

Policy Gradient

Lecturer: Pieter Abbeel

Scribe: Jan Biermeier

1 Recap

Recall:

$$U(\theta) = \mathbb{E}\left[\sum_{t=0}^H R(s_t, a_t); \pi_\theta\right] = \mathbb{E}[R(\tau); \pi_\theta] \quad (1)$$

Here τ is a sample path of states and actions, $s_0, a_0, \dots, s_H, a_H$.Example policy π_θ :

$$a \in \{0, 1\}, \pi_\theta(a|s_t) = \frac{e^{-\theta^\top \phi(s_t)}}{1 + e^{-\theta^\top \phi(s_t)}} \quad (2)$$

2 Policy Search

$$\max_{\theta} U(\theta) = \max_{\theta} \mathbb{E}[R(\tau); \pi_\theta] \quad (3)$$

$$= \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \quad (4)$$

taking the gradient w.r.t. θ gives

$$\nabla_{\theta} U(\theta) = \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \quad (5)$$

$$= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \quad (6)$$

$$= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau) \quad (7)$$

$$= \sum_{\tau} P(\tau; \theta) \frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)} R(\tau) \quad (8)$$

$$= \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) R(\tau) \quad (9)$$

Approximate with the empirical estimate for m sample paths under policy π_θ :

$$\hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) R(\tau^{(i)}) \quad (10)$$

$$\nabla_{\theta} \log P(\tau^{(i)}; \theta) = \nabla_{\theta} \log \left[\underbrace{\prod_{t=0}^H P(s_{t+1}^{(i)} | s_t^{(i)}, a_t^{(i)})}_{\text{dynamics model}} \cdot \underbrace{\pi_{\theta}(a_t^{(i)} | s_t^{(i)})}_{\text{policy}} \right] \quad (11)$$

$$= \nabla_{\theta} \left[\sum_{t=0}^H \log P(s_{t+1}^{(i)} | s_t^{(i)}, a_t^{(i)}) + \sum_{t=0}^H \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right] \quad (12)$$

$$= \nabla_{\theta} \sum_{t=0}^H \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \quad (13)$$

$$= \sum_{t=0}^H \underbrace{\nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})}_{\text{no dynamics model required!!}} \quad (14)$$

Note that:

$$\sum_{\tau} P(\tau; \theta) = 1 \quad (15)$$

$$\Rightarrow \nabla_{\theta} \sum_{\tau} P(\tau; \theta) = 0 \quad (16)$$

$$\Leftrightarrow \sum_{\tau} \nabla_{\theta} P(\tau; \theta) = 0 \quad (17)$$

$$\Leftrightarrow \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) = 0 \quad (18)$$

$$\Rightarrow \mathbb{E} \left[\sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) \right] = 0 \quad (19)$$

$$\Rightarrow \mathbb{E}_{\tau} [\log P(\tau; \theta)] = 0 \quad (20)$$

Unbiased gradient estimate:

$$\hat{g}_j = \frac{\partial U(\theta)}{\partial \theta_j} \quad (21)$$

$$= \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \theta_j} \log P(\tau^{(i)}; \theta) \cdot R(\tau^{(i)}) \quad (22)$$

$$= \frac{1}{m} \sum_{i=1}^m \left[\frac{\partial}{\partial \theta_j} \log P(\tau^{(i)}; \theta) \cdot R(\tau^{(i)}) - \frac{\partial}{\partial \pi_j} \log P(\tau^{(i)}; \theta) \cdot b_j^{(i)} \right] \quad (23)$$

$$(24)$$

$$\Leftrightarrow \hat{g}_j = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \theta_j} \log P(\tau^{(i)}; \theta) \cdot \left(R(\tau^{(i)}) - b_j^{(i)} \right) \quad (25)$$

\hat{g}_j is an unbiased estimate, with free parameters $b_j^{(i)}$. As $b_j^{(i)}$ has to be a constant, in practice, there is typically no reason to let it depend on i , as all trajectories $\tau^{(i)}$ are (presumably) sampled i.i.d. So in practice, we typically use b_j .

$$\hat{g}_j = \frac{\partial}{\partial \theta_j} \log P(\tau^{(i)}; \theta) \cdot (R(\tau) - b_j), \quad (26)$$

While our gradient estimates are unbiased, in that:

$$\mathbb{E} \hat{g}_j = \frac{\partial U(\theta)}{\partial \theta_j} \quad (27)$$

they are stochastic estimates and have a variance given by:

$$\mathbb{E} \left[(\hat{g}_j - \mathbb{E} [\hat{g}_j])^2 \right] \quad (28)$$

We will now describe how to choose b_j such that we are minimizing the variance of our gradient estimates:

$$\min_{b_j} \mathbb{E} \left[(\hat{g}_j - \mathbb{E} [\hat{g}_j])^2 \right] = \mathbb{E} \hat{g}_j^2 + \mathbb{E} \left[(\mathbb{E} \hat{g}_j)^2 \right] - 2\mathbb{E} [\hat{g}_j - \mathbb{E} [\hat{g}_j]] \quad (29)$$

