

## Remarks on deep RL

### Problem setup

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

s.t.  $X_{k+1} \sim f(x_{k+1} | x_k, u_k), x \in \mathcal{X} \subseteq \mathbb{R}^n$

$$U_k = p^\theta(X_k), u \in \mathcal{U} \subseteq \mathbb{R}^m$$

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

$$X_+^u \sim f(x_+ | x, u)$$

$Q$ -function:

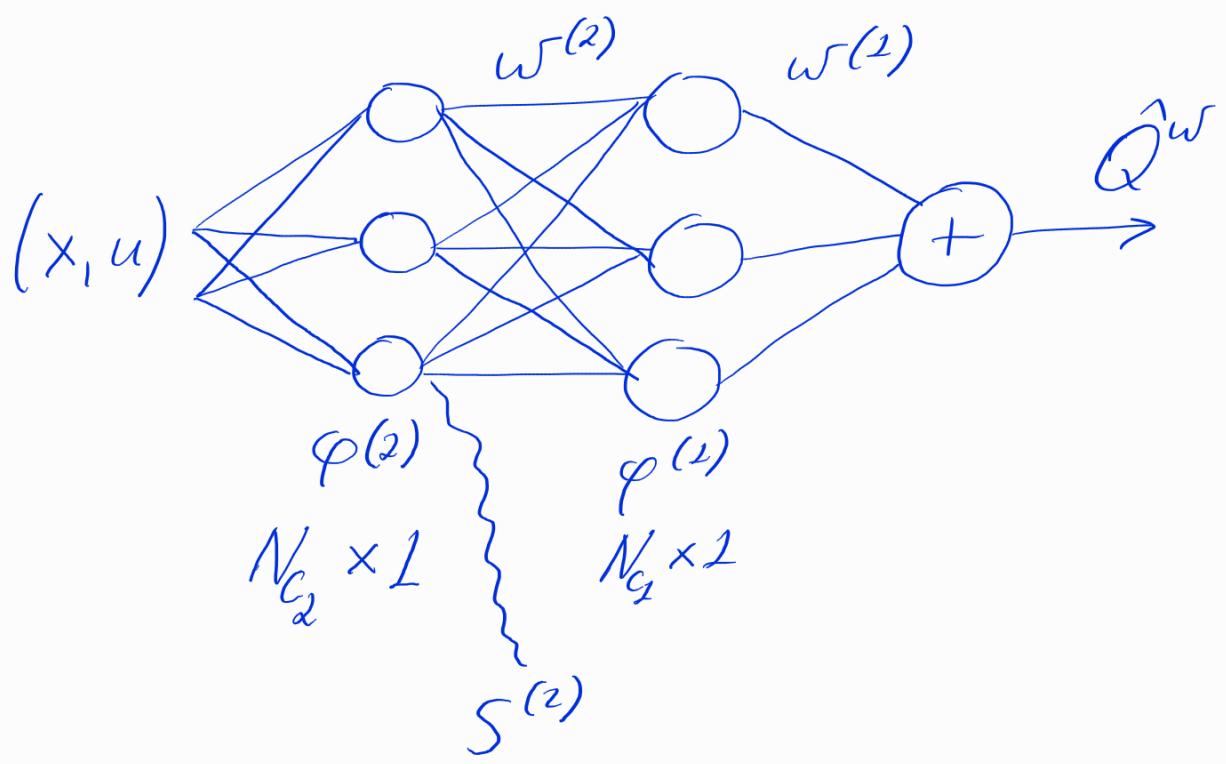
$$Q^*(x, u) = r(x, u) + \gamma \mathbb{E}[V^*(x_+^u)]$$

TD:

$$\mathcal{L}_{TD}^\omega(x, u, x_+, u_+) = \hat{Q}^\omega(x, u) - r(x, u) - \hat{Q}^\omega(x_+, u_+),$$

where

$$\hat{Q}^\omega(x, u) = \omega^{(2)T} \varphi^{(2)} \left( \omega^{(2)T} \varphi^{(2)}(x, u) \right)$$



$\omega^{(1)}$  is  $N_{C_1} \times 1$

$\omega^{(2)}$  is  $N_{C_1} \times N_{C_2}$

---

$z := (x, u)$

So,

$$\mathcal{L}_{TD}^\omega(z, z_+) = \omega^{(1)T} \varphi^{(1)}(\omega^{(2)T} \varphi^{(2)}(z)) - r(z) - \omega^{(1)T} \varphi^{(1)}(\omega^{(2)T} \varphi^{(2)}(z_+))$$

Let's linearize the TD wrt weights.

$$\nabla_{\omega^{(2)}} e^{\omega} \mathcal{L}_{TD}(z, z_+) = \varphi^{(2)} \left( \omega^{(2)T} \varphi^{(2)}(z) \right) - \varphi^{(2)} \left( \omega^{(2)T} \varphi^{(2)}(z_+) \right)$$

Now, gradient wrt hidden layer weights.

