

## Lecture 4. TD &amp; AC

$$\max_{\theta} J^\theta(x_0) = \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k r(X_k, U_k) \right]$$

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

$$U_k \sim p^\theta(u_k | x_k)$$

The HJB for the stated problem reads:

$$\forall x \in \mathcal{X} \quad V^*(x) = \max_u \{ r(x, u) + \gamma \mathbb{E}[V^*(X_+)] \}$$

where  $X_+^{(u)} \sim f(x_{k+1} | x_k, u)$

Verification principle: if a function  $V$  fulfills the HJB, then it is the optimal value function:  $V = V^*$

$$V^*(x_k) = \max_u \{ r(x_k, u) + \gamma \mathbb{E}[V^*(X_{k+1}^{(u)})] \}$$

The idea of TD is, roughly, to „equalize the lhs of the HJB with its rhs under a model  $\hat{V}^w$  with weights  $w$ “.

The goal is  $\hat{V}^w \approx V^*$

Let's recall tabular value iteration:

suppose  $\mathcal{X}_0 = \{x^j\}_{j=1}^N$

Init  $V_0 \equiv 0$ ,  $p_0 \equiv 0$

Loop over  $i$

$$\forall x^i \in \mathcal{X}_0 \quad V_{i+1}(x^i) := r(x^i, p_i(x^i)) + \gamma \mathbb{E} \left[ V_i(x^{p_i(x^i)}) \mid X = x^i \right]$$

$$\forall x^i \in \mathcal{X}_0 \quad p_{i+1}(x^i) := \arg \max_u \left\{ r(x^i, u) + \gamma \mathbb{E} \left[ V_{i+1}(x^u) \right] \right\}$$

Instead of  $V_i$ , let's use a model  $\hat{V}^w$  and so the above algorithm will be cast to, while also parametrizing the policy, i.e.,  $\rho^\theta$ :

Init  $w_0, \theta_0$

Loop over time steps  $k$

$$w_{k+1} := \arg \min_w \left\{ \frac{1}{2} \left( \hat{V}^w(x_k) - r(x_k, \rho^{\theta_k}(x_k)) - \gamma \mathbb{E} \left[ \hat{V}^{w_k} \left( X_{k+1}^{\rho^{\theta_k}(x_k)} \right) \right] \right)^2 \right\}$$

$$\theta_{k+1} := \arg \max_\theta \left\{ r(x_k, \rho^\theta(x_k)) + \gamma \mathbb{E} \left[ \hat{V}^{w_{k+1}} \left( X_{k+1}^{\rho^\theta(x_k)} \right) \right] \right\}$$

$\hat{V}^w$  is associated with the so-called **critic**

$\beta^\theta$  is associated with the so-called actor

The difference in the (...) is referred to as the temporal difference

## Upgrades and modifications

Action-value  $Q^\theta(x, u)$

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

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

$$\beta^*(x) = \underset{u}{\operatorname{argmax}} Q^*(x, u)$$

Casting this into a model

$$\hat{Q}^\omega(x, u), \text{ e.g., } \langle \omega, (\varphi_x(x), \varphi_u(u), \varphi_{xu}(x, u)) \rangle$$

say,  $\varphi_{xu}(x, u) = x \cdot u + x^2 \cdot u^2 + \dots$

## Action-value update

$$\omega_{k+2} := \arg \min_{\omega} \frac{1}{2} TD_{k-2}^2$$

$$TD_{k-2} = \hat{Q}^{\omega}(x_{k-2}, u_{k-2}) - r(x_{k-2}, u_{k-2}) \\ - \hat{Q}^{\omega_k}(x_k, u_k)$$

$$TD_{k-2}$$

...

over an experience replay

## policy update

$$\rho^{\theta_{k+2}}(x_k) := \arg \max_{\theta} \hat{Q}^{\omega_{k+2}}(x_k, \rho^\theta(x_k))$$