$$= \mathbb{E} \hat{g}_j^2 + (\mathbb{E} \hat{g}_j)^2 - 2\mathbb{E} [\hat{g}_j] \mathbb{E} [\hat{g}_j] \quad (30)$$

$$= \mathbb{E} \hat{g}_j^2 - \underbrace{(\mathbb{E} \hat{g}_j)^2}_{= \frac{\partial U(\theta)}{\partial \theta_j} - \text{independent of } b_j} \quad (31)$$

$$\min_{b_j} \mathbb{E} \hat{g}_j^2 = \min_{b_j} \mathbb{E} \left[\left(\frac{\partial}{\partial \theta_i} \log P(\tau; \theta) \cdot (R(\tau) - b_j) \right)^2 \right] \quad (32)$$

$$= \min_{b_j} \mathbb{E} \left[\left(\frac{\partial}{\partial \theta_i} \log(\tau; \theta) \right)^2 \cdot (R(\tau)^2 + b_j^2 - 2b_j R(\tau)) \right] \quad (33)$$

$$= \min_{b_j} \mathbb{E} \left[\underbrace{\left(\frac{\partial}{\partial \theta_i} \log P(\tau; \theta) \right)^2 \cdot R(\tau)^2}_{\text{independent of } b_j} + \mathbb{E} \left[\left(\frac{\partial}{\partial \theta_i} \log P(\tau; \theta) \right)^2 \cdot b_j^2 \right] \right] \quad (34)$$

$$- 2\mathbb{E} \left[\left(\frac{\partial}{\partial \theta_i} \log P(\tau; \theta) \right)^2 \cdot b_j R(\tau) \right] \quad (35)$$

$$= \min_{b_j} b_j^2 \cdot \mathbb{E}_\tau \left[\left(\frac{\partial}{\partial \theta_i} \log P(\tau; \theta) \right)^2 \right] - 2b_j \mathbb{E}_\tau \left[\left(\frac{\partial}{\partial \theta_i} \log P(\tau; \theta) \right)^2 R(\tau) \right] \quad (36)$$

$$\frac{\partial}{\partial b_j} = 0 \Rightarrow 2b_j \mathbb{E}_\tau [\dots] - 2\mathbb{E}_\tau [\dots] = 0 \quad (37)$$

$$\Rightarrow b_j = \frac{\mathbb{E}_\tau \left[\left(\frac{\partial}{\partial \theta_i} \log P(\tau; \theta) \right)^2 \cdot R(\tau) \right]}{\mathbb{E}_\tau \left[\left(\frac{\partial}{\partial \theta_i} \log P(\tau; \theta) \right)^2 \right]} \quad (38)$$

If we need to minimize the variance, we can compute b_j as above from samples for a good estimate.

3 Gradient descent

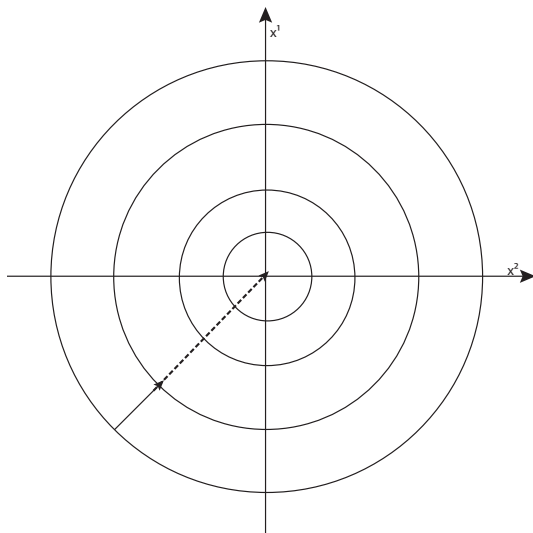


Figure 1: $f(x) = -x_1^2 - x_2^2$

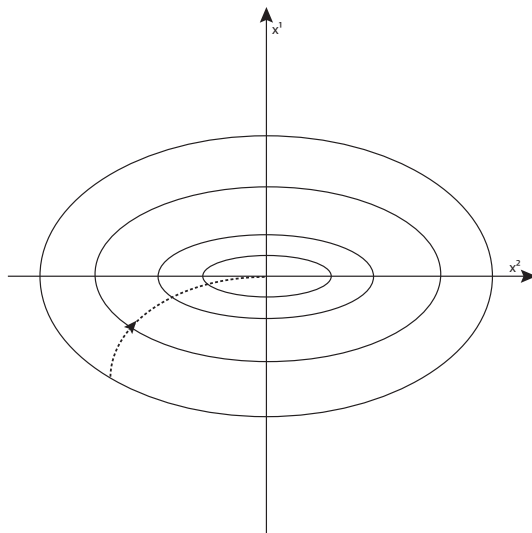


Figure 2: $\bar{f}(x) = -x_1^2 - 100x_2^2$

For function $f(x) = -x_1^2 - x_2^2$ in Figure 1, $\nabla f = -x_1 - x_2$ leads directly to the maximum. Note that $\bar{f}(x) = -x_1^2 - 100x_2^2$ in Figure 2 is essentially the same function! We just scaled the variables (e.g. could correspond to different measurements units).

