

# Lecture notes

## Part XV

# Policy Gradient (REINFORCE)

We will present a model-free algorithm called REINFORCE that does not require the notion of value functions and  $Q$  functions. It turns out to be more convenient to introduce REINFORCE in the finite horizon case, which will be assumed throughout this note: we use  $\tau = (s_0, a_0, \dots, s_{T-1}, a_{T-1}, s_T)$  to denote a trajectory, where  $T < \infty$  is the length of the trajectory. Moreover, REINFORCE only applies to learning a **randomized policy**. We use  $\pi_\theta(a|s)$  to denote the probability of the policy  $\pi_\theta$  outputting the action  $a$  at state  $s$ . The other notations will be the same as in previous lecture notes.

The advantage of applying REINFORCE is that we only need to assume that we can sample from the transition probabilities  $\{P_{sa}\}$  and can query the reward function  $R(s, a)$  at state  $s$  and action  $a$ ,<sup>1</sup> but we do not need to know the analytical form of the transition probabilities or the reward function. We do not explicitly learn the transition probabilities or the reward function either.

Let  $s_0$  be sampled from some distribution  $\mu$ . We consider optimizing the expected total payoff of the policy  $\pi_\theta$  over the parameter  $\theta$  defined as.

$$\eta(\theta) \triangleq \mathbb{E} \left[ \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t) \right] \quad (1)$$

Recall that  $s_t \sim P_{s_{t-1}a_{t-1}}$  and  $a_t \sim \pi_\theta(\cdot|s_t)$ . Also note that  $\eta(\theta) = \mathbb{E}_{s_0 \sim P} [V^{\pi_\theta}(s_0)]$  if we ignore the difference between finite and infinite horizon.

---

<sup>1</sup>In this notes we will work with the general setting where the reward depends on both the state and the action.

We aim to use gradient ascent to maximize  $\eta(\theta)$ . The main challenge we face here is to compute (or estimate) the gradient of  $\eta(\theta)$  without the knowledge of the form of the reward function and the transition probabilities.

Let  $P_\theta(\tau)$  denote the distribution of  $\tau$  (generated by the policy  $\pi_\theta$ ), and let  $f(\tau) = \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t)$ . We can rewrite  $\eta(\theta)$  as

$$\eta(\theta) = \mathbb{E}_{\tau \sim P_\theta} [f(\tau)] \quad (2)$$

We face a similar situations in the variational auto-encoder (VAE) setting covered in the previous lectures, where we need to take the gradient w.r.t to a variable that shows up under the expectation — the distribution  $P_\theta$  depends on  $\theta$ . Recall that in VAE, we used the re-parametrization techniques to address this problem. However it does not apply here because we do not know how to compute the gradient of the function  $f$ . (We only have an efficient way to evaluate the function  $f$  by taking a weighted sum of the observed rewards, but we do not necessarily know the reward function itself to compute the gradient.)

The REINFORCE algorithm uses an another approach to estimate the gradient of  $\eta(\theta)$ . We start with the following derivation:

$$\begin{aligned} \nabla_\theta \mathbb{E}_{\tau \sim P_\theta} [f(\tau)] &= \nabla_\theta \int P_\theta(\tau) f(\tau) d\tau \\ &= \int \nabla_\theta (P_\theta(\tau) f(\tau)) d\tau \quad (\text{swap integration with gradient}) \\ &= \int (\nabla_\theta P_\theta(\tau)) f(\tau) d\tau \quad (\text{because } f \text{ does not depend on } \theta) \\ &= \int P_\theta(\tau) (\nabla_\theta \log P_\theta(\tau)) f(\tau) d\tau \\ &\quad (\text{because } \nabla \log P_\theta(\tau) = \frac{\nabla P_\theta(\tau)}{P_\theta(\tau)}) \\ &= \mathbb{E}_{\tau \sim P_\theta} [(\nabla_\theta \log P_\theta(\tau)) f(\tau)] \end{aligned} \quad (3)$$

Now we have a sample-based estimator for  $\nabla_\theta \mathbb{E}_{\tau \sim P_\theta} [f(\tau)]$ . Let  $\tau^{(1)}, \dots, \tau^{(n)}$  be  $n$  empirical samples from  $P_\theta$  (which are obtained by running the policy  $\pi_\theta$  for  $n$  times, with  $T$  steps for each run). We can estimate the gradient of  $\eta(\theta)$  by

$$\nabla_\theta \mathbb{E}_{\tau \sim P_\theta} [f(\tau)] = \mathbb{E}_{\tau \sim P_\theta} [(\nabla_\theta \log P_\theta(\tau)) f(\tau)] \quad (4)$$

$$\approx \frac{1}{n} \sum_{i=1}^n (\nabla_\theta \log P_\theta(\tau^{(i)})) f(\tau^{(i)}) \quad (5)$$

