

---

# CS 6101 Week #4 Notes: Actor-Critic Introduction, Value Functions and Q-Learning

---

Note taking: Alexandre Gravier, Joel Lee

L<sup>A</sup>T<sub>E</sub>X transcription: Alexandre Gravier <al.gravier@gmail.com>

## Contents

### 1 Recap about policy gradients

We define  $J(\theta) \doteq E_{\tau \sim p_\theta(\tau)} [\sum_t r(\mathbf{s}_t, \mathbf{a}_t)]$  so that the objective of RL can be defined as an optimization exercise consisting in finding an assignment of policy parameters  $\theta^* = \arg \max_\theta E_{\tau \sim p_\theta(\tau)} J(\theta)$ .

$J(\theta)$  is not usually optimizable as such due to f.e. dimensionality issues, so we use a sample-based unbiased estimate:  $J(\theta) \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$ . Taking the gradient of  $J(\theta)$  along  $\theta$  allows maximizing the expected reward as per the policy.

#### 1.1 Policy differentiation with a “convenient identity”

Let  $r(\tau) \doteq \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t)$  the total reward of a trajectory  $\tau$ .

$$\nabla_\theta J(\theta) = \nabla_\theta E_{\tau \sim p_\theta(\tau)} [r(\tau)] \quad (1a)$$

$$= \nabla_\theta \int \pi_\theta(\tau) r(\tau) d\tau \quad \text{by definition of expectation} \quad (1b)$$

$$= \int \nabla_\theta \pi_\theta(\tau) r(\tau) d\tau \quad \text{by linearity} \quad (1c)$$

At this point, the expression  $\int \nabla_\theta \pi_\theta(\tau) r(\tau) d\tau$  seems rather intractable. This is where the following “convenient identity” can be used to derive a tractable expression of  $\nabla_\theta J(\theta)$ .

#### A convenient identity

$$\pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau) = \pi_\theta(\tau) \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} = \nabla_\theta \pi_\theta(\tau) \quad (2)$$

Furthermore, we can expand the definition of  $\pi_\theta(\tau)$  and take its logarithm:

$$\pi_\theta(\tau) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \quad (3a)$$

$$\Leftrightarrow \log \pi_\theta(\tau) = \log p(\mathbf{s}_1) + \sum_{t=1}^T \log \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) + \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \quad (3b)$$

Using (??) and (??) in (??), the gradient of the objective becomes:

$$\nabla_{\theta} J(\theta) = \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau) d\tau \quad \text{using (??)} \quad (4a)$$

$$= E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)] \quad \text{by definition of expectation} \quad (4b)$$

$$= E_{\tau \sim p_{\theta}(\tau)} \left[ \nabla_{\theta} \left[ \log p(\mathbf{s}_1) + \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) + \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \right] r(\tau) \right] \quad (4c)$$

We note that in the expression  $\nabla_{\theta} \left[ \log p(\mathbf{s}_1) + \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) + \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \right]$  of the gradient w.r.t.  $\theta$  in (??), the terms  $\log p(\mathbf{s}_1)$  and  $\log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$  are independent of  $\theta$ , so we are left with:

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[ \nabla_{\theta} \left[ \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right] r(\tau) \right] \quad (5a)$$

$$= E_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) r(\tau) \right] \quad \text{by linearity} \quad (5b)$$

$$= E_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right] \quad \text{by definition of } r(\tau) \quad (5c)$$

In Equation (??), the gradient of  $J$  is now a computable function of  $\pi_{\theta}$  only.

We earlier mentioned the sample estimate of  $J(\theta) \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$ ; similarly  $\nabla_{\theta} J(\theta)$  is approximated with samples, leading us to the algorithm:

$$\text{REINFORCE algorithm: } \begin{cases} \text{sample } \{\tau^i\} \text{ from } \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \\ \nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \\ \theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta) \end{cases} \quad (6)$$

In (??), there is no use of the Markov property, so the algorithm can be used as such on POMDPs.

## 1.2 The bad news

We introduce some simplifying notation:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \quad (7a)$$

$$\approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau) \quad (7b)$$

The big problem with vanilla REINFORCE lies in variance: if you were to repeatedly collect a small finite number  $N$  of samples, estimating each time the gradient based on these samples, you would observe that these estimates of the gradient vary a lot.

That means that the gradient descent step will likely take us away from the goal, and convergence will be very slow, or may not happen.

Additionally, given two sets of sample trajectories differing only by a constant factor, but differently centered around a total reward of 0, the basic policy gradient algorithm in (??) may change the parameters very differently.

There are two tricks that help with the large variance of the samples, both of which should always be used because they have no drawbacks.

### 1.3 Variance reduction with causality: the “rewards to go” trick

In (??), the gradient approximation can be improved by making use of causality: the policy at time  $t'$  cannot affect the reward at time  $t$  when  $t < t'$ .

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \quad (8a)$$

$$\approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \left( \underbrace{\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}_{\text{gradient at } t} \underbrace{\sum_{t'=1}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})}_{\text{this total reward should be computed from time } t} \right) \right) \quad \text{by distributivity} \quad (8b)$$

$$\approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \left( \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \underbrace{\sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})}_{\text{reward to go}} \right) \right) \quad \text{by causality} \quad (8c)$$

Intuitively, the “rewards to go” trick comes from the observation that at time  $t$ , all past rewards cannot be affected by policy decisions.

We define the “reward to go” function  $\hat{Q} : \llbracket 1, N \rrbracket \times \llbracket 1, T \rrbracket \mapsto \mathbb{R}$  as:

$$\hat{Q}_{i,t} \doteq \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \quad (9)$$

Therefore, the gradient of the objective function approximation with rewards to go is:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \hat{Q}_{i,t} \quad (10)$$

We should note that most often, in policy gradient, we search the optimal policy within a class of policy functions that is time-invariant. However, to get the optimal policy in a finite horizon problem ( $T < \infty$ ), we may need a time-dependent policy (e.g. in the last time step, there is no more need to care about energy preservation).

