

Convergence of actor-critic

Consider

$$\max_{\theta} V^{\theta}(x) = E \left[\sum_{k=0}^{\infty} \gamma^k r(x_k, u_k) \middle| X_0 = x \right]$$

s. t. $X_{k+1} \sim f(x_{k+1} | x_k, u_k)$

$$V_k = \rho^{\theta}(x_k)$$

Recall the HJB:

$$V^*(x) = \max_u \{ r(x, u) + \gamma E[V^*(x^u)] \}$$

Recall the Q-function:

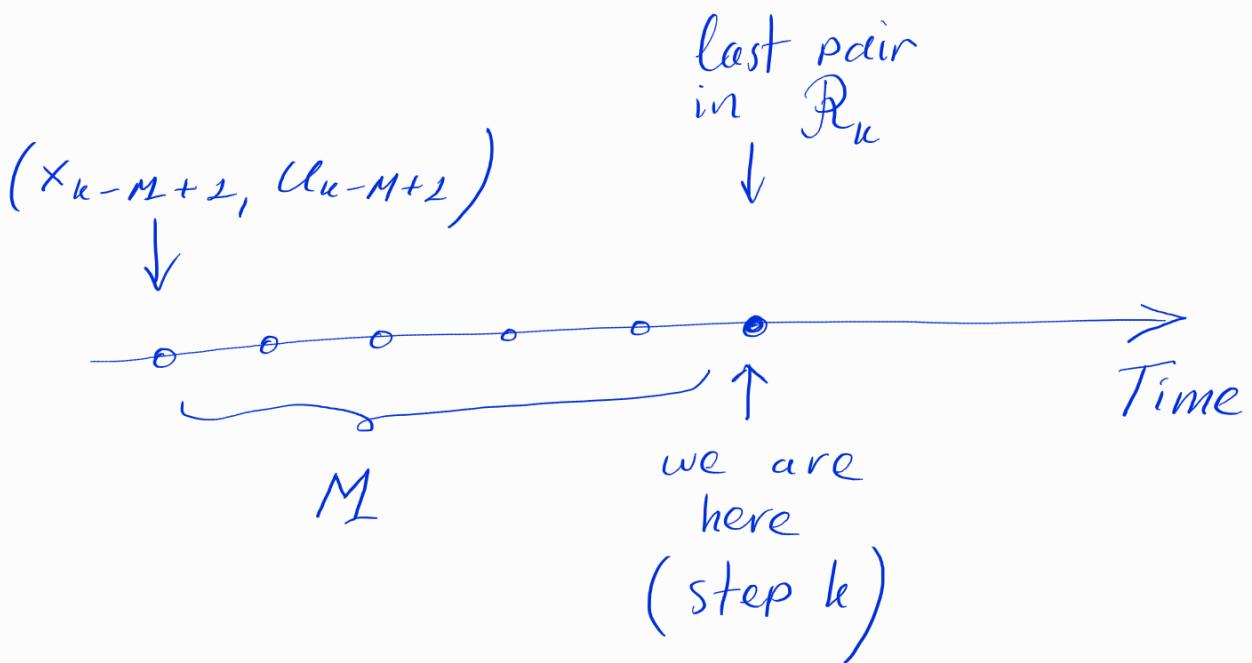
$$Q^*(x, u) = r(x, u) + \gamma E[V^*(x^u)]$$

Assuming a NN model of the Q-function $\hat{Q}^{\omega}(x, u)$, state the TD:

$$e_{TD}^{\omega}(x, u, x_t, u_t) = \hat{Q}^{\omega}(x, u) - r(x, u) - \hat{Q}^{\omega}(x_t, u_t)$$

Let's say we have a data buffer (experience replay) :

$$\mathcal{R}_k := \left\{ (x_{k-j+1}, u_{k-j+1}) \right\}_{j=1}^M$$



First, let's pick a model :

$$\hat{Q}^\omega(x, u) = \omega^T \varphi(x, u)$$

↑
feature vector

Getting back to the TD in terms of this model :

$$e_{TD}^\omega(x, u, x_t, u_t) = \omega^T \varphi(x, u) - r(x, u) - \omega^T \varphi(x_t, u_t)$$

By the approximation theorem,

$$Q^*(x, u) = \omega^* \varphi(x, u) + e_Q(x, u)$$

↑
ideal weights (say, in terms of \sup_{norm})

