

# On the Sample Complexity of $(\varepsilon, \delta)$ -PAC Learning with the Entropic Risk Measure

Amer Essakine under the supervision of Dr. Claire Vernade

ENS Paris Saclay, MVA master

October 22, 2025

# Introduction

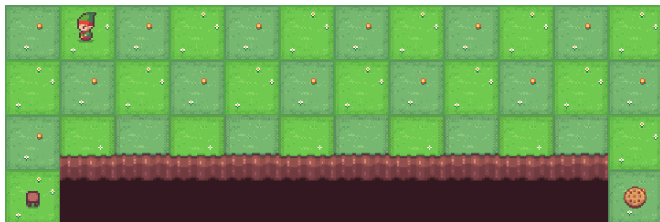


Figure: The cliff environment [6]

# Markov decision process

- Consider an finite episodic MDP  $(\mathcal{S}, \mathcal{A}, H, \{p_h\}_{h \in [H]}, \{r_h\}_{h \in H})$
- Assume the rewards are deterministic, bounded in  $[0,1]$ .
- we want to solve the problem :

$$\pi^* = \arg \max_{\pi \in \Pi_{\text{Markov, det}}} \rho(R^\pi)$$

Where  $\rho$  is functional called risk measure

In practice, the risk measure satisfies certain natural properties from a measure of risk like monotonicity and translation invariance.

- **Threshold probability**

$$\Pr(R^\pi \geq T)$$

- **Value-at-Risk** At level  $\alpha \in (0, 1)$  :

$$\text{VaR}_\alpha[R^\pi] = \inf \left\{ x \in \mathbb{R} : \Pr(R^\pi \leq x) \geq \alpha \right\}.$$

- **Conditional-Value-at-Risk** At level  $\alpha \in (0, 1)$  :

$$\text{CVaR}_\alpha[R^\pi] = \frac{1}{\alpha} \int_0^\alpha \text{VaR}_\gamma[R^\pi] d\gamma$$

# Entropic risk measure

For a random variable  $X$  and parameter  $\beta \in \mathbb{R}$ , the entropic risk measure is defined as

$$\rho_{\beta}(X) = \begin{cases} \frac{1}{\beta} \log \mathbb{E} [e^{\beta X}] & \beta \neq 0 \\ \mathbb{E}[X] & \beta = 0 \end{cases}$$

- Second order expansion of the entropic risk measure

$$\rho_{\beta}(X) = \mathbb{E}[X] + \frac{\beta}{2} \text{Var}(X) + O(\beta^2), \quad \beta \rightarrow 0.$$

- Can be exactly optimized using dynamic programming [4]
- Apart from the expectation, it's the only continuous objective that can be optimized exactly using dynamic programming[10]

# Bellman equations for entropic risk measure

- **Planning in MDPs:** Entropic risk measure admit a *log-sum-exp* Bellman update with policy/value iteration for fixed risk parameter [7]: **Policy evaluation (reward, sums):**

$$V_t^\pi(s) = r_t(s, a) + \frac{1}{\beta} \log \left( \sum_{a \in \mathcal{A}(s)} \pi_t(a | s) \sum_{s' \in \mathcal{S}} P_t(s' | s, a) \exp \left\{ \beta (V_{t+1}^\pi(s')) \right\} \right)$$

**Optimality (reward, sums):**

$$V_t(s) = r_t(s, a) + \max_{a \in A} \frac{1}{\beta} \log \left( \sum_{s' \in \mathcal{S}} P_t(s' | s, a) \exp \left\{ \beta (V_{t+1}(s')) \right\} \right)$$

- **Learning (RL):** Asymptotic convergence for risk-sensitive actor-critic and Q-learning [3, 2]

# Best policy identification

We consider a BPI algorithm with a policy sequence  $\{\pi^t\}_{t \in \mathbb{N}}$ , an exploration budget  $\tau$ , and an output policy  $\hat{\pi}$ .

Our goal is to design an  $(\varepsilon, \delta)$ -PAC algorithm (defined below) that minimizes the sample complexity, i.e., the number of exploration episodes  $\tau$ .