Solution: look at second order methods/ approximations. Indeed, if we use Newton’s method, which finds a step direction by considering the second order Taylor approximation, the resulting step directions always leads us directly to the minimum.

The gradient descent directly operates in the parameter space θ . However the same policy class π_θ can often be represented by various different parameterizations. Different parameterizations will lead to different gradients. The natural gradient method, described below, intends to get around this issue, by more directly optimizing in terms of the policy class (and distances within the policy class, i.e., distances between probability distributions) rather than the original parameters.

4 Natural gradients

We can approximate the derivative

$$\frac{\partial f}{\partial \theta_i} \approx \frac{f(\theta_i^0 + \Delta) - f(\theta_i^0 - \Delta)}{2\Delta} \tag{39}$$

where Δ is the “delta-width”. Note that the gradient computation normalizes by the distance traveled in θ space. However, θ is often an arbitrary way to index into the policy class, so it might be more natural to divide by a distance that is defined directly on the policy class, rather than on θ .

For example, rather than dividing by 2Δ , we could consider dividing by

$$KL(P_{\theta_i^0 + \Delta} || P_{\theta_i^0 - \Delta}) \tag{40}$$

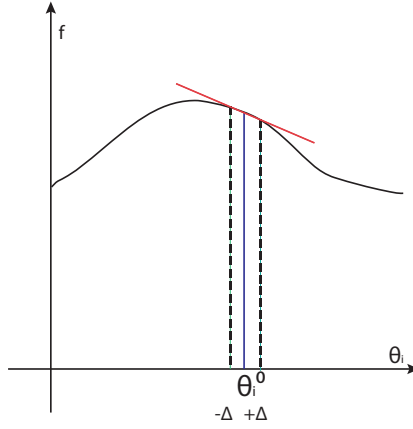


Figure 3: $f(\theta)$

The natural gradient essentially generalizes this idea, and finds the steepest ascent direction, when normalizing by a distance metric operating in the policy class space directly, rather than in the “arbitrary” θ space:

$$g_N = \arg \max_{\partial\theta_i^0: \|\delta\theta_i^0\|=\epsilon} \frac{f(\theta_i^0 + \delta\theta_i^0) - f(\theta_i^0 - \delta\theta_i^0)}{\text{distance}(\theta_i^0 + \delta\theta_i^0, \theta_i^0 - \delta\theta_i^0)} \quad (41)$$

When ϵ is “really small,” we can approximate the function f with a first-order Taylor expansion, and the distance metric by a 2nd order Taylor expansion:

$$f(\theta + \delta\theta) \approx f(\theta) + g_\theta^T \delta\theta \quad (42)$$

$$\text{distance}(\theta + \delta\theta, \theta - \delta\theta) \approx \sqrt{(\theta + \delta\theta - (\theta - \delta\theta))^T G_\theta (\theta + \delta\theta - (\theta - \delta\theta))} = 2\sqrt{\delta\theta^T G_\theta \delta\theta} \quad (43)$$

$$\Rightarrow g_N = \arg \max_{\partial\theta: \|\delta\theta\|=\epsilon} \frac{f(\theta) + g_\theta^T \delta\theta - (f(\theta) - g_\theta^T \delta\theta)}{2\sqrt{\delta\theta^T G_\theta \delta\theta}} \quad (44)$$

$$= \arg \max_{\partial\theta: \|\delta\theta\|=\epsilon} \frac{g_\theta^T \delta\theta}{\sqrt{\delta\theta^T G_\theta \delta\theta}} \quad (45)$$

$$\Rightarrow \boxed{\delta\theta = \alpha G_\theta^{-1} g_\theta}, \forall \alpha > 0 \quad (46)$$

The key is to pick a “clever” distance metric G_θ

$$G_\theta = I \Rightarrow g_N = \alpha \cdot g_\theta \quad (47)$$

A typical choice for distance metric between probability distributions would be a symmetric version of the KL-divergence:

$$\frac{1}{2}(KL(P\|Q) + KL(Q\|P)) = \frac{1}{2} \sum_x P(x) \log \frac{P(x)}{Q(x)} + \frac{1}{2} \sum_x Q(x) \log \frac{P(x)}{Q(x)} \quad (48)$$

For this choice, we have, for θ_1 and θ_2 close enough together:

$$KL(P_{\theta_1} || P_{\theta_2}) \approx (\theta_1 - \theta_2)^T G_{\theta_1} (\theta_1 - \theta_2) \quad (49)$$

For:

$$G_{\theta} = E_x \left[\frac{\partial \log P_{\theta}(x)}{\partial \theta} \frac{\partial \log P_{\theta}(x)}{\partial \theta}^T \right] \text{ (Fischer information matrix)} \quad (50)$$