

# Lecture 15: Policy Optimization with Trust Regions

## 1) Policy Gradient with value functions

PG w/ trajectories often has high variance. An alternative commonly used in practice uses an alternative estimate using Q functions.

Claim: for  $s, a \sim d_{\gamma_0}^{\pi_\theta}$ ,

the gradient of the log-likelihood is called the score

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

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

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

value fn. def.

$$= \mathbb{E}_{s_0 \sim \mathcal{M}} [\nabla_{\theta} \mathbb{E}_{a_0 \sim \pi_{\theta}(s_0)} [Q^{\pi_{\theta}}(s_0, a_0)]]$$

$s_0$  doesn't depend on  $\theta$  & Q fn. def.

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{a_0 \sim \pi_{\theta}(s_0)} [Q^{\pi_{\theta}}(s_0, a_0)] &= \sum_{a_0 \in \mathcal{A}} \nabla_{\theta} [\pi(a_0|s_0) Q^{\pi_{\theta}}(s_0, a_0)] \quad \text{defn. of expectation} \\ &\stackrel{\text{product rule}}{=} \sum_{a_0 \in \mathcal{A}} (\nabla_{\theta} \pi(a_0|s_0)) Q^{\pi_{\theta}}(s_0, a_0) + \pi(a_0|s_0) \nabla_{\theta} Q^{\pi_{\theta}}(s_0, a_0) \\ &\quad \left\{ \begin{array}{l} \text{importance weighting} \\ \text{trick} \end{array} \right. \quad \left\{ \begin{array}{l} r(s_0, a_0) \text{ doesn't} \\ \text{depend on } \theta \end{array} \right. \\ &= \mathbb{E}_{a_0 \sim \pi_{\theta}(s_0)} [\nabla_{\theta} \log \pi_{\theta}(a_0|s_0) Q^{\pi_{\theta}}(s_0, a_0)] + \gamma \mathbb{E}_{\substack{a_0 \sim \pi_{\theta}(s_0) \\ s_1 \sim P(s_0, a_0)}} [\nabla_{\theta} V^{\pi_{\theta}}(s_1)] \end{aligned}$$

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

$$\mathbb{E}_{s_0 \sim \mu_0} \left[ \nabla_\theta V^{\pi_\theta}(s_0) \right]$$

we can iterate!

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

expanding  
expectation

$$= \sum_{t=0}^{\infty} \sum_{\substack{s \in \mathcal{S} \\ a \in \mathcal{A}}} P_t^{\pi_\theta}(s, a; \mu_0) \gamma^t \cdot \nabla_\theta \log \pi_\theta(a | s) \cdot Q^{\pi_\theta}(s, a)$$

definition  
of  $d_{\mu_0}^{\pi_\theta}$

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

One final gradient estimate:  $s, a \sim d_{\mu_0}^{\pi_\theta}$

$$g = \frac{1}{1-\gamma} \underbrace{\nabla_\theta \log \pi_\theta(a | s)}_{\text{score}} \cdot (Q^{\pi_\theta}(s, a) - b(s))$$

Baseline function  $b(s)$  further helps in variance reduction. Most common  $b(s) = V^{\pi_\theta}(s)$  results in advantage function-based PG

$$A^{\pi_\theta}(s, a) = Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s).$$

Policy gradients that use estimate value (Q or A) functions are called

"Actor critic"  
 policy  $\nearrow$   $\nwarrow$  value fn.

To show that  $g$  with a baseline is unbiased, we show that

$$\mathbb{E}[\nabla_{\theta} \log \pi_{\theta}(a|s) \cdot b(s)] = 0$$

for any  $a \sim \pi_{\theta}(s)$  action-independent baseline.

$$\sum_a \pi_{\theta}(a|s) \cdot \frac{\nabla_{\theta} \pi_{\theta}(a|s)}{\pi_{\theta}(a|s)} \cdot b(s) \quad (\text{expanding exp \& grad})$$

$$= \nabla_{\theta} \sum_a \pi_{\theta}(a|s) \cdot b(s) \quad (\text{linearity of grad})$$

$$= \nabla_{\theta} [1 \cdot b(s)] = 0$$

doesn't depend on  $\theta$

( $\pi_{\theta}(\cdot|s)$  is probability distribution)

