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

Note taking: Alexandre Gravier, Joel Lee
L^AT_EX transcription: Alexandre Gravier <al.gravier@gmail.com>

Contents

1	Recap about policy gradients	1
1.1	Policy differentiation with a “convenient identity”	1
1.2	The bad news	2
1.3	Variance reduction with “rewards to go”	2

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 θ 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 τ .

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

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

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

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

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

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 (2) and (3b) in (1c), the gradient of the objective becomes:

$$\nabla_\theta J(\theta) = \int \pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau) r(\tau) d\tau \quad \text{using (2)} \quad (4a)$$

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

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

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. θ in (4c), the terms $\log p(\mathbf{s}_1)$ and $\log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$ are independent of θ , 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 (5c), the gradient of J is now a computable function of π_θ 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 (6), there is no use of the Markov property, so the algorithm can be used as such on POMDPs.

1.2 The bad news

1.3 Variance reduction with “rewards to go”

The “**rewards to go**” trick computes the expected future rewards based on from the current time t to reduce variance in J . It comes from the observation that at time t , all past rewards cannot be affected by policy decisions.

Rewards to go -

References