



# The suboptimality of $\mu_T^*$ with respect to $\mu^*$ is bounded

---

Damien Ernst, François Rozet and Valentin Vermeylen

February 11, 2021

University of Liège – School of Engineering



Let consider a **deterministic** environment described by the **dynamics**

$$f : X \times U \mapsto X$$

and the **reward** function

$$r : X \times U \mapsto \mathbb{R}^+$$

where  $X$  is the **state space** and  $U$  the **action space**.



The **return**  $J^\mu : X \mapsto \mathbb{R}^+$  of a **stationary policy**  $\mu : X \mapsto U$  is defined as

$$J^\mu(x) = \lim_{T \rightarrow \infty} \sum_{t=0}^T \gamma^t r(x_t, u_t) \quad (1)$$

with  $u_t = \mu(x_t)$ ,  $x_{t+1} = f(x_t, u_t)$  and  $x_0 = x$ . An interesting property is that

$$\|J^\mu\|_\infty \leq \frac{B_r}{1 - \gamma} \quad (2)$$

where  $B_r = \|r\|_\infty$  and iff  $\gamma \in [0; 1)$ . Indeed,

$$\lim_{T \rightarrow \infty} \sum_{t=0}^T \gamma^t = \lim_{T \rightarrow \infty} \frac{1 - \gamma^{T+1}}{1 - \gamma} = \frac{1}{1 - \gamma}$$

---

$$\|F\|_\infty = \sup_{y \in Y} |F(y)|, \quad \forall F : Y \mapsto \mathbb{R}$$



We define the **truncated** return  $J_T^\mu : X \mapsto \mathbb{R}^+$  as the recurrence

$$J_T^\mu(x) = r(x, \mu(x)) + \gamma J_{T-1}^\pi(f(x, \mu(x))), \quad \forall T \geq 1 \quad (3)$$

with  $J_0^\mu(x) \equiv 0$ .



By definition, we have

$$J^\mu(x) = \lim_{T \rightarrow \infty} J_T^\mu(x) \quad (4)$$

for all  $x \in X$ . Similarly, we can also write

$$J^\mu(x) = J_T^\mu(x) + \gamma^T J^\mu(x_T) \quad (5)$$

which leads, using (2),

$$\|J^\mu - J_T^\mu\|_\infty \leq \gamma^T \|J^\mu\|_\infty \leq \frac{\gamma^T B_r}{1 - \gamma} \quad (6)$$



We define the **state-action value**  $Q_T : X \times U \mapsto \mathbb{R}^+$  by the recurrence

$$Q_T(x, u) = r(x, u) + \gamma \max_{u' \in U} Q_{T-1}(f(x, u), u'), \quad \forall T \geq 1 \quad (7)$$

with  $Q_0(x) \equiv 0$ . Similarly to  $J^\mu$ , we have

$$Q(x, u) = \lim_{T \rightarrow \infty} Q_T(x, u) \quad (8)$$

which is the state-action value over an **infinite number of steps**.



A stationary policy  $\mu^*$  is **optimal** iff it selects an optimal action when there remains an **infinite** number of steps.

$$\mu^*(x) \in \arg \max_{u \in U} Q(x, u) \quad (9)$$

or, equivalently,

$$Q(x, \mu^*(x)) = \max_{u \in U} Q(x, u)$$

Thus,

$$\begin{aligned} Q(x, \mu^*(x)) &= r(x, \mu^*(x)) + \gamma \max_{u \in U} Q(f(x, \mu^*(x)), u) \\ &= r(x, \mu^*(x)) + \gamma Q(x', \mu^*(x')) \\ &= r(x, \mu^*(x)) + \gamma r(x', \mu^*(x')) + \gamma^2 Q(x'', \mu^*(x'')) \\ &= J^{\mu^*}(x) \end{aligned}$$

with  $x' = f(x, \mu^*(x))$ ,  $x'' = f(x', \mu^*(x'))$ , ...



In contrary, a stationary policy is  $T$ -optimal if it selects an optimal action when there remains exactly  $T$  steps.

$$\mu_T^*(x) \in \arg \max_{u \in U} Q_T(x, u) \quad (10)$$

or, equivalently,

$$Q_T(x, \mu_T^*(x)) = \max_{u \in U} Q_T(x, u)$$

Necessarily,  $\mu_T^*$  is suboptimal with respect to  $\mu^*$ , i.e.

$$J^{\mu^*}(x) \geq J^{\mu_T^*}(x) \quad (11)$$





However, choosing  $u_t = \mu_{T-t}^*(x_t)$  for all  $t \leq T$  is **optimal**.

