

# UCB Momentum Q-learning: Correcting the bias without forgetting

Pierre Ménard<sup>1</sup>, Omar Darwiche Domingues<sup>2</sup>, Xuedong Shang<sup>2,3</sup>, Michal Valko<sup>2,4</sup>

<sup>1</sup>Otto von Guericke University Magdeburg, <sup>2</sup>Inria Lille, <sup>3</sup>Université de Lille, <sup>4</sup>DeepMind

## Overview

- Analyze the benefits of adding a momentum term to Q-learning in the episodic setting.
- UCBMQ algorithm with regret bound that scales *linearly* with the number of states  $S$ .

## Setting

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

## Intuition

- Generic Q-learning step:

$$Q_h^n(s, a) = \alpha_n(r_h + p_h^n \bar{V}_h^{n-1})(s, a) + (1 - \alpha_n)Q_h^{n-1}(s, a)$$

$$\bar{Q}_h^n(s, a) = Q_h^n(s, a) + b_h^n(s, a) \quad \bar{V}_h^n(s) = \max_a \bar{Q}_h^n(s, a)$$

where the sample expectation  $(p_h^n f)(s, a) = f(s_{h+1}^n)$

- How to choose the *learning rate*  $\alpha_n$  and the *bonus*  $b_h^n$ ?
- learning rate*  $\alpha_n \approx 1/n$ , unfolding the formula for  $Q_h^n$  + Hoeffding inequality

$$Q_h^n(s, a) \approx r_h(s, a) + \frac{1}{n} \sum_{i=1}^n p_h^i \bar{V}_{h+1}^{i-1}(s, a)$$

$$\approx r_h(s, a) + p_h \left( \frac{1}{n} \sum_{i=1}^n \bar{V}_{h+1}^{i-1} \right) (s, a) \pm \underbrace{\sqrt{\frac{H^2}{n}}}_{\text{variance term} \rightarrow \text{bonus}}$$

$:= \bar{V}_{h,s,a}^n$  bias-value function

- learning rate*  $\alpha_n \approx H/n$  (OptQL [Jin et al., 2018])

$$Q_h^n(s, a) \approx r_h(s, a) + \frac{H}{n} \sum_{i \geq n-H/n}^n p_h^i \bar{V}_{h+1}^{i-1}(s, a)$$

$$\approx r_h(s, a) + p_h \left( \frac{H}{n} \sum_{i \geq n-H/n}^n \bar{V}_{h+1}^{i-1} \right) (s, a) \pm \underbrace{\sqrt{\frac{H^3}{n}}}_{\text{variance term}}$$

$:= \bar{V}_{h,s,a}^n$  bias-value function

## UCB Momentum Q-learning

**Idea** add a (negative) momentum to correct the bias [Azar et al., 2011]

*learning rate*  $\alpha_n \approx 1/n$  and *momentum rate*  $\gamma_n \approx H/n$ : UCBMQ

$$Q_h^n(s, a) = \alpha_n(r_h + p_h^n \bar{V}_{h+1}^{n-1})(s, a) + (1 - \alpha_n)Q_h^{n-1}(s, a) + \underbrace{\gamma_n p_h^n (\bar{V}_{h+1}^{n-1} - V_{h,s,a}^{n-1})(s, a)}_{\leq 0, \text{ momentum}}$$

where the bias-value function

$$V_{h,s,a}^n(s') = (\alpha_n + \gamma_n) \bar{V}_{h+1}^{n-1}(s') + (1 - \alpha_n - \gamma_n) V_{h,s,a}^{n-1}(s')$$

