

## Lecture 5: Policy gradient methods

*By Shipra Agrawal*

In Q-learning function approximation was used to approximate Q-function, and policy was a greedy policy based on estimated Q-function. In policy gradient methods, we approximate a stochastic policy directly using a parametric function approximator.

More formally, given an MDP  $(S, A, s_1, R, P)$ , let  $\pi_\theta : S \rightarrow \Delta^A$  denote a randomized policy parameterized by parameter vector  $\theta \in \mathbb{R}^d$ . For a scalable formulation, we want  $d \ll |S|$ . For example, the policy  $\pi_\theta$  might be represented by a neural network whose input is a representation of the state, whose output is action selection probabilities, and whose weights form the policy parameters  $\theta$ . (The architecture for such a [deep] neural network is similar to that for a multi-label classifier, with input being a state, and labels being different actions. And, the network should be trained to predict the probability of different actions given an input state).

For simplicity, assume that  $\pi_\theta$  is differentiable with respect to  $\theta$ , i.e.  $\frac{\partial \pi_\theta(s, a)}{\partial \theta}$  exists. This is true for example, if a neural network with differentiable activation functions is used to define  $\pi_\theta$ . Let  $\rho^{\pi_\theta}$  denote the gain of policy  $\pi_\theta$ . This may be defined as long term average reward, long term discounted reward or total reward in an episode or finite horizon. Therefore, solving for optimal policy reduces to the problem of solving

$$\max_{\theta} \rho^{\pi_\theta}(s_1)$$

In order to use stochastic gradient descent algorithm for finding a stationary point of the above problem, we need to compute (an unbiased) estimate of gradient of  $\rho^{\pi_\theta}(s_1)$  with respect to  $\theta$ .

## 1 Finite horizon MDP

Here performance measure to optimize is total expected reward over a finite horizon  $H$ .

$$\rho^\pi(s_1) = \mathbb{E}[\sum_{t=1}^H \gamma^{t-1} r_t | \pi, s_1]$$

Let  $\pi(s, a)$  denote the probability of action  $a$  in state  $s$  for randomized policy  $\pi$ . Let  $D^\pi(\tau)$  denote the probability distribution of a trajectory (state-action sequence)  $\tau = (s_1, a_1, s_2, \dots, a_{H-1}, s_H)$  of states on starting from state  $s_1$  and following policy  $\pi$ . That is,

$$D^\pi(\tau) := \prod_{i=1}^{H-1} \pi(s_i, a_i) P(s_i, a_i, s_{i+1})$$

**Theorem 1.** For finite horizon MDP  $(S, A, s_1, P, R, H)$ , let  $R(\tau)$  be the total reward for an sample trajectory  $\tau$ , on following  $\pi_\theta$  for  $H$  steps, starting from state  $s_1$ . Then,

$$\nabla_{\theta} \rho^{\pi_\theta}(s_1) = \mathbb{E}_{\tau} [R(\tau) \nabla_{\theta} \log(D^{\pi_\theta}(\tau)) | s_1] = \mathbb{E}_{\tau} \left[ R(\tau) \sum_{t=1}^{H-1} \nabla_{\theta} \log(\pi_{\theta}(s_t, a_t)) | s_1 \right]$$

*Proof.* Let  $R(\tau)$  be expected total reward for an entire sample trajectory  $\tau$ , on following  $\pi_\theta$  for  $H$  steps, starting from state  $s_1$ . That is, given a sample trajectory  $\tau = (s_1, a_1, s_2, \dots, a_{H-1}, s_H)$  from distribution  $D^{\pi_\theta}$ ,

$$R(\tau) := \sum_{t=1}^{H-1} \gamma^{t-1} R(s_t, a_t),$$

Then,

$$\rho^{\pi_\theta}(s) = \mathbb{E}_{\tau \sim D^{\pi_\theta}}[R(\tau)|s_1 = s]$$