## Definition (PAC algorithm for BPI)

An algorithm is  $(\varepsilon, \delta)$ -PAC for best policy identification if it returns a policy  $\hat{\pi}$  after some number of episodes  $\tau$  that satisfies

$$\mathbb{P} \left( V_1^*(s_1) - V_1^{\hat{\pi}}(s_1) \leq \varepsilon \right) \geq 1 - \delta$$

# BPI problem challenges

The three core challenges are:

1. the exploration rule: How to explore the environment
2. the stopping rule: when to stop so that the output is  $(\varepsilon, \delta)$ -PAC while minimizing the number of exploration episodes  $\tau$ .
3. which policy to output



# Contributions

- Literature review and problem framing
- Lower-bound on the BPI for entropic risk measure
- KL-driven exploration for entropic risk

1. Literature review
2. The Exponential Curse: Lower Bounds for the Entropic Risk Measure
3. On the Sample Complexity of  $(\epsilon, \delta)$ -PAC Learning with the Entropic Risk Measure
4. Conclusion

# Interaction with the Environment

- **Generative model (simulator):** Can simulate the next step from any  $(s, a)$
- **Forward / dynamics model:** Can only interact with the environment. After sampling a state  $s_0$ , we sample a trajectory
- **Reward-free model:** Can only sample trajectories and we do not have access to the reward

# Literature review for risk-neutral case

Perspective	Reference	Lower Bound	Upper Bound
Generative model	[1]	$\Omega\left(\frac{SAH^3}{\epsilon^2}\right)$	$\tilde{O}\left(\frac{SAH^4}{\epsilon^2}\right)$
Forward model	[5]	$\Omega\left(\frac{SAH^3}{\epsilon^2}\right)$	
	[9]		$\tilde{O}\left(\frac{SAH^4}{\epsilon^2}\right)$
	[11]		$\tilde{O}\left(\frac{SAH^3}{\epsilon^2}\right)$
Reward-free	[5]	$\Omega\left(\frac{SAH^3}{\epsilon^2} + S\right)$	
	[8]		$\tilde{O}\left(\frac{S^2AH^7}{\epsilon} + \frac{S^2AH^5}{\epsilon^2}\right)$
	[9]		$\tilde{O}\left(\frac{SAH^4}{\epsilon^2} + S\right)$
	[11]		$\tilde{O}\left(\frac{SAH^3}{\epsilon^2} + S\right)$

## Literature review for risk-sensitive case

Perspective	Reference	Lower Bound	Sampling complexity
Generative Model	[12]	$\Omega\left(\frac{SA\gamma^2}{c_1\varepsilon^2} \frac{e^{ \beta \frac{1}{1-\gamma}} - 3}{ \beta ^2}\right)$	$\tilde{O}\left(\frac{SA\left(S + \log(SA/\delta)\right)}{\varepsilon^2(1-\gamma)^2\beta^2} \cdot e^{2 \beta \frac{1}{1-\gamma}}\right)$

**Gap to lower bounds:** Between the lower bound of [12] and current best upper bounds, there remains at least a gap of  $H^2 (e^{|\beta|H} - 1)$

- **Planning vs. learning:** Planning for entropic-risk MDPs is well understood, but *learning* remains less developed and optimal PAC/BPI guarantees are still open.
- **Tail amplification & variance:** The entropic criterion applies an exponential tilt  $e^{\beta R}$  that magnifies tail outcomes—inflating high rewards when  $\beta > 0$  (penalizing low rewards when  $\beta < 0$ )—which raises estimator variance and induces difficulty that scales exponentially in  $|\beta|$  and the horizon  $H$ .
- **Regret results insufficient:** Even the strongest regret bounds to date do not close this gap; PAC guarantees tailored to entropic risk remain scarce.

1. Literature review
2. The Exponential Curse: Lower Bounds for the Entropic Risk Measure
3. On the Sample Complexity of  $(\epsilon, \delta)$ -PAC Learning with the Entropic Risk Measure
4. Conclusion

# Lower bound from the hard MDP

## Theorem (PAC lower bound for the entropic value objective)

Fix  $\beta \in \mathbb{R}$  and integers  $H \geq 3$ ,  $S \geq 4$ ,  $A \geq 2$ ,  $\delta \in (0, \frac{S-3}{2})$ , and sufficiently small  $\varepsilon > 0$ . There exists an episodic MDP with horizon  $H$ ,  $S$  states, and  $A$  actions such that any  $(\varepsilon, \delta)$ -PAC algorithm must, for some instance of this MDP, use an expected number of episodes  $T$  satisfying

$$T = \Omega\left(SAH^2 \frac{e^{|\beta|H} - 1}{\beta^2 \varepsilon^2} \log\left(\frac{S}{\delta}\right)\right)$$

