

From Dirichlet to Rubin: Optimistic Exploration in RL without Bonuses

Daniil Tiapkin^{1,2}, Denis Belomestny^{3,1}, Éric Moulines^{4,1}, Alexey Naumov¹, Sergey Samsonov¹, Yunhao Tang⁵, Michal Valko⁵, Pierre Ménard⁶

¹HSE University, ²Artificial Intelligence Research Institute, ³Duisburg-Essen University, ⁴École Polytechnique, ⁵DeepMind, ⁶Otto von Guericke University

Overview

- **Bayes-UCBVI** is Bayes-inspired algorithm with nearly minimax regret bound for large enough T ;
- Novel anti-concentration inequality for weighted sums of Dirichlet random vector;
- Extension of exploration mechanism of **Bayes-UCBVI** in Deep RL setting;

Setting

- Tabular MDP: H horizon, S states, A actions, $p_h(s'|s, a)$ unknown transitions, deterministic reward $r_h(s, a) \in [0, 1]$.
- Regret: $\sum_{t=1}^T V_1^*(s_1) - V_1^{\pi^t}(s_1)$

Bayes-UCBVI

- Bonus-based exploration (theoretically near optimal, poor empirical performance, not scalable)

$$\bar{Q}_h^t(s, a) = r_h(s, a) + \hat{p}_h^t \bar{V}_{h+1}^t(s, a) + B_h^t(s, a), \quad \bar{V}_h^t(s) = \max_a \bar{Q}_h^t(s, a),$$

where $\hat{p}_h^t(s, a)$ is empirical transition probabilities, and $\hat{p}_h^t f(s, a) \triangleq \sum_{s' \in \mathcal{S}} \hat{p}_h^t(s'|s, a) f(s')$.

- **Bayes-UCBVI** exploration (theoretically near optimal, great empirical performance, scalable)

$$\bar{Q}_h^t(s, a) = r_h(s, a) + \text{Quantile}_{p \sim \rho_h^t(s, a)}(p \bar{V}_{h+1}^t, \kappa), \quad \bar{V}_h^t(s) = \max_a \bar{Q}_h^t(s, a),$$

where $\rho_h^t(s, a) = \text{Dir}(n_0, n_h^t(s_1, s, a), \dots, n_h^t(s_S, s, a))$ is a posterior distribution of transition probabilities.

Optimistic prior for **Bayes-UCBVI**:

- Add an artificial isolated state s_0 with $r_h(s_0, a) > 1$;
- Add n_0 pseudo-transitions from each state s to s_0 into the history of visits.

Regret bounds

Algorithm	Upper bound (non-stationary)
UCBVI [Azar et al., 2017]	
UCB-Advantage [Zhang et al., 2020]	$\tilde{O}(\sqrt{H^3 SAT})$
RLSVI [Xiong et al., 2021]	
PSRL [Agrawal and Jia, 2017]	$\tilde{O}(H^2 S \sqrt{AT})$
BootNARL [Pacchiano et al., 2021]	
Bayes-UCBVI (this paper)	$\tilde{O}(\sqrt{H^3 SAT})$
Lower bound [Jin et al., 2018, Domingues et al., 2021]	$\Omega(\sqrt{H^3 SAT})$

Green: scalable, Orange: scalable under simplifications, Red: not scalable.