## 2) Trust Regions & KL-Divergence

Recall: motivation of EA by first order approximate maximization

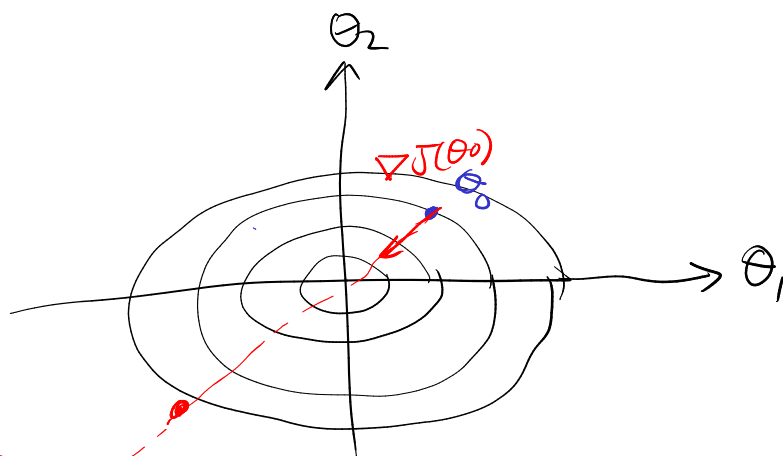
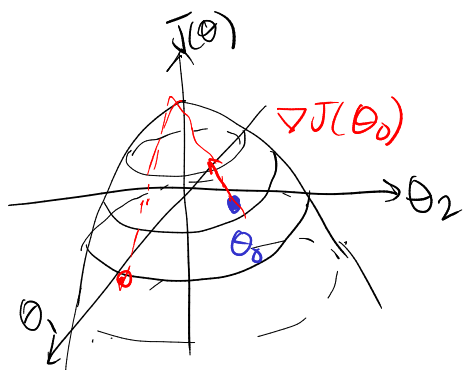
$$\max_{\theta} J(\theta) \approx \max_{\theta} J(\theta_0) + \nabla J(\theta_0)^T (\theta - \theta_0)$$

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

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

Question: why do we normally use a small step size  $\alpha$ ?  
 wouldn't as big  $\alpha$  as possible achieve a higher maximum value?

Answer: The linear approximation is only locally valid, so by choosing small step size  $\alpha$ , we ensure that  $\theta$  is close to  $\theta_0$ .



Following the gradient too far might even lead to decreasing  $J(\theta)$

A trust region approach makes the intuition about the step size more precise:

$$\begin{aligned} \max_{\theta} J(\theta) \\ \text{s.t. } d(\theta, \theta_0) < \delta \end{aligned}$$

trust region is described by bounded "distance" from  $\theta_0$

Another motivation for trust regions when it comes to RL: we might estimate  $J(\theta)$  using data collected with  $\theta_0$  (i.e. a policy  $\pi_{\theta_0}$ ). So our estimate might only be good close to  $\theta_0$ .

E.g. in conservative policy iteration, incremental update:

$$\begin{aligned}\pi'(s) &= \underset{a}{\operatorname{argmax}} \hat{Q}(s, a) \\ \pi^{t+1}(\cdot|s) &= (1-\alpha)\pi^t(\cdot|s) + \alpha\pi'(\cdot|s)\end{aligned}$$

## K-L Divergence:

In order to formulate a trust region problem for policy optimization, we need to decide how to measure the "distance" between  $\theta_t$  and  $\theta_{t+1}$ .

The K-L Divergence measures the "distance" between two distributions. Given  $P \in \Delta(\mathcal{X})$  and  $Q \in \Delta(\mathcal{X})$

Define (K-L Divergence)

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

Ex: if  $P = \mathcal{N}(\mu_1, \sigma^2 I)$  and  $Q = \mathcal{N}(\mu_2, \sigma^2 I)$  then

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

Fact:  $KL(P|Q) \geq 0$  and  $KL(P|Q) = 0 \Leftrightarrow P = Q$ .

KL divergence is a natural way to constrain policy updates because it directly considers the difference in the distributions.

