

# From Reflexes to Actor-Critic

May 16, 2015

We have developed a robot system that avoids collisions using a fixed stochastic policy with a given representation. The representation  $x$  has four components indexed from 0, (0) a bias unit, (1) a learned policy-contingent prediction (GVF) of the bumper, (2) the previous action was forward, (3) the previous action was to turn counterclockwise.

The policy is originally crafted by formalizing the real-world characteristics desired from the behavior. We have two action forward (F) and turn (T). The prediction  $p0$  anticipates when the robot bumper will be pressed, presumably from a collision with a wall. The ideal value of  $p0$  is 1 when the bumper will be pressed on the next timestep if the robot drives forward, and it is 0 when the robot is infinitely far away from a wall. The robot interacts with the world  $H = 100/3$  times per second. The prediction  $p0$  has a gamma of 0 when the bumper is pressed, and a gamma of 0.95 otherwise. Thus in the absence of any bumper presses, the prediction extends for  $20(= (1 - \gamma)^{-1})$  timesteps in expectation (i.e.  $2/3$  of a second). We want the robot to have an equal probability selecting a drive action or a turn action when the ideal prediction has a value of  $eq = 0.5$  (what is the timescale?  $0.95^{timescale} = 0.5$ ).

Now as we don't want to be limited by the decision rate of the robot, we instead specify the desired number of switches per second if the prediction was at this equilibrium point, and we set this value to be 5. Let  $TT = \text{decisions-per-second} / \text{expected-number-of-switches-per-second}$ . After some algebraic manipulation, we find that there is a constant  $K_s = \log_e((TT - 1)(\text{number} - \text{of} - \text{actions} - 1))$  which would govern switching at equilibrium in the absence of more information from the prediction.

There is one more constant used to scale the relative influence of the prediction on the action selection, we chose the number 4 (so  $p0 = .75$  means turning is 4 times more likely than going forward).

We can write the preferences and the policy in an intuitive form.

$$pref[F] = -4K_s(p0 - eq) + K_s 1(\text{LastAction} == F)$$

$$pref[T] = K_s 1(\text{LastAction} == T)$$

$$\pi(a = F|x) = \exp(pref[F]) / (\sum_{i=F,T} \exp(pref[i]))$$

$$\pi(a = T|x) = \exp(pref[T]) / (\sum_{i=F,T} \exp(pref[i]))$$

We rewrite the stochastic policy in a canonical quasi-linear form, with action selection probabilities given by exponentiated preferences for each action, where each preference is a linear function of the feature vector  $x$ , and the policy parameters  $u$  are initialized to the constants from above.

$$pref[F] = u_0 x_0 + u_1 x_1 + u_2 x_2$$

$$pref[T] = u_3 x_3$$

$$\pi(a = F|x) = \exp(pref[F]) / (\sum_{i=F,T} \exp(pref[i]))$$

$$\pi(a = T|x) = \exp(pref[T]) / (\sum_{i=F,T} \exp(pref[i]))$$

Now an average-reward actor-critic algorithm can be used to tune the parameters  $u$  so as to maximize a given reward signal over a long term, but it first requires us to calculate gradients, namely  ${}_u\pi(a = k|x)/\pi(a = k|x)$ . Given the above equations, the gradients can be found through the simple, tedious calculations shown below.

$$\frac{\nabla_u \pi}{\pi} = ((\partial_{u_0} \pi)/\pi, (\partial_{u_1} \pi)/\pi, (\partial_{u_2} \pi)/\pi)$$

To compute these, it is useful to precompute  $\partial_{u_j} \exp(\text{pref}[i])$  for all  $i$  and  $j$ .

$$\partial_{u_0} \exp(\text{pref}[F]) = \partial_{u_0} \exp(u_0 x_0 + u_1 x_1 + u_2 x_2) = \exp(u_0 x_0 + u_1 x_1 + u_2 x_2) \partial_{u_0} (u_0 x_0 + u_1 x_1 + u_2 x_2) = \exp(\text{pref}[F]) x_0$$

$$\partial_{u_1} \exp(\text{pref}[F]) = \exp(\text{pref}[F]) x_1$$

$$\partial_{u_2} \exp(\text{pref}[F]) = \exp(\text{pref}[F]) x_2$$

---


$$\partial_{u_0} \exp(\text{pref}[T]) = 0$$

$$\partial_{u_1} \exp(\text{pref}[T]) = 0$$

$$\partial_{u_2} \exp(\text{pref}[T]) = \exp(\text{pref}[T]) x_3$$

---


$$\text{Let } W = \sum_i \exp(\text{pref}(i)).$$

$$\partial_{u_0} W = \exp(\text{pref}(F)) x_0$$

$$\partial_{u_1} W = \exp(\text{pref}(F)) x_1$$

$$\partial_{u_2} W = \exp(\text{pref}[F]) x_2 + \exp(\text{pref}(T)) x_3$$

---

Now we can write the relevant gradients directly, recalling that  $\pi(a|x) = \exp(\text{pref}(a))/W$ .