Then, we define

$$J_T^{\pi^*}(x) = \sum_{t=0}^T \gamma^t r(x_t, \mu_{T-t}^*(x_t)) \quad (12)$$

or, recurrently,

$$J_T^{\pi^*}(x) = r(x, \mu_T^*(x)) + \gamma J_{T-1}^{\pi^*}(f(x, \mu_T^*(x))) \quad (13)$$

with  $J_0^{\pi^*}(x) \equiv 0$ . Interestingly,

$$J_T^{\pi^*}(x) = \max_{u \in U} Q_T(x, u) \geq J_T^{\mu^*}(x)$$

---

The notation  $\pi$  indicates a non-stationary policy.



- Optimal policy  $\mu^*$
- $T$ -optimal policy  $\mu_T^*$
- Return of the optimal policy  $J^{\mu^*}$
- Return of the  $T$ -optimal policy  $J^{\mu_T^*}$
- Truncated return of the optimal policy  $J_T^{\mu^*}$
- Optimal truncated return  $J_T^{\pi^*}$



The suboptimality of  $\mu_T^*$  with respect to  $\mu^*$  is **bounded**.

$$\|J^{\mu^*} - J^{\mu_T^*}\|_{\infty} \leq \frac{2\gamma^T B_r}{(1-\gamma)^2} \quad (14)$$



By definition and given (13), we have

$$J^\mu(x) = r(x, \mu(x)) + \gamma J^\mu(f(x, \mu(x)))$$

$$\begin{aligned} J_T^{\pi^*}(x) &= r(x, \mu_T^*(x)) + \gamma J_{T-1}^{\pi^*}(f(x, \mu_T^*(x))) \\ &\geq r(x, \mu^*(x)) + \gamma J_{T-1}^{\pi^*}(f(x, \mu^*(x))) \end{aligned}$$

Therefore,

$$\begin{aligned} J^{\mu^*}(x) - J^{\mu_T^*}(x) &\leq J^{\mu^*}(x) - [r(x, \mu^*(x)) + \gamma J_{T-1}^{\pi^*}(f(x, \mu^*(x)))] \\ &\quad + [r(x, \mu_T^*(x)) + \gamma J_{T-1}^{\pi^*}(f(x, \mu_T^*(x)))] - J^{\mu_T^*}(x) \\ &\leq \gamma [J^{\mu^*}(f(x, \mu^*(x))) - J_{T-1}^{\pi^*}(f(x, \mu^*(x)))] \\ &\quad + \gamma [J_{T-1}^{\pi^*}(f(x, \mu_T^*(x))) - J^{\mu_T^*}(f(x, \mu_T^*(x)))] \end{aligned}$$



Thus, in norm,

$$\begin{aligned}
 \|J^{\mu^*} - J^{\mu_T^*}\|_{\infty} &\leq \gamma \|J^{\mu^*} - J_{T-1}^{\pi^*}\|_{\infty} + \gamma \|J_{T-1}^{\pi^*} - J^{\mu_T^*}\|_{\infty} \\
 &\leq \gamma \|J^{\mu^*} - J_{T-1}^{\pi^*}\|_{\infty} \\
 &\quad + \gamma \|J_{T-1}^{\pi^*} - J^{\mu^*} + J^{\mu^*} - J^{\mu_T^*}\|_{\infty} \\
 &\leq 2\gamma \|J^{\mu^*} - J_{T-1}^{\pi^*}\|_{\infty} + \gamma \|J^{\mu^*} - J^{\mu_T^*}\|_{\infty} \\
 &\leq \frac{2\gamma}{1-\gamma} \|J^{\mu^*} - J_{T-1}^{\pi^*}\|_{\infty}
 \end{aligned}$$

But since

$$\begin{aligned}
 J^{\mu^*}(x) - J_T^{\pi^*}(x) &= J_T^{\mu^*}(x) + \gamma^T J^{\mu^*}(x_t) - J_T^{\pi^*}(x) \\
 &\leq \gamma^T J^{\mu^*}(x_t)
 \end{aligned}$$



We have

$$\begin{aligned}\|J^{\mu^*} - J^{\mu_T^*}\|_{\infty} &\leq \frac{2\gamma}{1-\gamma} \|J^{\mu^*} - J_{T-1}^{\pi^*}\|_{\infty} \\ &\leq \frac{2\gamma}{1-\gamma} \gamma^{T-1} \|J^{\mu^*}\|_{\infty} \\ &\leq \frac{2\gamma^T}{1-\gamma} \frac{B_r}{1-\gamma}\end{aligned}$$





Damien Ernst, Mevludin Glavic, Florin Capitanescu, and Louis Wehenkel. “Reinforcement learning versus model predictive control: a comparison on a power system problem”. In: *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 39.2 (2008), pp. 517–529.