Non-tabular extension: Bayes-UCBDQN

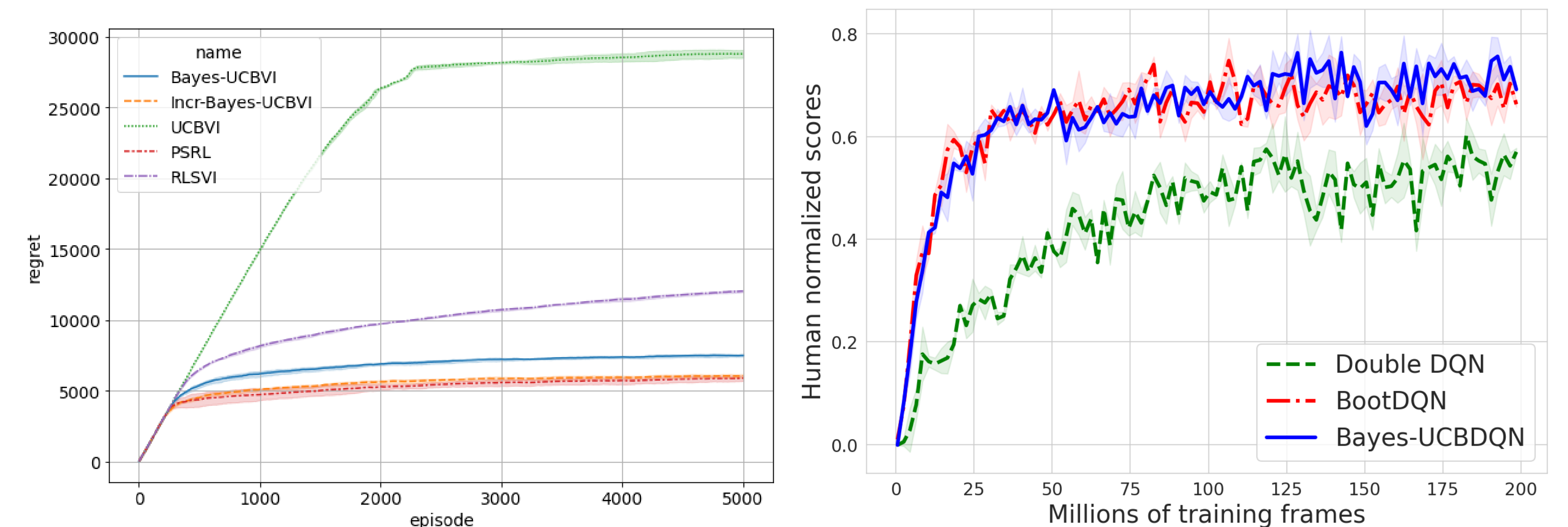


Figure 1: Left: Regret of **Bayes-UCBVI** and **Incr-Bayes-UCBVI** compared to baselines. Environment: grid-world with 5 rooms of size 5×5 ($S = 129$, $A = 4$).

Right: Evaluating deep RL algorithms with median human normalized scores across Atari-57 games.

- targets for Q-function estimation $y^n = r_h(s, a) + \bar{V}_{h+1}^t(s_{h+1}^n)$, $n = 1, \dots, n_h^t(s, a)$;
- targets from prior transitions $y^n = r_h(s, a) + \bar{V}_h^t(s_0)$, $n = -n_0 + 1, \dots, 0$.

Reformulate and approximate **Bayes-UCBVI** using B Bayesian bootstrap samples

$$\begin{aligned} \bar{Q}_h^t(s, a) &= \text{Quantile}_{w \sim \text{Dir}(\underbrace{1, \dots, 1}_{n_h^t(s, a) + n_0})} \left(\sum_{n=-n_0+1}^{n_h^t(s, a)} w^n y^n, \kappa \right) \\ &\approx \text{Quantile}_{b \sim \text{Unif}([1, B])} \left(\bar{Q}_h^{t, b}(s, a), \kappa \right), \\ \text{where } \bar{Q}_h^{t, b}(s, a) &= \sum_{n=-n_0+1}^{n_h^t(s, a)} w^{n, b} y^n \text{ Bayesian bootstrap sample with } w^{n, b} \sim \text{Dir}(\underbrace{1, \dots, 1}_{n_h^t(s, a) + n_0}). \end{aligned}$$

Since $w^{n, b}$ is a vector of normalized $z^{n, b} \sim \mathcal{E}(1)$ random variables

$$\bar{Q}_h^{t, b}(s, a) = \underset{x}{\text{argmin}} \sum_{n=-n_0+1}^{n_h^t(s, a)} z^{n, b} (x - y^n)^2.$$

Anti-concentration inequality for Dirichlet distribution

For any $\alpha = (\alpha_0, \alpha_1, \dots, \alpha_m) \in \mathbb{N}^{m+1}$ define $\bar{p} \in \Delta_m$ with $\bar{p}(\ell) = \alpha_\ell / \bar{\alpha}$, $\ell = 0, \dots, m$, where $\bar{\alpha} = \sum_{j=0}^m \alpha_j$. Under technical assumptions, for $f: \{0, \dots, m\} \rightarrow [0, b_0]$ and $\mu \in (\bar{p}f, b_0)$

$$\bar{\alpha}^{-3/2} \exp(-\bar{\alpha} \mathcal{K}_{\text{inf}}(\bar{p}, \mu, f)) \leq \mathbb{P}_{w \sim \text{Dir}(\alpha)}[wf \geq \mu] \leq \exp(-\bar{\alpha} \mathcal{K}_{\text{inf}}(\bar{p}, \mu, f)),$$

where $\mathcal{K}_{\text{inf}}(p, u, f)$ is given by

$$\mathcal{K}_{\text{inf}}(p, u, f) \triangleq \max_{\lambda \in [0, 1]} \mathbb{E}_{X \sim p} \left[\log \left(1 - \lambda \frac{f(X) - u}{b_0 - u} \right) \right].$$

- Lower bound is an essential part for optimism;
- Upper bound is important for the reduction to UCBVI.