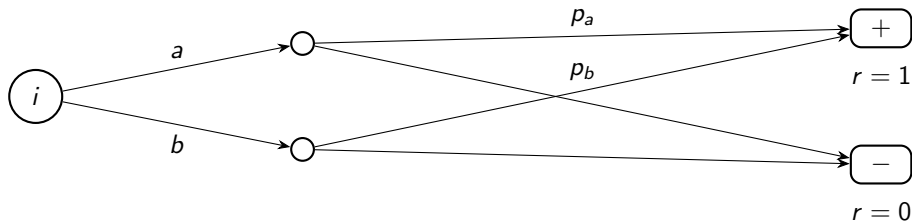


## Theorem

Fix  $\theta_0, \theta_1 \in \Theta$  such that  $\rho(\theta_0, \theta_1) \geq 2\varepsilon$ , then :

$$\inf_{\hat{\theta}} \sup_{i \in \{0,1\}} P_{\theta_i}^{\otimes n}(\rho(\hat{\theta}, \theta_i) \geq \varepsilon) \geq \frac{1 - \|P_{\theta_0}^{\otimes n} - P_{\theta_1}^{\otimes n}\|_{TV}}{2}$$

## A hard bandit for a lower bound



# Bandit Entropic Lower Bound

- **Create two hard instances:** they differ only in which action is better, and the advantage is tuned to be very small (exact target gap)
- **Bound information per episode:** the Kullback–Leibler divergence between the two instances for one episode is small so the two instances are hard to distinguish
- **Link indistinguishability to error** Apply the Bretagnolle–Huber inequality to show that if the two instances are hard to tell apart, any algorithm must make an error on at least one of them.
- **Trade accuracy for samples:** requiring small error on both instances forces a minimum total information budget.

## Create two hard instances

We define two MDPs  $\mathcal{M}^+$  and  $\mathcal{M}^-$  within the family of MDPs given before that only differ in which Left action is better  $a$  or  $b$ .

Denote

$$p = \frac{1}{2c} \text{ and } q = p(1 + \eta) \text{ and } c = e^{\beta R} - 1$$

Where  $\eta$  will be chosen to make the entropic gap between the two MDPs exactly  $\varepsilon$ , more precisely :

$$G(p) - G(q) = \frac{1}{\beta} \log \left( \frac{1 + cp}{1 + cq} \right) = \frac{1}{\beta} \log \left( 1 + \frac{\eta}{3} \right)$$

Hence, we chose

$$\eta = 3(e^{\beta\varepsilon} - 1) \quad \text{So that} \quad G(p) - G(q) \leq \varepsilon$$

In  $\mathcal{M}^+$ ,  $p_a = q$  and  $p_b = p$

## Bound information per episode

Let  $\mathbb{P}_{1:n}^+$  and  $\mathbb{P}_{1:n}^-$  denote the laws of the entire  $n$ -episode simulation under  $\mathcal{M}^+$  and  $\mathcal{M}^-$ . By decomposing the KL divergence we get:

$$\text{KL}(\mathbb{P}_{1:n}^+ \| \mathbb{P}_{1:n}^-) = \sum_{t=1}^n \mathbb{E}[\text{KL}(\text{episode } t \mid \text{history})] \leq n \max\{d(p, q), d(q, p)\} \leq n \frac{\eta^2}{2c - 1}$$

## Link in-distinguishability to error

We apply the Bretagnolle–Huber inequality to the event  $A$  "the algorithm outputs action  $a$ ". On  $\mathcal{M}^+$  the error is  $A^c$ ; on  $\mathcal{M}^-$  the error is  $A$  :

$$\Pr_{\mathcal{M}^+}(A^c) + \Pr_{\mathcal{M}^-}(A) \geq \frac{1}{2} \exp \left( - \text{KL} \left( \mathbb{P}_{1:n}^+ \parallel \mathbb{P}_{1:n}^- \right) \right)$$

If the learner is  $(\varepsilon, \delta)$ -correct on both instances, the LHS  $\leq 2\delta$ . Hence :