Now, (the calculations below implicitly assume finite state and action space, so that the distribution  $D^{\pi_\theta}(\tau)$  has a finite support)

$$\begin{aligned} \nabla_\theta \rho^{\pi_\theta}(s) &= \nabla_\theta \mathbb{E}_{\tau \sim D^{\pi_\theta}}[R(\tau)|s_1 = s] \\ &= \nabla_\theta \sum_{\tau: D^{\pi_\theta}(\tau) > 0} D^{\pi_\theta}(\tau|s_1 = s) R(\tau) \\ &= \sum_{\tau: D^{\pi_\theta}(\tau) > 0} D^{\pi_\theta}(\tau|s_1 = s) \nabla_\theta \log(D^{\pi_\theta}(\tau|s_1 = s)) R(\tau) \\ &= \mathbb{E}_{\tau \sim D^{\pi_\theta}}[\nabla_\theta \log(D^{\pi_\theta}(\tau)) R(\tau)|s_1 = s] \end{aligned}$$

Further, for a given sample trajectory  $\tau^i$ .

$$\begin{aligned} \nabla_\theta \log(D^{\pi_\theta}(\tau^i)) &= \sum_{t=1}^{H-1} \nabla_\theta \log(\pi_\theta(s_t^i, a_t^i)) + \nabla_\theta \log P(s_t^i, a_t^i, s_{t+1}^i) \\ &= \sum_{t=1}^{H-1} \nabla_\theta \log(\pi_\theta(s_t^i, a_t^i)) \end{aligned}$$

□

The gradient representation given by above theorem is extremely useful, as given a sample trajectory this can be computed only using the policy parameter, and does not require knowledge of the transition model  $P(\cdot, \cdot, \cdot)$ ! This does seem to require knowledge of reward model, but that can be handled by replacing  $R(\tau^i)$  by  $\hat{R}(\tau^i) = r_1 + \gamma r_2 + \dots, \gamma^{H-2} r_{H-1}$ , the total of sample rewards observed in this trajectory. Since, given a trajectory  $\tau$ , the quantity  $D^{\pi_\theta}(\tau)$  is determined, and  $\mathbb{E}[\hat{R}(\tau)|\tau] = R(\tau)$ ,

$$\begin{aligned} \nabla_\theta \rho^{\pi_\theta}(s_1) &= \mathbb{E}_\tau[R(\tau) \nabla_\theta \log(D^{\pi_\theta}(\tau))|s_1] \\ &= \mathbb{E}_\tau[\hat{R}(\tau) \nabla_\theta \log(D^{\pi_\theta}(\tau))|s_1] \\ &= \mathbb{E}_\tau\left[\hat{R}(\tau) \sum_{t=1}^{H-1} \nabla_\theta \log(\pi_\theta(s_t, a_t))|s_1\right] \end{aligned}$$

**Unbiased estimator of gradient from samples.** From above, given sample trajectories  $\tau^i, i = 1, \dots, m$ , an unbiased estimator for gradient  $\nabla_\theta \rho(\pi_\theta)$  is given as:

$$\hat{\mathbf{g}} = \frac{1}{m} \sum_{i=1}^m \hat{R}(\tau^i) \nabla_\theta \log(D^{\pi_\theta}(\tau^i)) = \frac{1}{m} \sum_{i=1}^m \hat{R}(\tau^i) \sum_{t=1}^{H-1} \nabla_\theta \log(\pi_\theta(s_t^i, a_t^i)) \quad (1)$$

**Baseline.** Note that for any constant  $b$  ( or  $b$  that is conditionally independent of sampling from  $\pi_\theta$  given  $\theta$ ), we have:

$$\mathbb{E}_\tau[b \frac{\partial}{\partial \theta_j} \log(D^{\pi_\theta}(\tau)) | \theta, s_1] = b \int_\tau \frac{\partial}{\partial \theta_j} D^{\pi_\theta}(\tau) = b \frac{\partial}{\partial \theta_j} \int_\tau D^{\pi_\theta}(\tau) = 0$$

Therefore, choosing any ‘baseline’  $b$ , following is also an unbiased estimator of the  $\nabla_\theta \rho(\pi_\theta)$ :

$$\hat{\mathbf{g}} = \frac{1}{m} \sum_{i=1}^m \sum_{t=1}^{H-1} (\hat{R}(\tau^i) - b) \nabla_\theta \log(\pi_\theta(s_t^i, a_t^i))$$