The next question is how to compute  $\log P_\theta(\tau)$ . We derive an analytical formula for  $\log P_\theta(\tau)$  and compute its gradient w.r.t  $\theta$  (using auto-differentiation). Using the definition of  $\tau$ , we have

$$P_\theta(\tau) = \mu(s_0)\pi_\theta(a_0|s_0)P_{s_0a_0}(s_1)\pi_\theta(a_1|s_1)P_{s_1a_1}(s_2)\cdots P_{s_{T-1}a_{T-1}}(s_T) \quad (6)$$

Here recall that  $\mu$  is used to denote the density of the distribution of  $s_0$ . It follows that

$$\begin{aligned} \log P_\theta(\tau) &= \log \mu(s_0) + \log \pi_\theta(a_0|s_0) + \log P_{s_0a_0}(s_1) + \log \pi_\theta(a_1|s_1) \\ &\quad + \log P_{s_1a_1}(s_2) + \cdots + \log P_{s_{T-1}a_{T-1}}(s_T) \end{aligned} \quad (7)$$

Taking gradient w.r.t to  $\theta$ , we obtain

$$\nabla_\theta \log P_\theta(\tau) = \nabla_\theta \log \pi_\theta(a_0|s_0) + \nabla_\theta \log \pi_\theta(a_1|s_1) + \cdots + \nabla_\theta \log \pi_\theta(a_{T-1}|s_{T-1})$$

Note that many of the terms disappear because they don't depend on  $\theta$  and thus have zero gradients. (This is somewhat important — we don't know how to evaluate those terms such as  $\log P_{s_0a_0}(s_1)$  because we don't have access to the transition probabilities, but luckily those terms have zero gradients!)

Plugging the equation above into equation (4), we conclude that

$$\begin{aligned} \nabla_\theta \eta(\theta) &= \nabla_\theta \mathbb{E}_{\tau \sim P_\theta} [f(\tau)] = \mathbb{E}_{\tau \sim P_\theta} \left[ \left( \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t) \right) \cdot f(\tau) \right] \\ &= \mathbb{E}_{\tau \sim P_\theta} \left[ \left( \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t) \right) \cdot \left( \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t) \right) \right] \end{aligned} \quad (8)$$

We estimate the RHS of the equation above by empirical sample trajectories, and the estimate is unbiased. The vanilla REINFORCE algorithm iteratively updates the parameter by gradient ascent using the estimated gradients.

**Interpretation of the policy gradient formula (8).** The quantity  $\nabla_\theta P_\theta(\tau) = \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t)$  is intuitively the direction of the change of  $\theta$  that will make the trajectory  $\tau$  more likely to occur (or increase the probability of choosing action  $a_0, \dots, a_{t-1}$ ), and  $f(\tau)$  is the total payoff of this trajectory. Thus, by taking a gradient step, intuitively we are trying to improve the likelihood of all the trajectories, but with a different emphasis or weight for each  $\tau$  (or for each set of actions  $a_0, a_1, \dots, a_{t-1}$ ). If  $\tau$  is very rewarding (that is,  $f(\tau)$  is large), we try very hard to move in the direction

that can increase the probability of the trajectory  $\tau$  (or the direction that increases the probability of choosing  $a_0, \dots, a_{t-1}$ ), and if  $\tau$  has low payoff, we try less hard with a smaller weight.

An interesting fact that follows from formula (3) is that

$$\mathbb{E}_{\tau \sim P_\theta} \left[ \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) \right] = 0 \quad (9)$$

To see this, we take  $f(\tau) = 1$  (that is, the reward is always a constant), then the LHS of (8) is zero because the payoff is always a fixed constant  $\sum_{t=0}^T \gamma^t$ . Thus the RHS of (8) is also zero, which implies (9).

In fact, one can verify that  $\mathbb{E}_{a_t \sim \pi_\theta(\cdot | s_t)} \nabla_\theta \log \pi_\theta(a_t | s_t) = 0$  for any fixed  $t$  and  $s_t$ .<sup>2</sup> This fact has two consequences. First, we can simplify formula (8) to

$$\begin{aligned} \nabla_\theta \eta(\theta) &= \sum_{t=0}^{T-1} \mathbb{E}_{\tau \sim P_\theta} \left[ \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \left( \sum_{j=0}^{T-1} \gamma^j R(s_j, a_j) \right) \right] \\ &= \sum_{t=0}^{T-1} \mathbb{E}_{\tau \sim P_\theta} \left[ \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \left( \sum_{j \geq t}^{T-1} \gamma^j R(s_j, a_j) \right) \right] \end{aligned} \quad (10)$$

where the second equality follows from

$$\begin{aligned} &\mathbb{E}_{\tau \sim P_\theta} \left[ \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot \left( \sum_{0 \leq j < t} \gamma^j R(s_j, a_j) \right) \right] \\ &= \mathbb{E} \left[ \mathbb{E} [\nabla_\theta \log \pi_\theta(a_t | s_t) | s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t] \cdot \left( \sum_{0 \leq j < t} \gamma^j R(s_j, a_j) \right) \right] \\ &= 0 \quad (\text{because } \mathbb{E} [\nabla_\theta \log \pi_\theta(a_t | s_t) | s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t] = 0) \end{aligned}$$