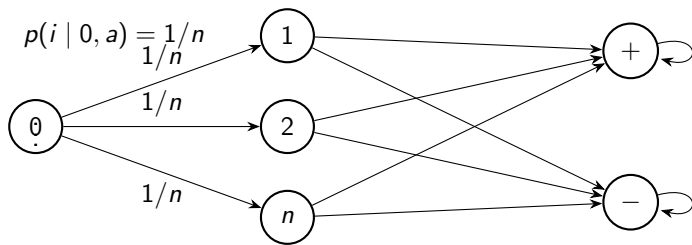
$$2\delta \geq \frac{1}{2} \exp \left( -n \cdot \frac{\eta^2}{2c-1} \right),$$

## Theorem

*Any algorithm that solves the  $(\epsilon, \delta)$ -PAC problem for  $\epsilon$  satisfying  $\beta\epsilon \leq \ln(2)$  must have a sampling complexity that satisfy :*

$$n \geq \frac{2(e^{\beta R} - 1) - 1}{72\beta^2\epsilon^2} \log \frac{1}{4\delta}$$

# A hard MDP



$$\begin{aligned} r(+) &= H - 2 \\ p(+ | i, a) &= p_+ \quad \text{if } a = a_i^* \\ p(+ | i, a) &= p_- \quad \text{if } a \neq a_i^* \end{aligned}$$



# Lower bound from the hard MDP

## Theorem (PAC lower bound for the entropic value objective)

Fix  $\beta \in \mathbb{R}$  and integers  $H \geq 3$ ,  $S \geq 4$ ,  $A \geq 2$ ,  $\delta \in (0, \frac{S-3}{2})$ , and sufficiently small  $\varepsilon > 0$ . There exists an episodic MDP with horizon  $H$ ,  $S$  states, and  $A$  actions such that any  $(\varepsilon, \delta)$ -PAC algorithm must, for some instance of this MDP, use an expected number of episodes  $T$  satisfying

$$T = \Omega\left(SAH^2 \frac{e^{|\beta|H} - 1}{\beta^2 \varepsilon^2} \log\left(\frac{S}{\delta}\right)\right)$$

1. Literature review
2. The Exponential Curse: Lower Bounds for the Entropic Risk Measure
3. On the Sample Complexity of  $(\varepsilon, \delta)$ -PAC Learning with the Entropic Risk Measure
4. Conclusion

# Empirical MDP

Let  $(s_h^i, a_h^i, s_{h+1}^i)$  be the state, the action, and the next state observed by an algorithm at step  $h$  of episode  $i$ .

For any step  $h \in [H]$  and any state–action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , define

$$n_h^t(s, a) \triangleq \sum_{i=1}^t \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \quad n_h^t(s, a, s') \triangleq \sum_{i=1}^t \mathbf{1}\{(s_h^i, a_h^i, s_{h+1}^i) = (s, a, s')\}.$$

These definitions permit us to define the empirical transitions:

$$\hat{p}_h^t(s' \mid s, a) \triangleq \begin{cases} \frac{n_h^t(s, a, s')}{n_h^t(s, a)}, & \text{if } n_h^t(s, a) > 0, \\ \frac{1}{S}, & \text{otherwise.} \end{cases}$$

# UCB algorithm

Sampling rule. At iteration  $t$ , act greedily w.r.t.  $\tilde{Q}^t$ :

$$\forall s \in \mathcal{S}, \forall h \in [H], \quad \pi_h^{t+1}(s) = \arg \max_{a \in \mathcal{A}} \tilde{Q}_h^t(s, a).$$

Stopping rule. Stop at the first  $t$  such that the certificate is below  $\varepsilon$ :

$$\tau = \inf \left\{ t \in \mathbb{N} : \pi_1^{t+1} G_1^t(s_1) \leq \varepsilon \right\}.$$

Prediction rule. Output the policy from the next iteration:

$$\hat{\pi} = \pi^{\tau+1}.$$

Given a confidence region set  $\mathcal{U}_h^t(s, a)$  such that for any  $(s, a)$ :

$$p_h(\cdot|s, a) \in \mathcal{U}_h^t(s, a)$$

An optimistic upper bound of the value function is then:

$$Q_h(s, a) = r_h(s, a) + \sup_{p \in \mathcal{U}_h^t(s, a)} \rho_p(V_h^t(s, a))$$

$$V_h^t(s, a) = \max_{a \in A} Q_h^t(s, a)$$

We chose the KL-divergence confidence region:

$$\mathcal{U}_h^t(s, a) = \left\{ q \in \Sigma_S \mid D_{KL}(\hat{p}_h^t(s, a) \parallel q(s, a)) \leq \frac{\lambda(n_h^t(s, a), \delta)}{n_h^t(s, a)} \right\}$$