$$\approx \frac{H}{n} \sum_{i \geq n-H/n}^n \bar{V}_{h+1}^{i-1}(s')$$

Unfolding+ Hoeffding inequality

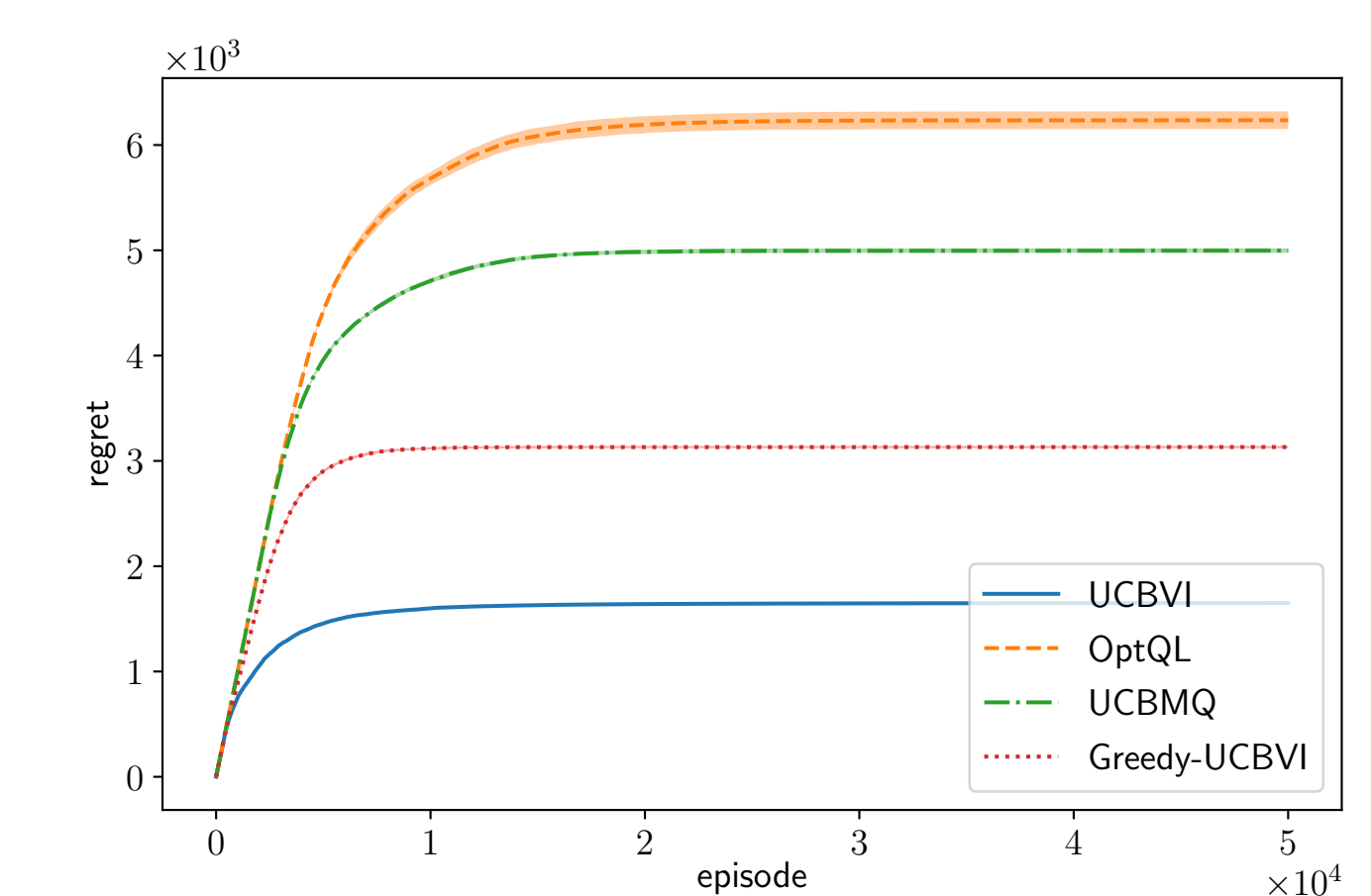
$$Q_h^n(s, a) \approx r_h(s, a) + \frac{1}{n} \sum_{i=1}^n p_h^i ((H+1) \bar{V}_{h+1}^{i-1} - V_{s,a,h}^{i-1})(s, a)$$

$$\approx r_h(s, a) + p_h \left( \frac{H}{n} \sum_{i \geq n-H/n}^n \bar{V}_{h+1}^{i-1} \right) (s, a) \pm \underbrace{\sqrt{\frac{H^2}{n}}}_{\text{variance term}} \pm \underbrace{\sqrt{\frac{H^3}{n} \sum_{i=1}^n p_h(V_{h,s,a}^{n-1} - \bar{V}_h^{n-1})(s, a)) \frac{1}{n}}}_{\text{momentum variance term}}$$

- keep only the last  $H/n$  fraction of the past targets: bound *polynomial in  $H$*

- $n$  samples to approximate the mean

- still an extra  $H$  in the bonus  $\rightarrow$  Bernstein inequality instead of Hoeffding



## Rates

Algorithm	Upper bound
UCBVI [Azar et al., 2017]	$\tilde{O}(\sqrt{H^3 SAT} + H^3 S^2 A)$
UBEV [Dann et al., 2017]	$\tilde{O}(\sqrt{H^4 SAT} + H^2 S^3 A^2)$
EULER [Zanette and Brunskill, 2019]	$\tilde{O}(\sqrt{H^3 SAT} + H^3 S^{3/2} A(\sqrt{S} + \sqrt{H}))$
OptQL [Jin et al., 2018] (Bernstein)	$\tilde{O}(\sqrt{H^4 SAT} + H^{9/2} S^{3/2} A^{3/2})$
UCB-Advantage [Zhang et al., 2020]	$\tilde{O}(\sqrt{H^3 SAT} + H^{33/4} S^2 A^{3/2} T^{1/4})$
UCBMQ (this paper)	$\tilde{O}(\sqrt{H^3 SAT} + H^4 SA)$