### 1.4 Variance reduction with baselines

Ideally, policy gradients should only depend on the relative values of the trajectory samples under the policy. But in REINFORCE, policy gradients dependent on the value of the total rewards of the policy samples: given the same set of samples shifted by a constant reward, the gradient direction and norm may change drastically. This is very undesirable.

The intuition behind the baseline trick is to subtract a constant “average reward” baseline from the trajectory rewards, so as to make trajectories that are less rewarding than average less likely, and those more rewarding more likely.

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau) (r(\tau) - b) \quad \text{where } b = \frac{1}{N} \sum_{i=1}^N r(\tau) \quad (11)$$

Regardless of the baseline  $b$  used in computing the gradient, and the policy gradient defined in (??) remains an unbiased estimator of the gradient of the objective.

*Proof that the gradient estimator in Eq. (??) is unbiased for any  $b$ .* To prove that the expression in (??) is an unbiased estimator of  $\nabla_{\theta} J(\theta)$ , we need to prove that  $E[\nabla_{\theta} \log \pi_{\theta}(\tau)(r(\tau) - b)] \stackrel{?}{=} E[\nabla_{\theta} \log \pi_{\theta}(\tau)r(\tau)]$ , which — by linearity of expectation — is equivalent to:

$$E[\nabla_{\theta} \log \pi_{\theta}(\tau)b] \stackrel{?}{=} 0 \quad (12)$$

We expand the left-hand side:

$$E[\nabla_{\theta} \log \pi_{\theta}(\tau)b] = \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) b d\tau \quad \text{by definition of expectation} \quad (13a)$$

$$= \int \nabla_{\theta} \pi_{\theta}(\tau) b d\tau \quad \text{by substitution with (??)} \quad (13b)$$

$$= b \nabla_{\theta} \int \pi_{\theta}(\tau) d\tau \quad \text{by linearity} \quad (13c)$$

$$= b \nabla_{\theta} \int 1 \quad \text{as } \pi_{\theta} \text{ is a probability distribution} \quad (13d)$$

$$= 0 \quad (13e)$$

Hence, any baseline  $b$  may be added to Eq. (??), and it will remain an unbiased estimator of the gradient of  $J(\theta)$ .  $\square$

Therefore, after subtracting an average reward baseline from the trajectories' reward, Eq. (??) remains an unbiased estimate of the gradient of the objective, but with a lower variance.

Choosing  $b = \frac{1}{N} \sum_{i=1}^N$  works well, but it is not the baseline that results in the lowest variance estimator. To derive the best baseline, we can express that at its minimum, the derivative of the variance in function of the baseline is null.

The variance of the gradient of the objective is, by definition:

$$\text{Var} = E[(\nabla_{\theta} J(\theta))^2] - E[\nabla_{\theta} J(\theta)]^2 \quad (14a)$$

$$= E_{\tau \sim \pi_{\theta}(\tau)} [(\nabla_{\theta} \log \pi_{\theta}(\tau)(r(\tau) - b))^2] - \underbrace{E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau)(r(\tau) - b)]^2}_{\text{as shown in Eqs. (??), this term is independent of } b} \quad (14b)$$

As we observe that only the first term depends on  $b$ ,

$$\frac{d\text{Var}}{db} = \frac{d}{db} E_{\tau \sim \pi_{\theta}(\tau)} [(\nabla_{\theta} \log \pi_{\theta}(\tau)(r(\tau) - b))^2] \quad (15a)$$

$$= \frac{d}{db} E [(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 (r(\tau) - b)^2] \quad (15b)$$

$$= \frac{d}{db} E [g(\tau)^2 (r(\tau) - b)^2] \quad (\text{we introduce } g(\tau) \doteq \nabla_{\theta} \log \pi_{\theta}(\tau)) \quad (15c)$$

$$= \frac{d}{db} (E [g(\tau)^2 r(\tau)^2] - 2bE [g(\tau)^2 r(\tau)] + b^2 E [g(\tau)^2]) \quad (15d)$$

$$= \frac{d}{db} (-2bE [g(\tau)^2 r(\tau)] + b^2 E [g(\tau)^2]) \quad (15e)$$

$$= -2E [g(\tau)^2 r(\tau)] + 2bE [g(\tau)^2] \quad (15f)$$

We want to find  $b$  s.t.  $\frac{d\text{Var}}{db} = 0$ :

$$\frac{d\text{Var}}{db} = 0 \tag{16a}$$

$$\implies -2E[g(\tau)^2 r(\tau)] + 2bE[g(\tau)^2] = 0 \tag{16b}$$

$$\implies b = \frac{E[g(\tau)^2 r(\tau)]}{E[g(\tau)^2]} \tag{16c}$$

We have derived the baseline that minimizes the variance of the objective function:

**The optimal baseline**

$$\hat{b} = \frac{E\left[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2 r(\tau)\right]}{E\left[(\nabla_{\theta} \log \pi_{\theta}(\tau))^2\right]} \tag{17}$$

This the expected reward weighted by the square of the gradient magnitude, **element-wise**: the baseline  $\hat{b}$  is also defined per element of the vector of parameters, and the  $g(\tau)^2$  is an element-wise square of the gradient vector. Each scalar in the gradient baseline vector is weighted by its contribution to the gradient.

It is possible to implement  $\hat{b}$ , and it is guaranteed to be the optimal baseline w.r.t variance reduction. However, researchers often choose the simpler solution of the expected reward baseline, as it is easier to implement.

## 2 On-policy and off-policy algorithms