Or, more generally, one could even use a state and time dependent baseline  $b_t(s_t^i)$  conditionally is independent of sampling from  $\pi_\theta$  given  $s_t^i, \theta$ , to get estimator:

$$\hat{\mathbf{g}} = \frac{1}{m} \sum_{i=1}^m \sum_{t=1}^{H-1} (\hat{R}(\tau^i) - b_t(s_t^i)) \nabla_\theta \log(\pi_\theta(s_t^i, a_t^i)) \quad (2)$$

Below we show this is unbiased. The expectations below are over trajectories  $(s_1, a_1, \dots, a_{H-1}, s_H)$ , where  $a_t \sim \pi(s_t, \cdot)$ , given  $s_t$ . For any fixed  $\theta, t$ , the baseline  $b_t(s_t)|_{s_t}$  needs to be deterministic or independent of  $a_t|s_t$ . For simplicity we assume it is deterministic.

$$\begin{aligned} \mathbb{E}_\tau \left[ \sum_{t=1}^{H-1} b_t(s_t) \frac{\partial}{\partial \theta_j} \log(\pi_\theta(s_t, a_t)) | \theta, s_1 \right] &= \mathbb{E} \left[ \sum_{t=1}^{H-1} \mathbb{E}[b_t(s_t) \frac{\partial}{\partial \theta_j} \log(\pi_\theta(s_t, a_t)) | s_t] | \theta, s_1 \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^{H-1} b_t(s_t) \mathbb{E} \left[ \frac{\partial}{\partial \theta_j} \log(\pi_\theta(s_t, a_t)) | s_t \right] | \theta, s_1 \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^{H-1} b_t(s_t) \sum_a \pi_\theta(s_t, a) \frac{\partial}{\partial \theta_j} \log(\pi_\theta(s_t, a)) | \theta, s_1 \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^{H-1} b_t(s_t) \sum_a \frac{\partial}{\partial \theta_j} \pi_\theta(s_t, a) | \theta, s_1 \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^{H-1} b_t(s_t) \frac{\partial}{\partial \theta_j} \sum_a \pi_\theta(s_t, a) | \theta, s_1 \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^{H-1} b_t(s_t) \frac{\partial}{\partial \theta_j} (1) | \theta, s_1 \right] \\ &= 0 \end{aligned}$$

An example of such state dependent baseline  $b_t(s)$ , given  $s$  and  $\theta$ , is  $V_{H-t}^{\pi_\theta}(s)$ , i.e., the value of policy  $\pi_\theta$ , starting from state  $s$  at time  $t$ . We will see later that such a baseline is useful in reducing the variance of gradient estimates.

**Vanilla policy gradient algorithm** Initialize policy parameter  $\theta$ , and baseline.  
In each iteration,

- Execute current policy  $\pi^\theta$  to obtain several sample trajectories  $\tau^i, i = 1, \dots, m$ .
- Use these sample trajectories and chosen baseline to compute the gradient estimator  $\hat{\mathbf{g}}$  as in (2).
- Update  $\theta \leftarrow \theta + \alpha \hat{\mathbf{g}}$
- Update baseline as required.

Above is essentially same as the **REINFORCE** algorithm introduced by [Williams, 1988, 1992].

## 2 Infinite horizon average reward case

Let's consider the case when  $\rho^\pi(s)$  is long-term average reward of (randomized) policy  $\pi$ .  $\rho^\pi(s)$  is assumed to exist and independent of the starting state  $s_1$  for all policies  $\pi$  (refer to the introductory lecture notes for conditions under which this assumption holds). Therefore,

$$\rho^\pi(s_1) = \rho^\pi = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}[r_1 + r_2 + \dots + r_T | \pi] = \sum_s d^\pi(s) \sum_a \pi(s, a) R(s, a)$$

where  $d^\pi(s) = \lim_{t \rightarrow \infty} \Pr(s_t = s | s_1, \pi)$ , the stationary distribution for policy  $\pi$ . In the average reward case, recall that the bias of a state given a policy  $\pi$  is defined as:

$$h^\pi(s) = \lim_{T \rightarrow \infty} \mathbb{E}[\sum_{t=1}^T (r_t - \rho(\pi)) | s_1 = s, a_t = \pi(s_t), t = 1, \dots, T] = \sum_a \pi(s_1, a) A^\pi(s_1, a)$$

where

$$A^\pi(s, a) := \lim_{T \rightarrow \infty} \mathbb{E}[\sum_{t=1}^T (r_t - \rho(\pi)) | s_1 = s, a_1 = a, a_t = \pi(s_t), t = 2, \dots, T]$$

**Theorem 2. [Policy gradient theorem [Sutton et al., 1999]]** *For infinite horizon MDP average reward case,*

$$\nabla_\theta \rho(\pi_\theta, s_1) = \sum_s d^{\pi_\theta}(s) \sum_a A^{\pi_\theta}(s, a) \nabla_\theta \pi_\theta(s, a) = \sum_s d^{\pi_\theta}(s) (\mathbb{E}_{a \sim \pi(s)} [A^{\pi_\theta}(s, a) \nabla_\theta \log(\pi_\theta(s, a))])$$

*That is gradient of gain with respect to  $\theta$  can be expressed in terms of gradient of policy function with respect to  $\theta$ .*

*Proof.*

$$\begin{aligned} A^\pi(s, a) &= R(s, a) + \sum_{s'} P(s, a, s') V^\pi(s') - \rho(\pi) \\ V^\pi(s) &= \sum_a \pi(s, a) Q^\pi(s, a) \\ \nabla_\theta V^\pi(s) &= \sum_a Q^\pi(s, a) \nabla_\theta \pi(s, a) + \pi(s, a) \nabla_\theta A^\pi(s, a) \\ &= \sum_a A^\pi(s, a) \nabla_\theta \pi(s, a) + \sum_a \pi(s, a) \left[ \sum_{s'} P(s, a, s') \nabla_\theta V^\pi(s') - \nabla_\theta \rho(\pi) \right] \\ \nabla_\theta \rho(\pi) &= \sum_a A^\pi(s, a) \nabla_\theta \pi(s, a) + \sum_a \pi(s, a) \left[ \sum_{s'} P(s, a, s') \nabla_\theta V^\pi(s') \right] - \nabla_\theta V^\pi(s) \\ \sum_s d^\pi(s) \nabla_\theta \rho(\pi) &= \sum_s d^\pi(s) \sum_a A^\pi(s, a) \nabla_\theta \pi(s, a) + \sum_s d^\pi(s) \sum_a \pi(s, a) \left[ \sum_{s'} P(s, a, s') \nabla_\theta V^\pi(s') \right] - \sum_s d^\pi(s) \nabla_\theta V^\pi(s) \\ &\quad \text{since } d^\pi \text{ is the stationary distribution of } \pi \\ &= \sum_s d^\pi(s) \sum_a A^\pi(s, a) \nabla_\theta \pi(s, a) + \sum_{s'} d^\pi(s') \nabla_\theta V^\pi(s') - \sum_s d^\pi(s) \nabla_\theta V^\pi(s) \\ &= \sum_s d^\pi(s) \sum_a A^\pi(s, a) \nabla_\theta \pi(s, a) \end{aligned}$$

□

### 3 Infinite horizon discounted rewards

$$\rho^\pi(s_1) = \lim_{T \rightarrow \infty} \mathbb{E}[\sum_{t=1}^T \gamma^{t-1} r_t | \pi, s_1] = \sum_s d^\pi(s) \sum_a \pi(s, a) R(s, a)$$

where  $d^\pi(s) = \lim_{T \rightarrow \infty} \sum_{t=1}^T \gamma^{t-1} \Pr(s_t = s | s_1, \pi)$ , and the value of a state given a policy  $\pi$  is defined as:

$$V^\pi(s) = \lim_{T \rightarrow \infty} \mathbb{E}[\sum_{t=1}^T \gamma^{t-1} r_t | s_1 = s, a_t = \pi(s_t), t = 1, \dots, T] = \sum_a \pi(s_1, a) Q^\pi(s_1, a)$$

where

$$Q^\pi(s, a) := \lim_{T \rightarrow \infty} \mathbb{E} \left[ \sum_{t=1}^T \gamma^{t-1} r_t \mid s_1 = s, a_1 = a, a_t = \pi(s_t), t = 2, \dots, T \right]$$