# Why KL planning ?

1. For  $\lambda(n, \delta) = \log\left(\frac{3SAH}{\delta}\right) + S \log(8e(n+1))$  with probability  $1 - \delta$ [11]:

$$p_h(s, a) \in \mathcal{U}_h^t(s, a)$$

2. Controls trajectory mismatch. The chain rule of KL yields :

$$\text{KL}(\hat{P}^{t,\pi}, P^\pi) = \sum_{h,s,a} d_h^{t,\pi}(s, a) \text{KL}(\hat{p}_h^t(\cdot|s, a), p_h^*(\cdot|s, a))$$

so shrinking local KL terms reduces global error.

# A Bernstein's inequality

To turn the optimistic upper bound to a computable bound we use Bernstein's inequality [11]:

## Theorem

For any  $q \in \mathcal{U}_h^t(s, a)$ , for any bounded  $f : \mathcal{S} \rightarrow [0, b]$ :

$$|E_p[f] - E_q[f]| \leq \sqrt{2 \text{Var}_q(f) \frac{\lambda(n_h^t(s, a), \delta)}{n_h^t(s, a)}} + \frac{2}{3} b \frac{\lambda(n_h^t(s, a), \delta)}{n_h^t(s, a)}$$

## Bonus terms

Using Bernstein's inequality, we define bonus terms :

$$\begin{aligned}\tilde{B}_h^t(s, a) &\triangleq 2\sqrt{2} \sqrt{\text{Var}_{\hat{p}_h^t}(e^{\beta \tilde{V}_{h+1}^t})(s, a) \frac{\beta^*(n_h^t(s, a), \delta)}{n_h^t(s, a)}} + e^{\beta H} (8H + 2\sqrt{2} + 3) \frac{\beta(n_h^t(s, a), \delta)}{n_h^t(s, a)} \\ \tilde{s}_h^t(s, a) &\triangleq \frac{1}{\max\{1, \hat{p}_h^t V_{h+1}^t(s, a) - \tilde{B}_h^t(s, a)\}}\end{aligned}$$



# Optimistic/pessimistic state-value function

Upper and lower Q- and value functions:

$$\tilde{Q}_h^t(s, a) \triangleq \min \left\{ H, r_h(s, a) + \frac{1}{|\beta|} \log(1 + \tilde{s}_h^t(s, a) \tilde{B}_h^t(s, a)) \right. \\ \left. + \frac{1}{\beta H} \left( \rho_{\beta}^{\hat{p}_h^t}(\tilde{V}_{h+1}^t(s, a)) - \rho_{\beta}^{\hat{p}_h^t}(\underline{V}_{h+1}^t(s, a)) + \rho_{\beta}^{\hat{p}_h^t}(\tilde{V}_{h+1}^t(s, a)) \right) \right\},$$

$$\underline{Q}_h^t(s, a) \triangleq \max \left\{ 0, r_h(s, a) - \frac{1}{|\beta|} \log(1 + \tilde{s}_h^t(s, a) \tilde{B}_h^t(s, a)) \right. \\ \left. - \left( \rho_{\beta}^{\hat{p}_h^t}(\tilde{V}_{h+1}^t(s, a)) - \rho_{\beta}^{\hat{p}_h^t}(\underline{V}_{h+1}^t(s, a)) + \rho_{\beta}^{\hat{p}_h^t}(\underline{V}_{h+1}^t(s, a)) \right) \right\},$$

$$\tilde{V}_h^t(s) \triangleq \max_{a \in \mathcal{A}} \tilde{Q}_h^t(s, a), \quad \underline{V}_h^t(s) \triangleq \max_{a \in \mathcal{A}} \underline{Q}_h^t(s, a), \quad \tilde{V}_{H+1}^t = \underline{V}_{H+1}^t \equiv 0.$$

# Stopping rule

We define the stopping rule bounding the gap between the optimal policy and the policy at instant  $t$ :

$$G_h^t(s) = \min \left\{ H, \frac{2}{\beta} \log(1 + \tilde{s}_h^t \tilde{B}_h^t(s, \pi_h^{t+1}(s))) + \left(1 + \frac{3}{H}\right) \hat{p}_h^t \pi_{h+1}^{t+1} G_{h+1}^t(s, a) \right\} \quad G_{H+1}^t \equiv 0$$

