# Revisiting Peng's Q($\lambda$): Good Old-Fashioned Algorithm for Modern Reinforcement Learning

Tadashi Kozuno[1]*, Yunhao Tang[2]*, Mark Rowland[3], Remi Munos[4], Steven Kapturowski[3], Will Dabney[3], Michal Valko[4], David Abel[3]

*: equal contribution

1: Univ of Alberta, 2: Columbia Univ, 3: DeepMind London, 4: DeepMind Paris

# Value Iteration (VI)

Initialize a Q-function $Q$;

Initialize a behavior policy $\mu$;

For $t = 1, \ldots, T$ do;

    Observe $s_t$ and take $a_t \sim \mu(\cdot | a_t)$.

    Observe $s_{t+1}$ and receive $r_t$.

    Store $(s_t, a_t, s_{t+1}, r_t)$ in data buffer $\mathcal{D}$.

    If $t \bmod T_{update} = 0$ then

$$Q \leftarrow \operatorname{argmin}_f \mathbb{E}_{\mathcal{D}} \left[ \left( f(s, a) - \hat{Q}_{VI} \right)^2 \right].$$

$$\mu \leftarrow \text{NewBehaviorPolicy}(Q).$$

Initialize a Q-function $Q$;

Initialize a behavior policy $\mu$;

For $t = 1, \ldots, T$ do;

    Observe $s_t$ and take $a_t \sim \mu(\cdot | a_t)$.

    Observe $s_{t+1}$ and receive $r_t$.

    Store $(s_t, a_t, s_{t+1}, r_t)$ in data buffer $\mathcal{D}$.

    If $t \bmod T_{update} = 0$ then

      $Q \leftarrow \operatorname{argmin}_f \mathbb{E}_{\mathcal{D}} \left[ \left( f(s, a) - \hat{Q}_{VI} \right)^2 \right]$.

      $\mu \leftarrow \text{NewBehaviorPolicy}(Q)$.

$$Q_{VI} := r + \gamma \max_{a'} Q(s', a')$$

$\lambda$ Policy Iteration:

$$Q_{\lambda PI} := (1 - \lambda) \sum_{N=1}^{\infty} \lambda^{N-1} \left( \sum_{n=1}^{N} \gamma^{n-1} r_{t+n} + \gamma^N \max_a Q(s_{t+N+1}, a) \right)$$

$$= Q(s_t, a_t) + \sum_{N=0}^{\infty} \gamma^N \lambda^N \left( r_{t+N} + \gamma \max_a Q(s_{t+N+1}, a) - Q(s_{t+N}, a_{t+N}) \right)$$

$(s_{t+n}, a_{t+n})_{n=0,1,...,}$ is obtained by following a greedy policy!

Importance sampling $\lambda$ -PI:

$$Q_{IS\,\lambda PI} := Q(s_t, a_t) + \sum_{N=0}^{\infty} c_N \gamma^N \lambda^N \left( r_{t+N} + \gamma \max_a Q(s_{t+N+1}, a) - Q(s_{t+N}, a_{t+N}) \right)$$

where $c_0 := 1$, and $c_N := \prod_{n=0}^{N} \frac{\pi(a_{t+n}|s_{t+n})}{\mu(a_{t+n}|s_{t+n})}$.

Importance sampling $\lambda$ -PI:

$$Q_{IS\,\lambda PI} := Q(s_t, a_t) + \sum_{N=0}^{\infty} c_N \gamma^N \lambda^N \left( r_{t+N} + \gamma \max_a Q(s_{t+N+1}, a) - Q(s_{t+N}, a_{t+N}) \right)$$

where $c_0 := 1$, and $c_N := \prod_{n=0}^{N} \frac{\pi(a_{t+n}|s_{t+n})}{\mu(a_{t+n}|s_{t+n})}$.

Because of $c_N$, traces are cut when $\pi(a_{t+n}|s_{t+n}) \approx 0$.

Conservative Algorithms

- Importance sampling $\lambda$-PI
- Retrace (Munos et al., 2016)
- Watkin's Q($\lambda$) (Watkins, 1989)
- Tree-backup (Precup et al., 2000)

Convergence thanks to trace cuts.

# Conservative and Non-Conservative Algorithms

Conservative Algorithms

- Importance sampling $\lambda$ -PI
- Retrace (Munos et al., 2016)
- Watkin's Q($\lambda$) (Watkins, 1989)
- Tree-backup (Precup et al., 2000)

Convergence thanks to trace cuts.

Non-conservative Algorithms

- Uncorrected n-step return
- Peng's Q($\lambda$) (Peng & Williams, 1994)
- Harutyunyan's Q($\lambda$) (Harutyunyan et al., 2016)

No generic convergence guarantee.

But some of them are known to yield a better performance.

Why do non-conservative algorithms outperform theoretically-sound conservative ones?

# Theoretical Result

- Convergence guarantee of Peng's Q($\lambda$) with fixed $\mu$
  - Fast convergence (to a biased optimal solution)

  - More robustness to error (compared to VI, $\lambda$-PI)

- Convergence guarantee of Peng's Q($\lambda$) with changing $\mu$
  - Convergence under appropriate conditions on $\mu$

  - Slower convergence (to the optimal solution)

  - Less robustness to error (compared to VI, $\lambda$-PI)