**Theorem 3.** [Policy gradient theorem [Sutton et al., 1999]] *For infinite horizon MDP discounted reward case,*

$$\nabla_\theta \rho(\pi_\theta, s_1) = \sum_s d^{\pi_\theta}(s) \sum_a Q^{\pi_\theta}(s, a) \nabla_\theta \pi_\theta(s, a) = \sum_s d^{\pi_\theta}(s) \left( \mathbb{E}_{a \sim \pi(s)} [Q^{\pi_\theta}(s, a) \nabla_\theta \log(\pi_\theta(s, a))] \right)$$

*That is gradient of gain with respect to  $\theta$  can be expressed in terms of gradient of policy function with respect to  $\theta$ .*

*Proof.*

$$Q^\pi(s, a) = R(s, a) + \gamma \sum_{s'} P(s, a, s') V^\pi(s')$$

By Bellman equations:

$$\begin{aligned} V^\pi(s) &= \sum_a \pi(s, a) Q^\pi(s, a) \\ \nabla_\theta V^\pi(s) &= \sum_a Q^\pi(s, a) \nabla_\theta \pi(s, a) + \pi(s, a) \nabla_\theta Q^\pi(s, a) \end{aligned}$$

Let  $\Pr(s \rightarrow x, k, \pi)$  is the probability of going from state  $s$  to state  $x$  in  $k$  steps under policy  $\pi$ .

$$\begin{aligned} \sum_s d^\pi(s) \nabla_\theta V^\pi(s) &:= \sum_s \left( \sum_{k=0}^{\infty} \gamma^k \Pr(s_1 \rightarrow s, k, \pi) \right) \nabla_\theta V^\pi(s) \\ &= \sum_s \left( \sum_{k=0}^{\infty} \gamma^k \Pr(s_1 \rightarrow s, k, \pi) \right) \left( \sum_a Q^\pi(s, a) \nabla_\theta \pi(s, a) + \gamma \sum_a \pi(s, a) \sum_{s'} P(s, a, s') \nabla_\theta V^\pi(s') \right) \\ &= \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla_\theta \pi(s, a) + \sum_{s'} \left( \sum_s \sum_{k=0}^{\infty} \gamma^k \Pr(s_1 \rightarrow s, k, \pi) \right) \gamma \sum_a \pi(s, a) P(s, a, s') \nabla_\theta V^\pi(s') \\ &= \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla_\theta \pi(s, a) + \sum_{s'} \left( \sum_{k=0}^{\infty} \gamma^{k+1} \Pr(s_1 \rightarrow s', k+1, \pi) \right) \nabla_\theta V^\pi(s') \\ &= \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla_\theta \pi(s, a) + \sum_{s'} d^\pi(s') - \Pr(s_1 \rightarrow s', 0, \pi) \nabla_\theta V^\pi(s') \\ &= \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla_\theta \pi(s, a) + \sum_{s'} d^\pi(s') \nabla_\theta V^\pi(s') - \nabla_\theta V^\pi(s_1) \end{aligned}$$

Moving the terms around

$$\nabla_\theta V^\pi(s_1) = \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla_\theta \pi(s, a)$$

That is,

$$\nabla_\theta \rho(\pi, s_1) = \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla_\theta \pi(s, a)$$

□

**Remarks on Theorem 3.** The key aspect of the expression for the gradient is that there are no terms of the form  $\frac{\partial d^{\pi_\theta}(s)}{\partial \theta}$ : the effect of policy changes on the distribution of states does not appear. This is convenient for approximating the gradient by sampling. For example, if  $s$  was sampled from the distribution obtained by following  $\pi_\theta$ , then  $\sum_a Q^{\pi_\theta}(s, a) \nabla_\theta \pi_\theta(s, a)$  would be an unbiased estimate of  $\nabla_\theta \rho(\pi_\theta)$ . Of course,  $Q^{\pi_\theta}(s, a)$  is also not normally known and must be estimated (in an unbiased way).