Let's introduce the weight error:

$$\tilde{\omega} := \omega - \omega^*$$

Get back to the TD once again:

$$e_{TD}^{\omega}(x, u, x_+, u_+) = \omega^T \varphi(x, u) - r(x, u) - \omega^T \varphi(x_+, u_+)$$

$$= \omega^T (\varphi(x, u) - \varphi(x_+, u_+)) - r(x, u)$$

$$= \tilde{\omega}^T (\varphi(x, u) - \varphi(x_+, u_+)) - r(x, u)$$

$$+ \omega^{*T} \varphi(x, u) - \omega^{*T} \varphi(x_+, u_+)$$

$$= \tilde{\omega}^T (\varphi(x, u) - \varphi(x_+, u_+)) - r(x, u)$$

$$+ Q^*(x, u) - e_Q(x, u)$$

$$- Q^*(x_+, u_+) + e_Q(x_+, u_+)$$

Let's denote

$$\ell^*(x, u, x_+, u_+) := Q^*(x, u) - r(x, u) - Q^*(x_+, u_+)$$

If we followed the optimal policy precisely, ℓ^* would equal zero.

So, all in all

$$\ell_{TD}^\omega(x, u, x_+, u_+) = \underbrace{\tilde{\omega}^T d(x, u, x_+, u_+)}_{\text{data vector}}$$

$$- \ell_Q(x, u) + \ell_Q(x_+, u_+) \\ + \ell^*(x, u, x_+, u_+)$$

Now, finally, the critic loss:

$$J^c(\omega | R_u) = \frac{1}{2} \sum_{j=k-M+1}^{k-1} \frac{\left(\ell_{TD}^\omega(x_j, u_j, x_{j+1}, u_{j+1}) \right)^2}{(d_j^T d_j + 1)^2}$$

Shorthand $d_j := d(x_j, u_j, x_{j+1}, u_{j+1})$

Let's learn weights via
gradient descent on the critic loss:

$$\omega_{k+1} := \omega_k - \alpha \nabla_{\omega} J^C(\omega_k | R_k)$$

↓
learning rate

Let's work out the gradient:

$$\nabla_{\omega} J^C(\omega_k | R_k) = \sum_{j=k-M+1}^{k-1} \frac{e^{\omega_k(x_j, u_j, x_{j+1}, u_{j+1})} \nabla_{\omega} e^{\omega_k(x_j, u_j, x_{j+1}, u_{j+1})}}{(d_j^T d_j + 1)^2}$$

In turn,

$$\nabla_{\omega} e^{\omega_k(x_j, u_j, x_{j+1}, u_{j+1})} = d_j$$

So,

$$\nabla_{\omega} J^C(\omega_k | R_k) = \sum_{j=k-M+1}^{k-1} \frac{e^{\omega_k(x_j, u_j, x_{j+1}, u_{j+1})} d_j}{(d_j^T d_j + 1)^2}$$

$$\nabla_{\omega} J^c(\omega_k / R_k) =$$

$$\sum_{j=k-M+1}^{k-1} \frac{d_j d_j^T}{(d_j^T d_j + 1)^2} \tilde{\omega}_k +$$

$$\sum_{j=k-M+1}^{k-1} \frac{d_j (\ell_Q(x_{j+1}, u_{j+2}) - \ell_Q(x_j, u_j) + \ell^*(x_j, u_j, x_{j+1}, u_{j+2}))}{(d_j^T d_j + 1)^2}$$

Denote a data matrix:

$$\mathcal{E}_k := \sum_{j=k-M+1}^{k-1} \frac{d_j d_j^T}{(d_j^T d_j + 1)^2} \geq 0$$

Let's pretend that ℓ_Q was zero
and ℓ^* was zero

We have the evolution of
the weights described by:

$$\omega_{k+1} = \omega_k - d \mathcal{E}_k \tilde{\omega}_k$$

Evidently

$$\tilde{\omega}_{k+1} = \tilde{\omega}_k - d E_k \tilde{\omega}_k$$

We will compare now how the weight error changes from a time step to a time step in norm square:

$$\begin{aligned}\tilde{\omega}_{k+1}^T \tilde{\omega}_{k+1} - \tilde{\omega}_k^T \tilde{\omega}_k &= \\ (\tilde{\omega}_k^T - d E_k \tilde{\omega}_k^T)(\tilde{\omega}_k - d E_k \tilde{\omega}_k) - \tilde{\omega}_k^T \tilde{\omega}_k &= \\ -2d E_k \|\tilde{\omega}_k\|^2 + d^2 E_k^2 \|\tilde{\omega}_k\|^2\end{aligned}$$