- Update target

was  $r(x, u) + \gamma \underbrace{\mathbb{E}[\hat{V}^{\omega_-}(x_+)]}_{Q^{\omega_-}(x, u)}$

$$(TD(\cdot))$$

$TD(N) :$

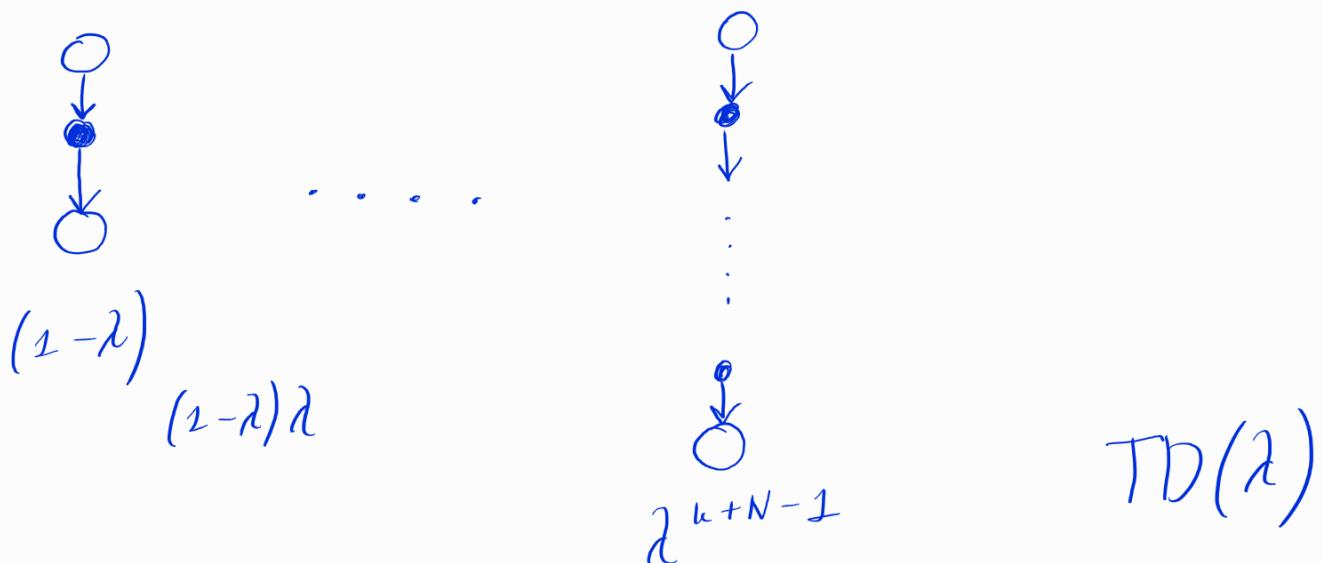
$$O_N^{\omega^-}(x, u) := r(x, u) + \mathbb{E} \left[ \sum_{k=1}^{N-1} \gamma^k r(X_{+k}, U_{+k}) + \gamma^N \hat{V}^{\omega^-}(X_{+N}) \right],$$

where  $X_{+k}$  is the state  $k$  steps ahead of  $x$

There is also

$$O_\lambda^{\omega^-}(x, u) = (1-\lambda) \sum_{k=1}^{N-1} \lambda^{k-1} O_k^{\omega^-}(x, u) + \lambda^{N-1} \sum_{j=1}^N \gamma^{j-1} \mathbb{E}[r(X_{+j}, U_{+j})]$$

Visualization : (backup diagrams)



## • Eligibility traces

Introduce  $\mathcal{Z}_k$  (a vector)

$$\mathcal{Z}_0 := \emptyset$$

Then, for  $k$

$$\mathcal{Z}_{k+1} := \gamma \lambda \mathcal{Z}_k + \nabla_{\omega} \hat{V}_{\omega=\omega_k}$$