

# CHAPTER 4

---

## Risk-Sensitive Actor-Critic Algorithm

---

Since we found a suitable risk-adjusted value function in Chapter 3 that satisfies a Bellman equation, we now aim to connect sections 2.2 and 3.4. The goal is to incorporate the Bellman equation (3.7) directly into a reinforcement learning algorithm and exploit the fact that risk-sensitive strategies can then be learned. Since we want to experiment in environments with continuous action and state spaces later, the authors in [14] propose a slightly modified version of the (soft) actor-critic algorithm.

---

### Algorithm 1 Multiple risk aversion actor-critic

---

```

1: Initialize actor network  $\pi^\theta$  and critic network  $V^\phi$  with weights  $\theta, \phi$ 
2: for Episode 1,..., M do
3:   Sample a batch of N initial states  $s_0$  and risk aversion levels  $\lambda$ 
4:   for  $t = 0, \dots, T$  do
5:     Select actions from current policy  $a_t = \pi(s_t)$ 
6:     Execute actions, receive rewards  $r_t$  and transition into new states  $s_{t+1}$ 
7:     Compute the target
8:      $y_t = r_t + \bar{V}(s_{t+1})$ 
9:     Update the critic by gradient descent with respect to the loss
10:     $\mathcal{L}(\phi) = \frac{1}{N} \sum_{i=1}^N \left[ \frac{1}{\lambda} \exp(-\lambda(y_t - V(s_t))) - V(s_t) \right]$ 
11:    Update the actor by gradient descent with respect to the loss
12:     $\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda} \exp(-\lambda(r_t + V(s_{t+1})))$ 
13:    Update target critic network  $\bar{\phi} = \tau\phi + (1 - \tau)\bar{\phi}$ 
14:   end for
15: end for

```

---

The goal of this chapter is to investigate whether this proposal is truly compatible with our risk-adjusted Bellman equation. Therefore, we focus on the loss functions

$$\mathcal{L}_{\text{actor}}(\theta) = \mathbb{E}_{s_t, a_t, s_{t+1}} \left[ \frac{1}{\lambda} \exp(-\lambda(r_t + V_\phi(s_{t+1}))) \right], \quad (4.1)$$

$$\mathcal{L}_{\text{critic}}(\phi) = \mathbb{E}_{s_t, a_t, s_{t+1}} \left[ \frac{1}{\lambda} \exp(-\lambda(y_t - V_\phi(s_t))) - V_\phi(s_t) \right], \quad (4.2)$$

where

$$y_t := r_t + V_{\text{target}}(s_{t+1})$$

is the bootstrapped target (with a slowly updated target network  $V_{\text{target}}$ ), and to show that these losses correspond to the risk-averse Bellman equation based on exponential utility.

We recall the risk-averse value function  $V_\lambda^\pi$ , which satisfies the Bellman equation (3.7) (and which our policy should maximize), and define the exponentiated value function  $W_\lambda^\pi(s)$ .<sup>1</sup>

$$\begin{aligned} V_\lambda^\pi(s_t) &= -\frac{1}{\lambda} \log \mathbb{E} \left[ \exp(-\lambda(r_t + V_\lambda^\pi(s_{t+1}))) \middle| s_t \right] \\ W_\lambda^\pi(s) &:= \exp(-\lambda V_\lambda^\pi(s)) = \mathbb{E} \left[ \exp(-\lambda(r_t + V_\lambda^\pi(s_{t+1}))) \middle| s_t \right]. \end{aligned} \quad (4.3)$$