Recall:

$$E_k = \sum_{j=k-M+1}^{k-1} \frac{d_j d_j^T}{(d_j^T d_j + 1)^2}$$

$$\text{First, } \|E_k\|^2 \leq \frac{M}{4}$$

Second, we need a

persistence of excitation condition
to hold:

$$+k \quad E_k \geq \epsilon I_{N_c}$$



of critic features

Then,

$$\tilde{w}_{k+1}^T \tilde{w}_{k+1} - \tilde{w}_k^T \tilde{w}_k \leq -2\lambda \epsilon \|\tilde{w}_k\|^2 + \lambda^2 \frac{M}{4} \|\tilde{w}_k\|^2$$

For $\|\tilde{w}_k\|^2$ to "squeeze", we
need to ensure

$$(need) -2\lambda \epsilon \|\tilde{w}_k\|^2 + \lambda^2 \frac{M}{4} \|\tilde{w}_k\|^2 < 0$$

$$\Rightarrow (need) -2\lambda \epsilon + \lambda^2 \frac{M}{4} < 0$$

$$\Rightarrow (need) -2\epsilon + \lambda \frac{M}{4} < 0$$

$$\Rightarrow (\text{need}) \quad d \frac{M}{\bar{y}} < 2\epsilon$$

$$\Rightarrow d < \frac{8\epsilon}{M}$$

This is an ideal scenario, but in reality there are errors ℓ_Q, ℓ^* .

How do they influence the learning?

Looking at

$$\sum_{j=k-M+1}^{k-1} \frac{d_j(\ell_Q(x_{j+1}, u_{j+1}) - \ell_Q(x_j, u_j) + \ell^*(x_j, u_j, x_{j+1}, u_{j+1}))}{(d_j^T d_j + 1)^2}$$

we may write

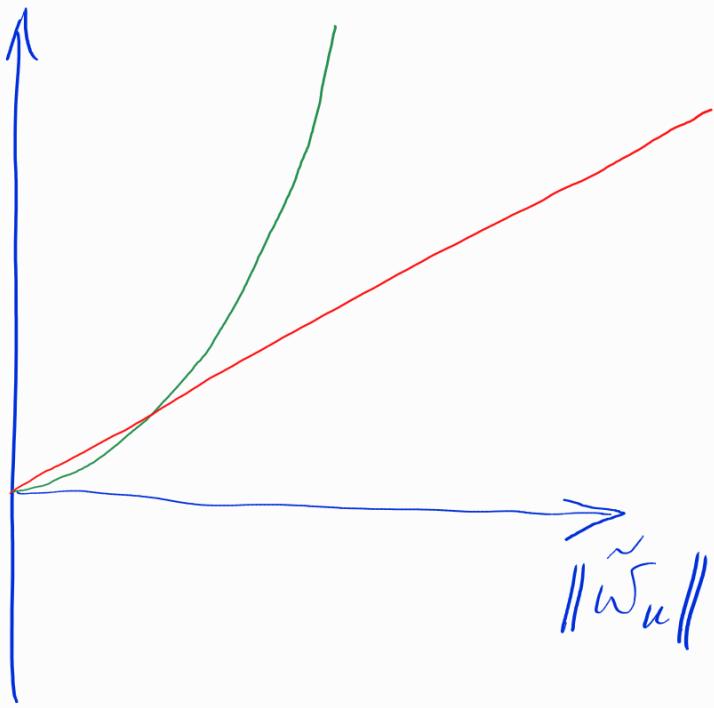
$$\tilde{\omega}_{k+1}^T \tilde{\omega}_{k+1} - \tilde{\omega}_k^T \tilde{\omega}_k =$$

NB!

NB!

$$(\tilde{\omega}_k^T - d\mathcal{E}_k \tilde{\omega}_k + C)(\tilde{\omega}_k - d\mathcal{E}_k \tilde{\omega}_k + C) - \tilde{\omega}_k^T \tilde{\omega}_k \leq$$

$$-2d\mathcal{E} \|\tilde{\omega}_k\|^2 + d^2 \frac{M}{\bar{y}} \|\tilde{\omega}_k\|^2 + \mathcal{O}(d \|\tilde{\omega}_k\|)$$



Actor: $\rho^\theta(x) = \theta^T \varphi(x)$

$$\max_{\theta} \hat{Q}^\omega(x, \theta^T \varphi(x))$$

$$\theta_{u+1} := \theta_u + \beta \nabla_{\theta} \hat{Q}^{\omega_u}(x_u, \theta_u^T \varphi(x_u))$$

$$= \theta_u + \beta \nabla_{\theta} \omega_u^T \varphi(x_u, \theta_u^T \varphi(x_u))$$