

# Human in the loop RL

Sandip Patel

28 July 2020

## 1 Variance of the Return

Variance of any variable  $X$  can be given using the following expression:

$$V[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \quad (1)$$

We already know the learning equation for expectation of the returns in the form of value functions and hence we have learning equation for  $\mathbb{E}[x]$ . In this document I am presenting the work to calculate  $\mathbb{E}[X^2]$ (second moment) term and finally calculate the Variance of returns in the form of Learning equations (or update rule).

## 2 Expectation of second moment

Let the variable  $x$  be the return by taking an action  $a \in A$  from state  $s \in S$  and then following the policy. Hence,

$$\mathbb{E}[x] = Q(s, a) = J(s, a) \quad (2)$$

$$Q(s, a) = r(s, a) + \gamma * \sum P(s'|s, a)Q(s', a') \quad (3)$$

our well known action value function.

Let the expectation of second moment be represented by  $M(s, a)$ . So,

$$\mathbb{E}[x^2] = M(s, a) \quad (4)$$

$$\begin{aligned} M(s, a) &= \mathbb{E}[\text{returns}^2 | s, a] \\ &= \mathbb{E}[(\sum \gamma^k * r(s, a))^2 | s \in S, a \in A] \\ &= \mathbb{E}[(r(s_o, a_o) + \sum \gamma^k * r(s', a'))^2 | s_o, s' \in S, a_o, a' \in A] \\ &= r(s_o, a_o)^2 + 2\gamma r(s_o, a_o) \sum P(s'|s_o, a_o)Q(s', a') + \gamma^2 \sum P(s'|s_o, a_o)M(s', a') \end{aligned} \quad (5)$$

Using Equation (4) one can easily find expectation of the square of the *returns*.

Update rule for  $Q(s, a)$  from equation (3) is well known and can be given using the following equation.

$$Q_t(s, a) = Q_{t-1}(s, a) + \alpha_q [r(s, a) + \gamma * \text{argmax}_a (Q_{t-1}(s', a')) - Q_{t-1}(s, a)] \quad (6)$$

Similarly, we can write the update rule for equation (5) as following.

$$M_t(s, a) = M_{t-1}(s, a) + \alpha_m [r(s, a)^2 + 2\gamma r(s, a) * \text{argmax}_a (Q_{t-1}(s', a')) + \gamma^2 M_{t-1}(s', a') - M_{t-1}(s, a)] \quad (7)$$

## 3 Variance Update

Once we have the new estimate of  $Q(s, a)$  and  $M(s, a)$  at time step  $t$  using (6) and (7) we can use (1) to update the variance.

Since,

$$V_t(s, a) = M_t(s, a) - Q_t(s, a)^2 \quad (8)$$

We can update the Variance using the simple Temporal difference method. Which come out to be

$$V_t(s, a) = V_{t-1}(s, a) + \alpha_v [M_t(s, a) - Q_t^2(s, a) - V_{t-1}(s, a)] \quad (9)$$

The final eqns here have the max over 'a', which I think should also be present in eqns 3 and 5

I believe argmax should be replaced by max

M\_t-1 should also have an max around it.

This term has to be squared