## 4.1 Actor Loss

With a distribution  $\mu$  over initial states  $s_0 \in S$  in our environment, the natural risk-averse objective is to simply maximize

$$J(\pi) := \mathbb{E}_{s_0 \sim \mu} [V_\lambda^\pi(s_0)]. \quad (4.4)$$

Equivalently, the optimal value function at  $t = 0$  is  $V_\lambda^*(s_0) = \sup_\pi V_\lambda^\pi(s_0)$  and we want  $\pi^* \in \arg \max_\pi J(\pi)$ . For fixed  $\lambda > 0$ , the map  $f(x) := \frac{1}{\lambda} e^{-\lambda x}$  is strictly decreasing in  $x$ . Hence, maximizing (4.4) is, in terms of policies, equivalent to minimizing

$$\tilde{J}(\pi) := \mathbb{E}_{s_0 \sim \mu} \left[ \frac{1}{\lambda} \exp(-\lambda V_\lambda^\pi(s_0)) \right] = \mathbb{E}_{s_0 \sim \mu} \left[ \frac{1}{\lambda} W_\lambda^\pi(s_0) \right]. \quad (4.5)$$

such that we will continue with our equivalent control problem (4.5) which is related to the exponentiated Bellman equation (4.3). For any  $t \in \mathbb{N}$ , taking expectation over  $s_t$  under the state distribution induced by the current policy  $\pi$  (the on-policy distribution or  $P(\cdot | s_t, a_t)$  from chapter 2.1), gives

$$\mathbb{E}_{s_t} \left[ \frac{1}{\lambda} W_\lambda^\pi(s_t) \right] = \mathbb{E}_{s_t} \left[ \frac{1}{\lambda} \mathbb{E}_{s_{t+1} | a_t} \left[ \exp(-\lambda(r_t + V_\lambda^\pi(s_{t+1}))) \middle| s_t \right] \right] = \frac{1}{\lambda} \mathbb{E}_{s_t, a_t, s_{t+1}} \left[ \exp(-\lambda(r_t + V_\lambda^\pi(s_{t+1}))) \right]. \quad (4.6)$$

Thus the exponentiated value function is equal (in expectation) to the exponentiated one-step reward plus next value, and this Bellman identity also shows that  $J(\pi)$  is maximized if it pointwise minimizes  $\mathbb{E}_{s_t} [\frac{1}{\lambda} W_\lambda^\pi(s_t)]$ . We now define parameterizations of our networks:

- $\pi_\theta$  is a deterministic policy network with parameters  $\theta$ .
- $V_\phi$  is a value network (critic) with parameters  $\phi$ , approximating  $V_\lambda^{\pi_\theta}$ .

Now our on-policy transformations  $(s_t, a_t, s_{t+1})$  are sampled using the current policy of our actor  $\pi_\theta$ . If  $V_\phi$  were exactly equal to  $V_\lambda^{\pi_\theta}$ , then, from (4.6),

$$\mathbb{E}_{s_t} \left[ \frac{1}{\lambda} W_\lambda^\pi(s_t) \right] = \frac{1}{\lambda} \mathbb{E}_{s_t, a_t, s_{t+1}} \left[ \exp(-\lambda(r_t + V_\lambda^{\pi_\theta}(s_{t+1}))) \right] = \frac{1}{\lambda} \mathbb{E}_{s_t, a_t, s_{t+1}} \left[ \exp(-\lambda(r_t + V_\phi(s_{t+1}))) \right]. \quad (4.7)$$

This motivates defining the actor loss

$$\mathcal{L}_{\text{actor}}(\theta) := \mathbb{E}_{s_t, a_t^{\pi_\theta}, s_{t+1}} \left[ \frac{1}{\lambda} \exp(-\lambda(r_t + V_\phi(s_{t+1}))) \right], \quad (4.8)$$

which can be approximated using one-step transitions and the critic  $V_\phi$ . From our above derivation follows that a policy  $\pi$  that minimizes  $\mathcal{L}_{\text{actor}}$  is optimal for the risk-averse control problem (4.4) under the assumption that  $V_\phi$  is a decent approximation of  $V_\lambda^{\pi_\theta}$ , which we will investigate in the next section.

<sup>1</sup>Since we are in the one-step Bellman environment, if we do not specify  $\mathbb{E}$  we mean the one-step  $\mathbb{E}_{s_{t+1} | a_t}$ .

<sup>2</sup>In [14] they use single transitions, which is why their  $\mathbb{E}_{s_t, a_t, s_{t+1}}$  reduces to  $\mathbb{E}_{s_t}$ .