### 3.1 Vanilla Policy Gradient Algorithm (for discounted case)

**Estimation using samples.** Suppose we run policy  $\pi$  several times starting from  $s_1$  to observe sample trajectories  $\{\tau^i\}$ . Then, at each time step  $t$  in a trajectory  $\tau$ , set

$$\hat{Q}_t := \sum_{t'=t}^{T-1} \gamma^{t'-t} r_{t'}$$

where  $r_{t'}$  is the observed reward at time  $t'$ . Then,  $\hat{Q}_t$  is an unbiased estimate of  $Q(s_t, a_t)$  at time  $t$  (almost, assuming large enough  $T$ ), i.e.,

$$\mathbb{E}[\hat{Q}_t | s_t, a_t] = Q(s_t, a_t)$$

Let

$$F_t := \hat{Q}_t \nabla_\theta \log \pi_\theta(s_t, a_t)$$

Then,

$$\mathbb{E}[F_t | s_t] = \mathbb{E}[\mathbb{E}[\hat{Q}_t | s_t, a_t] \nabla_\theta \log(\pi_\theta(s_t, a_t)) | s_t] = \mathbb{E}[Q(s_t, a_t) \nabla_\theta \log(\pi_\theta(s_t, a_t)) | s_t] = \sum_a Q(s_t, a) \nabla_\theta \pi_\theta(s_t, a)$$

And,

$$\begin{aligned} \mathbb{E}\left[\sum_{t=1}^{\infty} \gamma^{t-1} F_t\right] &= \sum_{t=1}^{\infty} \gamma^{t-1} \sum_s \mathbb{E}[F_t | s_t = s] \Pr(s_t = s | s_1) \\ &= \sum_{t=1}^{\infty} \gamma^{t-1} \sum_s \left( \sum_a Q(s, a) \nabla_\theta \pi_\theta(s, a) \right) \Pr(s_t = s | s_1) \\ &= \sum_s \left( \sum_a Q(s, a) \nabla_\theta \pi_\theta(s, a) \right) \sum_{t=1}^{\infty} \gamma^{t-1} \Pr(s_t = s | s_1) \\ &= \sum_s \left( \sum_a Q(s, a) \nabla_\theta \pi_\theta(s, a) \right) d^\pi(s) \\ &= \nabla_\theta \rho(\pi_\theta) \end{aligned}$$

where the last step follows from the policy gradient theorem. Therefore, following is an unbiased (almost, for large  $T$ ) estimate of gradient of  $\nabla_\theta \rho(\pi, s_1)$  of policy  $\pi$  at  $\hat{\theta}$ , starting in state  $s_1$ :

$$\hat{\mathbf{g}} = \sum_{t=1}^T \gamma^{t-1} F_t = \sum_{t=1}^T \gamma^{t-1} \hat{Q}_t \nabla_\theta \log(\pi_\theta(s_t, a_t))$$

From  $N$  sample trajectories  $\{\tau^i, i = 1, \dots, N\}$  we can develop sample average estimate

$$\hat{\mathbf{g}} = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \gamma^{t-1} F_t^i = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \gamma^{t-1} \hat{Q}_t^i \nabla_\theta \log(\pi_\theta(s_t^i, a_t^i)) \quad (3)$$

**Baseline.** We can obtain another unbiased estimate by replacing  $F_t$  by

$$F'_t = (\hat{Q}_t - b_t(s_t)) \nabla_{\theta} \log(\pi_{\theta}(s_t, a_t))$$

for an arbitrary baseline function  $b_t(s)$ . Then,

$$\mathbb{E}[\sum_t \gamma^{t-1} F'_t] = \mathbb{E}[\sum_t \gamma^{t-1} F_t] - \mathbb{E}[\sum_t \gamma^{t-1} b_t(s_t) \nabla_{\theta} \log(\pi_{\theta}(s_t, a_t))] = \mathbb{E}[\sum_t \gamma^{t-1} F_t]$$