Note that here we used the law of total expectation. The outer expectation in the second line above is over the randomness of  $s_0, a_0, \dots, a_{t-1}, s_t$ , whereas the inner expectation is over the randomness of  $a_t$  (conditioned on  $s_0, a_0, \dots, a_{t-1}, s_t$ .) We see that we've made the estimator slightly simpler. The second consequence of  $\mathbb{E}_{a_t \sim \pi_\theta(\cdot | s_t)} \nabla_\theta \log \pi_\theta(a_t | s_t) = 0$  is the following: for any value  $B(s_t)$  that only depends on  $s_t$ , it holds that

$$\begin{aligned} &\mathbb{E}_{\tau \sim P_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t) \cdot B(s_t)] \\ &= \mathbb{E} [\mathbb{E} [\nabla_\theta \log \pi_\theta(a_t | s_t) | s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t] B(s_t)] \\ &= 0 \quad (\text{because } \mathbb{E} [\nabla_\theta \log \pi_\theta(a_t | s_t) | s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t] = 0) \end{aligned}$$

---

<sup>2</sup>In general, it's true that  $\mathbb{E}_{x \sim p_\theta} [\nabla \log p_\theta(x)] = 0$ .

Again here we used the law of total expectation. The outer expectation in the second line above is over the randomness of  $s_0, a_0, \dots, a_{t-1}, s_t$ , whereas the inner expectation is over the randomness of  $a_t$  (conditioned on  $s_0, a_0, \dots, a_{t-1}, s_t$ .) It follows from equation (10) and the equation above that

$$\begin{aligned}\nabla_{\theta}\eta(\theta) &= \sum_{t=0}^{T-1} \mathbb{E}_{\tau \sim P_{\theta}} \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot \left( \sum_{j \geq t}^{T-1} \gamma^j R(s_j, a_j) - \gamma^t B(s_t) \right) \right] \\ &= \sum_{t=0}^{T-1} \mathbb{E}_{\tau \sim P_{\theta}} \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot \gamma^t \left( \sum_{j \geq t}^{T-1} \gamma^{j-t} R(s_j, a_j) - B(s_t) \right) \right]\end{aligned}\tag{11}$$

Therefore, we will get a different estimator for estimating the  $\nabla\eta(\theta)$  with a difference choice of  $B(\cdot)$ . The benefit of introducing a proper  $B(\cdot)$  — which is often referred to as a **baseline** — is that it helps reduce the variance of the estimator.<sup>3</sup> It turns out that a near optimal estimator would be the expected future payoff  $\mathbb{E} \left[ \sum_{j \geq t}^{T-1} \gamma^{j-t} R(s_j, a_j) | s_t \right]$ , which is pretty much the same as the value function  $V^{\pi_{\theta}}(s_t)$  (if we ignore the difference between finite and infinite horizon.) Here one could estimate the value function  $V^{\pi_{\theta}}(\cdot)$  in a crude way, because its precise value doesn't influence the mean of the estimator but only the variance. This leads to a policy gradient algorithm with baselines stated in Algorithm 1.<sup>4</sup>

---

<sup>3</sup>As a heuristic but illustrating example, suppose for a fixed  $t$ , the future reward  $\sum_{j \geq t}^{T-1} \gamma^{j-t} R(s_j, a_j)$  randomly takes two values  $1000 + 1$  and  $1000 - 2$  with equal probability, and the corresponding values for  $\nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$  are vector  $z$  and  $-z$ . (Note that because  $\mathbb{E}[\nabla_{\theta} \log \pi_{\theta}(a_t | s_t)] = 0$ , if  $\nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$  can only take two values uniformly, then the two values have to be two vectors in an opposite direction.) In this case, without subtracting the baseline, the estimators take two values  $(1000 + 1)z$  and  $-(1000 - 2)z$ , whereas after subtracting a baseline of 1000, the estimator has two values  $z$  and  $2z$ . The latter estimator has much lower variance compared to the original estimator.

<sup>4</sup>We note that the estimator of the gradient in the algorithm does not exactly match the equation 11. If we multiply  $\gamma^t$  in the summand of equation (13), then they will exactly match. Removing such discount factors empirically works well because it gives a large update.

---

**Algorithm 1** Vanilla policy gradient with baseline

---

**for**  $i = 1, \dots$  **do**

Collect a set of trajectories by executing the current policy. Use  $R_{\geq t}$  as a shorthand for  $\sum_{j \geq t}^{T-1} \gamma^{j-t} R(s_j, a_j)$

Fit the baseline by finding a function  $B$  that minimizes

$$\sum_{\tau} \sum_t (R_{\geq t} - B(s_t))^2 \quad (12)$$

Update the policy parameter  $\theta$  with the gradient estimator

$$\sum_{\tau} \sum_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot (R_{\geq t} - B(s_t)) \quad (13)$$


---