

Notes on NEURAL ARCHITECTURE SEARCH WITH REINFORCEMENT LEARNING

Barret Zoph, Quoc V. Le

David Meyer

dmm@{1-4-5.net,uoregon.edu,brocade.com,...}

November 28, 2016

Notes

- $\pi_\theta(s, a) = P[A_t = a | S_t = s, \theta] = P_{(a_1:T;\theta_c)}$ (policy parameterized by θ)
- $J(\theta_c) = \mathbb{E}_{P_{(a_1:T;\theta_c)}}[R]$ (hook the RNN loss function to the RL reward)
- $J_1(\theta) = V^{\pi_\theta}(s_1) = \mathbb{E}_{\pi_\theta}[v_1] = \mathbb{E}_{P_{(a_1:T;\theta_c)}}[R]$ (episodic environments)
- $\nabla_{\theta_c} J(\theta_c) = \sum_{t=1}^T E_{P_{(a_1:T;\theta_c)}}[\nabla_{\theta_c} \log P(a_t | a_{(t-1):1}; \theta_c) R]$ (REINFORCE policy gradient)
- a_t is the predicted action (a) and $a_{(t-1):1}$ is the state s up to step $t - 1$ encoded in the RNN

1 Computing the gradient analytically

First, we assume that the policy π_θ is differentiable wherever it is non-zero (this is a softer requirement than requiring π_θ be differentiable *everywhere*). In addition, we know the gradient: $\nabla_\theta J(\theta)$. In this case, let $p(\mathbf{x}; \theta)$ be the likelihood parameterized by θ and let $\log p(\mathbf{x}; \theta)$ be the *log likelihood*. Then

$$y = p(\mathbf{x}; \theta) \quad \# \text{ definition; see above} \quad (1)$$

$$z = \log y = \log p(\mathbf{x}; \theta) \quad \# \text{ definition; } z \text{ is the log likelihood} \quad (2)$$

$$\frac{dz}{d\theta} = \frac{dz}{dy} \cdot \frac{dy}{d\theta} \quad \# \text{ chain rule definition} \quad (3)$$

$$\frac{dz}{dy} = \frac{1}{p(\mathbf{x}; \theta)} \quad \# \frac{\log(X)}{dX} \approx \frac{1}{X} \quad (4)$$

$$\frac{dy}{d\theta} = \frac{d p(\mathbf{x}; \theta)}{d\theta} = \nabla_{\theta} p(\mathbf{x}; \theta) \quad \# \text{ definition (chain rule, again)} \quad (5)$$

$$\frac{dz}{d\theta} = \frac{dz}{dy} \cdot \frac{dy}{d\theta} = \frac{\nabla_{\theta} p(\mathbf{x}; \theta)}{p(\mathbf{x}; \theta)} \quad \# \text{ chain rule} \quad (6)$$

$$= \nabla_{\theta} \log p(\mathbf{x}; \theta) \quad \# \text{ using the identity } \nabla_{\theta} \log(w) = \frac{1}{w} \nabla_{\theta} w \quad (7)$$

and setting $w = p(\mathbf{x}; \theta)$. Here $\nabla_{\theta} \log p(\mathbf{x}; \theta)$ is known as the score or sometimes the *Fischer* information. So the *log derivative trick* (sometimes *likelihood ratio*) is

$$\nabla_{\theta} \log p(\mathbf{x}; \theta) = \frac{\nabla_{\theta} p(\mathbf{x}; \theta)}{p(\mathbf{x}; \theta)}$$

Setting $\pi_{\theta}(s, a) = p(\mathbf{x}; \theta)$ we see that

$$\nabla_{\theta} \pi_{\theta}(s, a) = \pi_{\theta}(s, a) \frac{\nabla_{\theta} \pi_{\theta}(s, a)}{\pi_{\theta}(s, a)} \quad (8)$$

$$= \pi_{\theta}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a) \quad \# \text{ log derivative trick} \quad (9)$$

and the score function is $\nabla_{\theta} \log \pi_{\theta}(s, a)$.

Now, since here $\pi_{\theta}(s, a) = P(a_{1:T}; \theta_c)$, we have

$$\nabla_{\theta_c} J(\theta_c) = \sum_{s \in S} d(s) \sum_{a \in A} \nabla_{\theta_c} \pi_{\theta_c}(s, a) \mathcal{R}_{s,a} \quad \# \text{ defn policy gradient} \quad (10)$$

$$= \sum_{s \in S} d(s) \sum_{a \in A} \pi_{\theta_c}(s, a) \nabla \log \pi_{\theta_c}(s, a) \mathcal{R}_{s,a} \quad \# \text{ log derivative trick (Eqn 9)} \quad (11)$$

$$= \mathbb{E}_{\pi_{\theta_c}} [\nabla_{\theta_c} \log \pi_{\theta_c}(s, a) R] \quad \# \text{ defn expectation} \quad (12)$$

$$= \mathbb{E}_{a_t \sim P} [\nabla_{\theta_c} \log P(a_t | a_{(t-1):1}; \theta_c) R] \quad \# \pi_{\theta}(s, a) = P(a_{1:T}; \theta_c) \quad (13)$$

$$= \sum_{t=1}^T P_{(a_{1:T}; \theta_c)} [\nabla_{\theta_c} \log P(a_t | a_{(t-1):1}; \theta_c) R] \quad \# \text{ REINFORCE pg} \quad (14)$$