

Dudén Moreno-Bate dates

C/Neuro 2024

Beijing

Bellman eq

$$V_A(s) = E_{\pi} \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \quad z = (z_0 = s_0, a_0, z_1 = s_1, a_1, z_2 = s_2, a_2, \dots)$$

$$p(z|s_0) = \pi(a_0|s_0) p(s_1|s_0, a_0) \pi(a_1|s_1) p(s_2|s_1, a_1) \dots$$

$$V_{\pi}(s_0) = \sum_z p(z|s_0) \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$$

$$= \sum_z p(z|s_0) [r(s_0, a_0) + \gamma r(s_1, a_1) + \gamma^2 r(s_2, a_2) + \dots]$$

$$= \sum_z p(z|s_0) r(s_0, a_0) + \gamma \sum_z p(z|s_0) [r(s_1, a_1) + \gamma r(s_2, a_2) + \dots]$$

$$= \sum_{a_0} \pi(a_0|s_0) r(s_0, a_0) + \gamma \sum_{a_0, s_1, z_1} \pi(a_0|s_0) p(s_1|s_0, a_0) p(z_1|s_1) [r(s_1, a_1) + \gamma r(s_2, a_2) + \dots]$$

$$= \sum_{a_0} \pi(a_0|s_0) r(s_0, a_0) + \gamma \sum_{s_1} \pi(a_0|s_0) p(s_1|s_0, a_0) \sum_{z_1} p(z_1|s_1) [r(s_1, a_1) + \gamma r(s_2, a_2) + \dots]$$

$$= \sum_{a_0} \pi(a_0|s_0) r(s_0, a_0) + \gamma \sum_{s_1} \pi(a_0|s_0) p(s_1|s_0, a_0) V_{\pi}(s_1)$$

$s_0 \rightarrow s$
 $s_1 \rightarrow s'$

$$V_{\pi}(s) = \sum_a \pi(a|s) [r(s, a) + \gamma \sum_{s'} p(s'|s, a) V_{\pi}(s')]$$

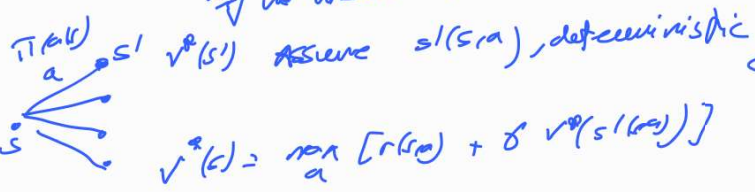
$$= E_{\pi} [r(s, a)] + \gamma E_{\pi} [V_{\pi}(s')]$$

Bellman equation

optimality Bellman equation

$$V^*(s) = \max_a [r(s, a) + \gamma \sum_{s'} p(s'|s, a) V^*(s')]$$

if no zero $V^*(s')$ & s' children of s



may be explicit the first Max the optimality Bellman eq.

$$V^*(s) = \max_a [r(s, a) + \gamma \sum_{s'} p(s'|s, a) V^*(s')]$$

How do we find $V^*(s)$?

$0 \leq r(s, a) < \infty$, start with initial guess $V_0(s) = 0 \forall s$

$$V_1(s) = \max_a [r(s, a) + \gamma \sum_{s'} p(s'|s, a) V_0(s')] = \max_a [r(s, a)] \geq 0$$

$$V_2(s) \geq V_1(s) \forall s$$

Assume $V_{t+1}(s) \geq V_t(s)$.
Do it over that $V_{t+1}(s) \geq V_t(s)$?

$$V_{t+1}(s) = \max_a [r(s, a) + \gamma \sum_{s'} p(s'|s, a) V_t(s')]$$

$$\geq \max_a [r(s, a) + \gamma \sum_{s'} p(s'|s, a) V_{t-1}(s')] = V_t(s)$$

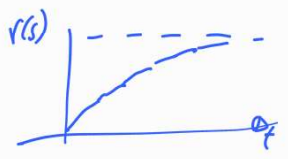
By induction

$$\Rightarrow V_{t+1}(s) \geq V_t(s) \quad \forall s, t$$

↓ non-decreasing.

$V^*(s)$ is bounded for $\gamma < 1$

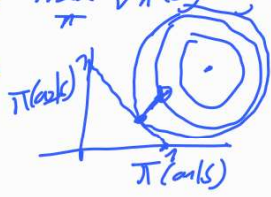
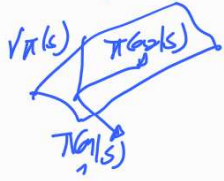
$$\Rightarrow V^*(s) \leq \frac{\Gamma_{max}}{1-\gamma}$$



Optimal value function and Policy in Soft RL

$$V_{\pi}(s) = \sum_a \pi(a|s) \left[r(s,a) - \alpha \log \frac{\pi(a|s)}{\pi_0(a|s)} + \gamma \sum_{s'} p(s'|s,a) V_{\pi}(s') \right]$$

$$V^*(s) = \max_{\pi} V_{\pi}(s) \quad \pi^* = \arg \max_{\pi} V_{\pi}(s)$$



$$\frac{\partial V_{\pi}(s)}{\partial \pi(a|s)} = \lambda(s) \quad \forall a$$

$$\frac{\partial V_{\pi}(s)}{\partial \pi(a|s)} = \lambda(s)$$

Critical π

$$\frac{\partial V_{\pi}(s)}{\partial \pi(a|s)} = r(s,a) - \alpha \log \frac{\pi(a|s)}{\pi_0(a|s)} + \gamma \sum_{s'} p(s'|s,a) V_{\pi}(s') - \alpha \frac{\pi(a|s)}{\pi_0(a|s)} + \gamma \sum_{a'} \pi(a'|s) \sum_{s'} p(s'|s,a') \frac{\partial V_{\pi}(s')}{\partial \pi(a|s)} = \lambda(s)$$

$$\Rightarrow \pi^*(a|s) \propto \pi_0(a|s) e^{\frac{1}{\alpha} \left[r(s,a) + \gamma \sum_{s'} p(s'|s,a) V^*(s') \right]}$$

critical \Rightarrow (critical points of boundaries of simplex) uniqueness