We define a measure of "distance" between  $\pi_{\theta_0}(\cdot|s)$  and  $\pi_{\theta}(\cdot|s)$  averaged over states  $s$  from the discounted-steady-state distribution of  $\pi_{\theta_0}$

$$d_{KL}(\theta_0, \theta) = \mathbb{E}_{s \sim d_{\gamma_0}^{\pi_{\theta_0}}(s)} [KL(\pi_{\theta_0}(\cdot|s) | \pi_{\theta}(\cdot|s))]$$

↑  
marginalized over  $a$

$$= \mathbb{E}_{s \sim d_{\gamma_0}^{\pi_{\theta_0}}(s)} \left[ \mathbb{E}_{a \sim \pi_{\theta_0}(a|s)} \left[ \log \left( \frac{\pi_{\theta_0}(a|s)}{\pi_{\theta}(a|s)} \right) \right] \right]$$

$$= \mathbb{E}_{s, a \sim d_{\gamma_0}^{\pi_{\theta_0}}} \log \left( \frac{\pi_{\theta_0}(a|s)}{\pi_{\theta}(a|s)} \right)$$

# 3) Natural Policy Gradient

Alg: Natural PG

initialize  $\theta_0$

for  $t=0, 1, \dots$

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

Estimate Fisher information matrix by

$$F_t = \underbrace{\nabla \log(\pi_{\theta_t}(a|s))}_{\text{score}} \underbrace{\nabla \log(\pi_{\theta_t}(a|s))}_{\text{score}}^T \text{ for } s, a \sim d_{\pi_{\theta_t}}$$

Natural Gradient step:

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

The gradient is preconditioned by the Fischer information matrix.

Derive as approximating constrained optimization

$$\max_{\theta} J(\theta) \longrightarrow \text{Gradient Ascent: first order approx}$$

$$\text{s.t. } d_{KL}(\theta_0, \theta) \leq \delta$$

idea: second order approx!

A second order approximation to the divergence

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

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

Claim:  $\ell(\theta_0) = 0$ ,  $\nabla \ell(\theta_0) = 0$ , and

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

$\nwarrow$  fisher information matrix  $F_{\theta_0}$

Proof:  $\ell(\theta_0) = KL(p_{\theta_0} | p_{\theta_0}) = 0 \checkmark$

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

$$= \mathbb{E}_{s, a \sim \pi_{\theta_0}} \left( - \frac{\nabla_{\theta} \pi_{\theta}(a|s)}{\pi_{\theta}(a|s)} \right)$$

$$\nabla \ell(\theta_0) = \mathbb{E}_{s \sim d_{\pi_0}} \left[ \sum_a \cancel{\pi_{\theta_0}(a|s)} \cdot \frac{-\nabla \pi_{\theta_0}(a|s)}{\cancel{\pi_{\theta_0}(a|s)}} \right]$$

$$= - \mathbb{E}_s \left[ \nabla_{\theta} \underbrace{\sum_a \pi_{\theta}(a|s)}_{=1} \Big|_{\theta=\theta_0} \right]$$

$$= - \mathbb{E}_s [\nabla_{\theta} (1)] = 0 \quad \checkmark$$



$$\nabla^2 \ell(\theta) = \mathbb{E}_{S, \text{and } \theta_0} \left[ \frac{-\nabla_{\theta_0}^2 \pi_{\theta_0}(a|s)}{\pi_{\theta_0}(a|s)} + \frac{\nabla_{\theta} \pi_{\theta_0}(a|s) \nabla_{\theta} \pi_{\theta_0}(a|s)^T}{\pi_{\theta_0}(a|s)^2} \right]$$

$$\nabla^2 \ell(\theta_0) = \mathbb{E}_S \sum_a \pi_{\theta_0}(a|s) \frac{-\nabla_{\theta_0}^2 \pi_{\theta_0}(a|s)}{\pi_{\theta_0}(a|s)} + \mathbb{E}_{S, \text{and } \theta_0} \left[ \nabla \log(\pi_{\theta_0}(a|s)) \nabla \log(\pi_{\theta_0}(a|s))^T \right]$$

0 by same logic as above

Therefore, the Trust Region constrained approximate maximization:

$$\begin{aligned} \max_{\theta} \quad & \nabla J(\theta_0)^T (\theta - \theta_0) \\ \text{s.t.} \quad & (\theta - \theta_0)^T F_{\theta_0} (\theta - \theta_0) \leq \delta \end{aligned}$$

Claim: This maximization can be solved in closed form:

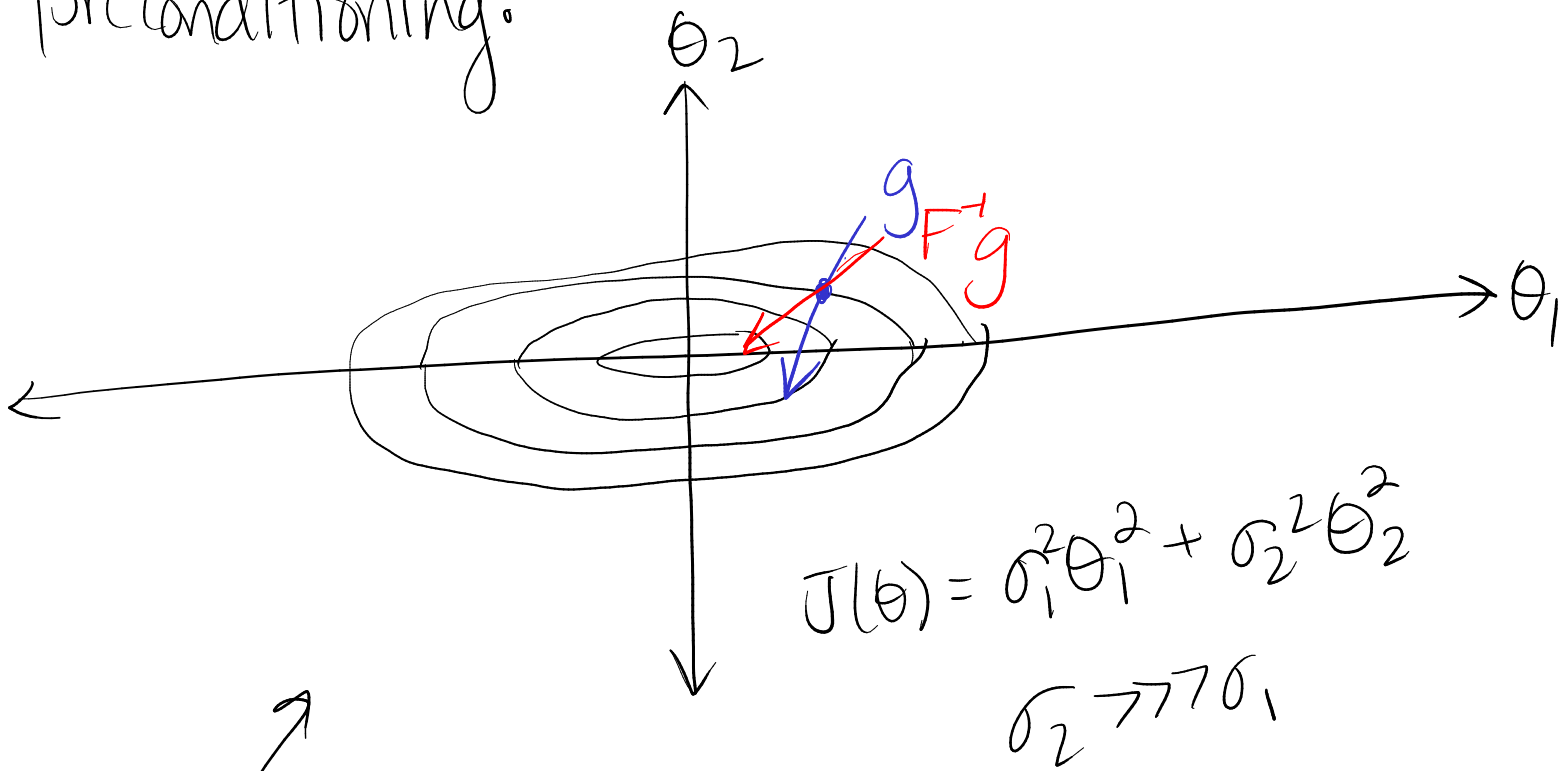
$$\theta = \theta_0 + \alpha F_{\theta_0}^{-1} \nabla J(\theta_0)$$

$$\text{where } \alpha = \left( \frac{\delta}{\nabla J(\theta_0)^T F_{\theta_0}^{-1} \nabla J(\theta_0)} \right)^{1/2}$$

Exercise: show that this is true.

Hint: let  $v = F_{\theta_0}^{1/2} (\theta - \theta_0)$  and  $c = F_{\theta_0}^{-1/2} \nabla J(\theta_0)$  and consider  $\max c^T v$  s.t.  $\|v\|_2^2 \leq \delta$

Intuitive explanation of the benefit  
preconditioning:



Step along vertical

axis. preconditioning by

$$F = \begin{bmatrix} \sigma_1 & \\ & \sigma_2 \end{bmatrix}$$

accounts for this  
and adjusts the stepsize on  $\theta_1$  vs.  $\theta_2$ .