The last step follows because:

$$\mathbb{E}[b_t(s_t) \nabla_{\theta} \log(\pi_{\theta}(s_t, a_t)) | s_t] = b_t(s_t) \sum_a \nabla_{\theta} \pi_{\theta}(s_t, a) = b_t(s_t) \nabla_{\theta} (1) = 0$$

That is, following is also an unbiased gradient estimate

$$\hat{\mathbf{g}} = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \gamma^{t-1} (\hat{Q}_t^i - b_t(s_t^i)) \nabla_{\theta} \log(\pi(s_t^i, a_t^i)) \quad (4)$$

The difference  $\hat{Q}_t - b_t(s_t)$  is also referred to as **Advantage**. This terminology appears in more recent algorithms like ‘Asynchronous Advantage Actor Critic Algorithm (A3C)’ [Mnih et al., 2016] and ‘Generalized Advantage Estimation (GAE)’ [Schulman et al., 2015].

**Vanilla Policy Gradient Algorithm.** Initialize policy parameter  $\theta$ , and baseline function  $b_t(s), \forall s$ . In each iteration,

1. Execute current policy  $\pi^{\theta}$  to obtain several sample trajectories  $\tau^i, i = 1, \dots, m$ .
2. For any given sample trajectory  $i$ , use observed rewards  $r_1, r_2, \dots$ , to compute  $\hat{Q}_t^i := \sum_{t'=t}^{T-1} \gamma^{t'-t} r_{t'}$ .
3. Use  $\hat{Q}_t^i$  and baseline function  $b_t(s)$  to compute the gradient estimator  $\hat{\mathbf{g}}$  as in (4).
4. Update  $\theta \leftarrow \theta + \alpha \hat{\mathbf{g}}$ .
5. Re-optimize baseline.

## 4 Examples

A key contribution of policy gradient theorem is to reduce the computation of gradient of gain with respect to  $\theta$  to computation of gradient of policy function (or log of policy function  $\nabla_{\theta} \log(\pi(s, a; \theta))$ ) with respect to  $\theta$ . So, it is desirable that the gradient of policy function is efficiently computable. Here are some examples of policy functions where this is efficiently implementable.

### 4.1 Softmax policies

Consider policy set parameterized by  $\theta \in \mathbb{R}^d$  such that given  $s \in S$ , probability of picking action  $a \in A$  is given by:

$$\pi_{\theta}(s, a) = \frac{e^{\theta^{\top} \phi_{sa}}}{\sum_{a' \in A} e^{\theta^{\top} \phi_{sa'}}}$$

where each  $\phi_{sa}$  is an  $d$ -dimensional feature vector characterizing state-action pair  $s, a$ . This is a popular form of policy space called softmax policies. Here,

$$\nabla_{\theta} \log(\pi_{\theta}(s, a)) = \phi_{sa} - \left( \sum_{a' \in A} \phi_{sa'} \pi_{\theta}(s, a') \right) = \phi_{sa} - \mathbb{E}_{a' \sim \pi(s)}[\phi_{sa'}]$$

## 4.2 Gaussian policy for continuous action spaces

In continuous action spaces, it is natural to use Gaussian policies. Given state  $s$ , the probability of action  $a$  is given as:

$$\pi_\theta(s, a) = \mathcal{N}(\phi(s)^T \theta, \sigma^2)$$

for some constant  $\sigma$ . Here  $\phi(s)$  is a feature representation of  $s$ . Then,

$$\nabla_\theta \log(\pi_\theta(s, a)) = \nabla_\theta \frac{-(a - \theta^\top \phi(s))^2}{2\sigma^2} = \frac{(\theta^\top \phi(s) - a)}{\sigma^2} \phi(s)$$

## 4.3 Deep neural networks: backpropagation for gradient computation

In deep reinforcement learning, the policy function is computed by a multi-layer neural network. The independent-layer structure of deep neural network allows the gradient computations efficiently through **backpropagation**. (In the reinforcement learning context) Backpropagation refers to simply a computational implementation of the chain rule : an algorithm that determines the relevant partial derivatives by means of the backward pass.

