabla_{\theta} \log(\pi_{\theta}(a|s)) (Q^{\pi_{\theta}}(s,a) - \underbrace{V^{\pi_{\theta}}(s)}_{A^{\pi_{\theta}}(s,a)})$

## 2) Trust Regions & KL-Divergence

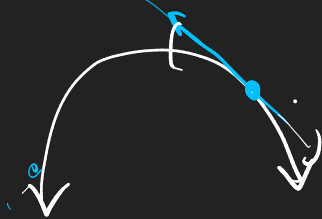
$$\max_{\theta} J(\theta) \approx \left[ \max_{\theta} J(\theta_0) + \underbrace{(\nabla J(\theta_0))^T (\theta - \theta_0)}_{\text{first order approx}} \right]$$

The maximum occurs when  $\theta - \theta_0$  is parallel to  $\nabla J(\theta_0)$

$$\theta - \theta_0 = \alpha \nabla J(\theta_0) \quad \alpha > 0$$

$$\theta = \theta_0 + \alpha \nabla J(\theta_0)$$

Choose small enough  $\alpha$  so that the linearization is good enough



A trust region approach makes this intuition formal:

$$\left[ \begin{array}{l} \max_{\theta} J(\theta) \\ \text{s.t. } d(\theta, \theta_0) < \delta \end{array} \right]$$

← bounded "distance" from  $\theta_0$

Another motivation: if  $J(\theta)$  is estimated it might only be a good approx. when  $\theta$  is near  $\theta_0$  e.g. conservative policy iteration

$$\pi^{t+1}(\cdot|s) = (1-\alpha)\pi^t(\cdot|s) + \alpha \pi^t(s)$$

## KL Divergence

way to measure "distance" between  $\theta_0$  and  $\theta$

KL Divergence: for  $P \in \Delta(\mathcal{X})$  and  $Q \in \Delta(\mathcal{X})$

$$KL(P|Q) = \mathbb{E}_{x \sim P} \left( \log \left( \frac{P(x)}{Q(x)} \right) \right) = \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{P(x)}{Q(x)} \right)$$

example:  $P = \mathcal{N}(\mu_1, \sigma^2 I)$      $Q = \mathcal{N}(\mu_2, \sigma^2 I)$

$$KL(P|Q) = \frac{\|\mu_1 - \mu_2\|_2^2}{\sigma}$$

Define our "distance" in terms of the KL divergence between  $\pi_{\theta_0}(\cdot|s) \in \Delta(\mathcal{S})$  and  $\pi_{\theta}(\cdot|s) \in \Delta(\mathcal{S})$  on average over  $d_{\mu_0}^{\pi_{\theta_0}}$

$$d_{KL}(\theta_0, \theta) = \mathbb{E}_{s \sim d_{\mu_0}^{\pi_{\theta_0}}(s)} [KL(\pi_{\theta_0}(\cdot|s) | \pi_{\theta}(\cdot|s))]$$

marginalized over actions

$$= \mathbb{E}_{s, a \sim d_{\mu_0}^{\pi_{\theta_0}}} \left[ \log \left( \frac{\pi_{\theta_0}(a|s)}{\pi_{\theta}(a|s)} \right) \right]$$

### 3) Natural Policy Gradient

Alg: Natural PG

Initiate  $\theta_0$

for  $t=0, 1, \dots$

Estimating  $\nabla J(\theta_t)$  with  $g_t$

Estimate Fisher Information matrix by

$$F_t = \nabla_{\theta} \log(\pi_{\theta_t}(a|s)) \nabla_{\theta} \log(\pi_{\theta_t}(a|s))^T$$

$s, a \sim d_{\mu}^{\pi_{\theta_t}}$

Natural gradient step  $\leftarrow$  pseudoinverse

$$\theta_{t+1} = \theta_t + \alpha F_t^+ g_t$$

trajectory,  $q$ , a function based

The gradient is preconditioned by the Fisher Info. matrix

Derive as an approximation to constrained optimization:

$$\begin{cases} \max_{\theta} J(\theta) \longrightarrow \text{first order approx} \\ \text{s.t. } d_{KL}(\theta_0, \theta) \leq \delta \longrightarrow \text{second order approx} \end{cases}$$

Second order approx

$$l(\theta) = d_{KL}(\theta_0, \theta) = \mathbb{E}_{S, a \sim d_{\theta_0}} \left( \log \left( \frac{\pi_{\theta_0}(a|s)}{\pi_{\theta}(a|s)} \right) \right)$$

$$l(\theta) \approx l(\theta_0) + \nabla l(\theta_0)^T (\theta - \theta_0) + (\theta - \theta_0)^T \nabla^2 l(\theta_0) (\theta - \theta_0)$$

Claim:  $l(\theta_0) = 0$ ,  $\nabla l(\theta_0) = 0$

$$\nabla^2 l(\theta_0) = \mathbb{E}_{S, a \sim d_{\theta_0}} \left[ \nabla_{\theta} \log(\pi_{\theta}(a|s)) \nabla_{\theta} \log(\pi_{\theta}(a|s))^T \Big|_{\theta=\theta_0} \right]$$

"Fisher Information matrix"  $\uparrow$   
 $F_{\theta_0}$

Proof:

$$\nabla_{\theta} l(\theta) = \mathbb{E}_{S, a \sim d_{\theta_0}} \left[ \nabla_{\theta} \left( \log \pi_{\theta_0}(a|s) - \log(\pi_{\theta}(a|s)) \right) \right]$$

$$\nabla_{\theta} l(\theta) \Big|_{\theta=\theta_0} = - \mathbb{E}_{S \sim d_{\theta_0}} \left[ \sum_{a \in \mathcal{A}} \frac{\pi_{\theta_0}(a|s) \nabla_{\theta} \pi_{\theta_0}(a|s)}{\pi_{\theta_0}(a|s)} \right] = 0 \quad \checkmark$$

Trust Region approximate maximization

$$\theta^* = \arg \max_{\theta} \nabla l(\theta_0)^T (\theta - \theta_0) + J(\theta_0)$$

$$\text{s.t. } (\theta - \theta_0)^T F_{\theta_0} (\theta - \theta_0) \leq \delta$$

Claim: Solve in closed form:

$$\Theta^* = \Theta_0 + \alpha F_{\Theta_0}^{-1} \nabla J(\Theta_0)$$

$$\alpha = \left( \frac{\delta}{\nabla J(\Theta_0)^T F_{\Theta_0}^{-1} \nabla J(\Theta_0)} \right)^{1/2}$$

Hint:  $V = F_{\Theta_0}^{1/2} (\Theta - \Theta_0)$  and  $C = F_{\Theta_0}^{1/2} \nabla J(\Theta_0)$

$$\begin{cases} \max & C^T V \\ \text{s.t.} & \|V\|_2 \leq \delta \\ & V = C \cdot \frac{\delta}{\|C\|_2} \end{cases}$$