$$\begin{aligned} \frac{\partial_{u_0} \pi(F|x)}{\pi(F|x)} &= \frac{W}{\exp(\text{pref}(F))} \partial_{u_0} \frac{\exp(\text{pref}(F))}{W} \\ &= \frac{W}{\exp(\text{pref}(F))} \frac{1}{W^2} ((\partial_{u_0} \exp(\text{pref}(F)))W - \exp(\text{pref}(F)) \partial_{u_0} W) \\ &= \frac{W}{\exp(\text{pref}(F))} \frac{1}{W^2} (\exp(\text{pref}(F)) x_0 W - \exp(\text{pref}(F)) \exp(\text{pref}(F)) x_0) \\ &= \frac{1}{1} \frac{1}{W} (x_0 W - \exp(\text{pref}(F)) x_0) \\ &= x_0 (1 - \frac{\exp(\text{pref}(F))}{W}) \\ &= x_0 (1 - \pi(F|x)) \end{aligned}$$

$$\begin{aligned} \frac{\partial_{u_1} \pi(F|x)}{\pi(F|x)} &= \frac{W}{\exp(\text{pref}(F))} \partial_{u_1} \frac{\exp(\text{pref}(F))}{W} \\ &= \frac{W}{\exp(\text{pref}(F))} \frac{1}{W^2} ((\partial_{u_1} \exp(\text{pref}(F)))W - \exp(\text{pref}(F)) \partial_{u_1} W) \\ &= \frac{W}{\exp(\text{pref}(F))} \frac{1}{W^2} (\exp(\text{pref}(F)) x_1 W - \exp(\text{pref}(F)) \exp(\text{pref}(F)) x_1) \\ &= x_1 (1 - \pi(F|x)) \end{aligned}$$

$$\frac{\partial_{u_2} \pi(F|x)}{\pi(F|x)} = x_2(1 - \pi(F|x))$$


---

$$\begin{aligned} \frac{\partial_{u_0} \pi(T|x)}{\pi(T|x)} &= \frac{W}{\exp(\text{pref}(T))} \partial_{u_0} \frac{\exp(\text{pref}(T))}{W} \\ &= \frac{W}{\exp(\text{pref}(T))} \frac{1}{W^2} ((\partial_{u_0} \exp(\text{pref}(T)))W - \exp(\text{pref}(T))\partial_{u_0} W) \\ &= \frac{1}{\exp(\text{pref}(T))} \frac{1}{W} (0 - \exp(\text{pref}(T)) \exp(\text{pref}(F))x_0) \\ &= -\frac{\exp(\text{pref}(F))}{W} x_0 \\ &= -\pi(F|x)x_0 \end{aligned}$$

$$\frac{\partial_{u_1} \pi(T|x)}{\pi(T|x)} = -\pi(F|x)x_1$$

$$\begin{aligned} \frac{\partial_{u_2} \pi(T|x)}{\pi(T|x)} &= \frac{W}{\exp(\text{pref}(T))} \partial_{u_2} \frac{\exp(\text{pref}(T))}{W} \\ &= \frac{W}{\exp(\text{pref}(T))} \frac{1}{W^2} ((\partial_{u_2} \exp(\text{pref}(T)))W - \exp(\text{pref}(T))\partial_{u_2} W) \\ &= \frac{1}{\exp(\text{pref}(T))} \frac{1}{W} ((\exp((T))x_3)W - \exp(\text{pref}(T))(\exp((F))x_2 + \exp(\text{pref}(T))x_3)) \\ &= \frac{1}{1} \frac{1}{W} (x_3W - (\exp((F))x_2 + \exp(\text{pref}(T))x_3)) \\ &= (x_3 - \pi(F)x_2 - \pi(T)x_3) \\ &= (1 - \pi(T))x_3 - \pi(F)x_2 \end{aligned}$$

---

These gradient terms are used in the actor-critic algorithm, evaluated for a given action (F or T) and the feature vector  $\mathbf{x}$ .

The average reward formulation is given in the Actor-Critic algorithms in practice paper, and it is listed below.

---

#### Average Reward Actor-Critic Algorithm

---

For each step with feature vector  $x(S)$

Choose  $A$  according to  $\pi_u(\cdot|S)$

Take action  $A$

Observe  $x(S')$ ,  $R$

$$\delta \leftarrow R - \bar{R} + v^\top x(S') - v^\top x(S)$$

$$\bar{R} \leftarrow \alpha_1 \delta$$

$$e_v \leftarrow \lambda e_v + x(S)$$

$$v \leftarrow v + \alpha_2 \delta e_v$$

$$e_u \leftarrow \lambda e_u + \frac{\nabla_u \pi(A|x(S))}{\pi(A|x(S))}$$

$$u \leftarrow u + \alpha_3 \delta e_u$$

---