First,

$$g(\omega^{(2)}, z, z_+) := \underbrace{\varphi^{(1)}\left(\omega^{(2)^\top} \varphi^{(2)}(z)\right)}_{N_{C_2} \times 1} - \varphi^{(1)}\left(\omega^{(2)^\top} \varphi^{(2)}(z_+)\right)$$

With this at hand, we can write:

$$e_{TD}^{\omega}(z, z_+) = \sum_{\ell=1}^{N_{c_2}} \omega_{\ell}^{(2)} g_{\ell}(\omega^{(2)}, z, z_+) -$$

$r(z)$       ↓

Let's introduce an operator:

`vec(•)` — conversion of a matrix into a column vector

by stacking up the matrix columns

Then,

$$\nabla_{\text{vec}(\omega^{(2)})} \mathcal{E}_{TD}^{\omega}(z, z_+) = \sum_{\ell=1}^{N_{C_2}} \nabla_{\text{vec}(\omega^{(2)})} g_{\ell}(\omega^{(2)}, z, z_+)$$

N\_{C\_2} \cdot N\_{C\_1} \times 1

\omega\_{\ell}^{(2)}

$$\begin{pmatrix} & | \\ \nabla_{\text{vec}(\omega^{(2)})} g_1(\dots) & \end{pmatrix} \omega_1^{(2)} + \dots + \begin{pmatrix} & | \\ & | \\ & \dots \\ & | \end{pmatrix} \omega_{N_{C_2}}^{(2)}$$

↓

$$G = \begin{pmatrix} | & | \\ | & | \\ \dots & \dots \\ | & | \end{pmatrix}$$

←

$$\omega^{(1)} = \begin{pmatrix} \omega_1^{(1)} \\ \vdots \\ \omega_{N_{C_1}}^{(1)} \end{pmatrix} \quad G \omega^{(1)}$$

$$\nabla_{\text{vec}(\omega^{(2)})} \mathcal{E}_{TD}^{\omega}(z, z_+) = G(\omega^{(2)}, z, z_+) \omega^{(1)}$$

Notice that  $G$  depends  
(in general) non-linearly on  
 $\omega^{(2)}$

With these gradients at hand,  
we can now linearize the  
TD :

first, the gradient reads:

$$\nabla_{\text{vec}(\omega)} \ell_{\text{TD}}^{\omega}(z, z_+) = \begin{pmatrix} g(\omega^{(2)}, z, z_+) \\ \varphi^{(1)} \left( \omega^{(2)} \varphi^{(2)}(z) \right) - \varphi^{(1)} \left( \omega^{(2)} \varphi^{(2)}(z_+) \right) \end{pmatrix}$$

Back to linearization:

$$\ell_{TD}^{\omega}(z, z_+) = \overbrace{\tilde{\omega}^{*(2)^\top} g(\omega^{*(2)}, z, z_+) - r(z)}^{\text{"pivot"}} - \\ (\nabla_{\text{vec}(\omega)} \ell_{TD}^{\omega}(z, z_+) \Big|_{\omega=\omega})^\top \text{vec}(\tilde{\omega}) + \\ \mathcal{O}(\|\text{vec}(\tilde{\omega})\|^2)$$

$$\omega = \{\omega^{(1)}, \omega^{(2)}\}$$

$$\Rightarrow \text{vec}(\omega) = (\omega^{(1)}, \text{vec}(\omega^{(2)}))$$

$$\text{vec}(\tilde{\omega}) := \text{vec}(\omega) - \text{vec}(\omega^*)$$

Following the recipe of the lecture on shallow actor-critic convergence, we can extract NN errors and error related to the „quality“ of the policy.

$$\ell_{TD}^{\omega}(z, z_+) = (\nabla_{\text{vec}(\omega)} \ell_{TD}^{\omega}(z, z_+) \Big|_{\omega})^\top \text{vec}(\tilde{\omega}) - \\ \ell_Q(z) + \ell_Q(z_+) + \ell^*(z, z_+) + \\ \mathcal{O}(\|\text{vec}(\tilde{\omega})\|^2) !$$

Now, data vector

$$d(z_1, z_t, \omega) := \nabla_{\text{vec}(\omega)} \ell_{TD}^{\omega}(z_1, z_t) / \omega$$

$\Rightarrow$

$$\ell_{TD}^{\omega}(z, z_t) = d^T(z, z_t, \omega) \text{vec}(\tilde{\omega}) - \dots$$



Is this bad?

< repeating the setup from the  
shallow NN case >

experience replay



$$\nabla_{\omega} J_c(\omega_u | R_u) =$$

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

$$\sum_{j=k-M+1}^{k-1} \frac{d_j (\ell_Q(z_{j+1}) - \ell_Q(z_j) + e^*(z_j, z_{j+1}))}{(d_j^T d_j + 1)^2} +$$

$$\sum_{j=k-M+1}^{k-1} \frac{d_j \mathcal{O}(\|\text{vec}(\tilde{\omega}_u)\|^2)}{(d_j^T d_j + 1)^2}$$

$$\begin{aligned}
 & \text{vec}(\tilde{\omega}_{k+1})^T \text{vec}(\tilde{\omega}_{k+1}) - \text{vec}(\tilde{\omega}_k)^T \text{vec}(\tilde{\omega}_k) \leq \\
 & \quad \underbrace{\leq -c \|\text{vec}(\tilde{\omega}_k)\|^2}_{-2\lambda E \|\text{vec}(\tilde{\omega}_k)\|^2 + \lambda^2 \frac{M}{q} \|\text{vec}(\tilde{\omega}_k)\|^2 +} \\
 & \quad O(\lambda \|\text{vec}(\tilde{\omega}_k)\|) + O(\lambda \|\text{vec}(\tilde{\omega}_k)\|^3)
 \end{aligned}$$

\*  $c$  comes from the learning rate choice for convergence