## 4.2 Critic Loss Fixed Point

We want a loss function whose minimizer (for each fixed  $s_t$ ) is exactly the Bellman value  $V_\lambda^\pi(s_t)$ . Fix a state  $s_t \in S$  and consider  $y$  as a random variable with distribution given by the next-step dynamics under policy  $\pi$ :

$$y := r(s, a) + V_\lambda^\pi(s_{t+1}), \quad s_{t+1} \sim P(\cdot | s_t, a_t).$$

We seek a scalar function  $f(\nu)$  such that  $\nu^* := V_\lambda^\pi(s_t)$  is the unique minimizer of  $f(\nu)$ . The paper [14] proposes for each state  $s_t$ ,

$$f_s(\nu) := \mathbb{E}_{s_{t+1}|a_t} \left[ \frac{1}{\lambda} \exp(-\lambda(y - \nu)) - \nu \mid s_t \right]. \quad (4.9)$$

For fixed  $y$ ,

$$\frac{\partial}{\partial \nu} \left( \frac{1}{\lambda} \exp(-\lambda(y - \nu)) - \nu \right) = \exp(-\lambda(y - \nu)) - 1.$$

Thus<sup>3</sup>  $f'_s(\nu) = \mathbb{E}_{s_{t+1}|a_t} [\exp(-\lambda(y - \nu)) - 1 \mid s_t]$ . A stationary point  $\nu^*$  must satisfy  $f'_s(\nu^*) = 0$ , i.e.

$$\mathbb{E}_{s_{t+1}|a_t} [\exp(-\lambda(y - \nu^*)) \mid s_t] = 1.$$

We make the following observation.

**Proposition 1.** For any real  $\nu$ , the following are equivalent:

$$\mathbb{E}[\exp(-\lambda(y - \nu)) \mid s_t] = 1, \quad (4.10)$$

$$\exp(-\lambda\nu) = \mathbb{E}[\exp(-\lambda y) \mid s_t]. \quad (4.11)$$

In particular,  $\nu^* = V_\lambda^\pi(s_t)$  satisfies both.

*Proof.* The equivalence is immediate.

$$1 = \mathbb{E}[\exp(-\lambda(y - \nu)) \mid s_t] = \exp(\lambda\nu) \mathbb{E}[\exp(-\lambda y) \mid s_t].$$

Moreover, note that  $\nu^* = V_\lambda^\pi(s_t)$  satisfies (4.11) by construction (4.3). Hence it also satisfies (4.10).  $\square$

We have proven that  $V_\lambda^\pi(s_t)$  minimizes  $f_s(\nu)$ . Arguing that it is strictly convex in  $\nu$ , it is also unique.

### Parametrized critic loss

We now generalize from a scalar  $\nu$  to a value function  $V_\phi : S \rightarrow \mathbb{R}$ . Since our actor evaluates  $V_\lambda^{\pi_\theta}(s_t)$  state-wise, or more precisely, its expectation. Then, the state-wise function  $f_s$  naturally leads to a loss

$$\begin{aligned} \mathcal{L}_{\text{critic}}(\phi) &:= \mathbb{E}_{s_t} [f_{s_t}(V_\phi(s_t))] = \mathbb{E}_{s_t} \left[ \mathbb{E}_{a_t^{\pi_\theta}, s_{t+1}} \left[ \frac{1}{\lambda} \exp(-\lambda(y_t - V_\phi(s_t))) - V_\phi(s_t) \mid s_t \right] \right] \\ &= \mathbb{E}_{s_t, a_t^{\pi_\theta}, s_{t+1}} \left[ \frac{1}{\lambda} \exp(-\lambda(y_t - V_\phi(s_t))) - V_\phi(s_t) \right]. \end{aligned} \quad (4.12)$$

Thus, for each state  $s_t \in S$ , the critic value  $V_\phi(s_t)$  being the Bellman value  $V_\lambda^{\pi_\theta}(s_t)$  uniquely minimizes the critic loss used in the algorithm. Note that in practice  $y_t$  is defined using a separate target network  $V_{\text{target}}$ , i.e.  $y_t := r_t + V_{\text{target}}(s_{t+1})$ , and gradients are not backpropagated through  $y_t$ .

---

<sup>3</sup>under standard integrability assumptions.