Suppose  $\pi_\theta$  is given by a deep neural network whose input is representation of state  $s$ , and output is a representation of action selection probabilities  $\pi(s, a)$ . (We can add another layer giving logarithmic of the action selection probabilities.) The weights of the neural network form the parameters  $\theta$  of the policy. Specifically let the quantity of interest  $G_\theta(s, \cdot) \in \mathbb{R}^A$  is given by a neural network with  $k-1$  hidden layers and weights  $\theta = (W_1, \dots, W_K)$  and input  $\mathbf{x}$  being a representation of state  $s$ , so that

$$G_\theta(s, \cdot) = \log(\pi(s, \cdot)) = \log(\sigma(W^L(\sigma(\dots W^3 \sigma(W^2 \cdot \sigma(W^1 \mathbf{x}))))))$$

Or, (let  $L(\cdot)$  denote the log function)

$$\begin{aligned} G(s, a) &= \log(\pi(s, a)) = L(h_a^K) \\ h^k &= \sigma(z^k), z^k = W^k h^{k-1} \\ h^{k-1} &= \sigma(z^{k-1}), z^{k-1} = W^{k-1} h^{k-2}, \\ &\dots, \\ h^1 &= \sigma(z^1), z^1 = W^1 \mathbf{x} \end{aligned}$$

Now, gradient with respect of  $G(s, a)$  with respect to parameter  $W_{ab}^\ell$  is:

$$\frac{\partial G(s, a)}{\partial W_{ab}^\ell} = \frac{\partial \log(h_a^K)}{\partial h_a^K} \frac{\partial h_a^K}{\partial W_{ab}^\ell}$$

Then, basic observation is as follows. For a neuron  $r$  in layer  $k$ ,  $\ell < k$ :

$$\begin{aligned} \frac{\partial h_r^k}{\partial W_{ab}^\ell} &= \sigma'(z_r^k) \frac{\partial z_r^k}{\partial W_{ab}^\ell} \\ &= \sigma'(z_r^k) \frac{\partial W_r^k h^{k-1}}{\partial W_{ab}^\ell} \\ &= \sigma'(z_r^k) \sum_{i \in \text{parents}(r)} W_{ri}^k \frac{\partial h_i^{k-1}}{\partial W_{ab}^\ell} \end{aligned}$$

Thus, the gradients can be back propagated over the network, until we reach layer  $\ell$ , so that

$$\frac{\partial h_i^\ell}{\partial W_{ab}^\ell} = \begin{cases} \sigma'(z_a^\ell) h_b^{\ell-1}, & \text{if } i = a \\ 0 & \text{otherwise} \end{cases}$$



This follows from layer structure of the neural network so that:

$$\begin{aligned}
\frac{\partial z_a^\ell}{\partial W_{ab}^\ell} &= \frac{\partial \sum_i W_{ai}^\ell h_i^{\ell-1}}{\partial W_{ab}^{\ell-1}} \\
&= h_b^{\ell-1} + \sum_i W_{ai}^\ell \frac{\partial h_i^{\ell-1}}{\partial W_{ab}^\ell} \\
&= h_b^{\ell-1}
\end{aligned}$$

as layer  $h^{\ell-1}$  is independent of layer  $\ell$  parameters.

## References

- Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. *CoRR*, abs/1602.01783, 2016. URL <http://arxiv.org/abs/1602.01783>.
- John Schulman, Philipp Moritz, Sergey Levine, Michael I. Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *CoRR*, abs/1506.02438, 2015. URL <http://arxiv.org/abs/1506.02438>.
- Richard S. Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Proceedings of the 12th International Conference on Neural Information Processing Systems*, pages 1057–1063, 1999.
- R. J. Williams. Toward a theory of reinforcement-learning connectionist systems. Technical Report NU-CCS-88-3, College of Comp. Sci., Northeastern University, Boston, MA, 1988.
- Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3):229–256, May 1992. ISSN 1573-0565. doi: 10.1007/BF00992696. URL <https://doi.org/10.1007/BF00992696>.