# Optimism results

## Lemma (Optimism and pessimism, entropic case)

On  $\mathcal{G}$ , for all  $t \geq 0$ ,  $h \in [H]$ , and  $(s, a)$ ,

$$\underline{Q}_h^t(s, a) \leq Q_h^*(s, a) \leq \tilde{Q}_h^t(s, a) \quad \underline{V}_h^t(s) \leq V_h^*(s) \leq \tilde{V}_h^t(s)$$

And the stopping rule :

## Lemma

*The greedy policy satisfies the standard gap domination:*

$$V_1^*(s_1) - V_1^{\pi^{t+1}}(s_1) \leq \pi_1^{t+1} G_1^t(s_1) \quad \forall t \geq 0$$

# UCB algorithm for entropic risk measure

---

**Algorithm 1** Entropic-BPI (greedy w.r.t. upper entropic confidence)

---

- 1: **Input:**  $\beta, \delta \in (0, 1), \varepsilon > 0$ .
- 2: Initialize counts  $n_h^0(\cdot) = 0$  and  $\hat{p}_h^0(\cdot|s, a) = 1/S$ .
- 3: **for**  $t = 0, 1, 2, \dots$  **do**
- 4:   For  $h = 1, \dots, H$  Compute  $\tilde{B}_h^t$  and  $\tilde{s}_h^t$ ; then  $(\tilde{Q}_h^t, \tilde{V}_h^t)$  and  $(\underline{Q}_h^t, \underline{V}_h^t)$ .
- 5:   **Sampling rule:**  $\pi_h^{t+1}(s) \in \arg \max_{a \in \mathcal{A}} \tilde{Q}_h^t(s, a)$  for all  $h, s$ .
- 6:   **Stopping potential:** for  $a = \pi_h^{t+1}(s)$ , set

$$G_h^t(s) = \min \left\{ H, \frac{2}{\beta} \log(1 + \tilde{s}_h^t \tilde{B}_h^t)(s, \pi_h^{t+1}(s)) + \left(1 + \frac{3}{H}\right) \hat{p}_h^t \pi_{h+1}^{t+1} G_{h+1}^t(s, a) \right\}$$

- 7:   **Stopping rule:**  $\tau = \inf \{t \in \mathbb{N} : \pi_1^{t+1} G_1^t(s_1) \leq \varepsilon\}$ . If  $t \geq \tau$ , **stop** and output  $\pi^{t+1}$ .
- 8:   Execute episode  $t + 1$  with  $\pi^{t+1}$ , update counts and  $\hat{p}_h^{t+1}$ .
- 9: **end for**

## Theorem

*For  $\varepsilon \in ]0, 1]$  and  $\delta > 0$  the algorithm Entropic-KL-BPI return an  $\varepsilon$ -optimal policy with probability at least  $1 - \delta$  after  $\tau$  steps, where with probability of  $1 - \delta$  and we have an upper bound on  $\tau$ :*

$$\tau = \tilde{O} \left( \frac{(e^{|\beta|H} - 1)^2}{\beta^2} \cdot \frac{H^2 SA}{\varepsilon^2} \log \left( \frac{3SAH}{\delta} \right) \right)$$

1. Literature review
2. The Exponential Curse: Lower Bounds for the Entropic Risk Measure
3. On the Sample Complexity of  $(\epsilon, \delta)$ -PAC Learning with the Entropic Risk Measure
4. Conclusion

# A promising approach

Planning using Renyi divergence:

$$Q_h^t(s, a) = r_h(s, a) + \sup_{p \in \mathcal{U}_{h,\kappa}^t(s, a)} \rho_\beta^p(V_{h+1}^t) \quad V_h^t(s) = \max_{a \in \mathcal{A}} Q_h^t(s, a), \quad V_{H+1}^t \equiv 0$$

$$\mathcal{U}_{h,\kappa}^t(s, a) := \left\{ p : R_\kappa(p \| \hat{p}_h^t(\cdot | s, a)) \leq r_h^t(s, a) \right\}$$

We get by using the variational formula

$$\sup_{p: R_\kappa(p \| \hat{p}) \leq r} \rho_\beta^p(X) \leq \rho_\gamma^{\hat{p}}(X) + \frac{r}{\gamma - \beta} \quad \gamma > \beta, \quad \kappa = \frac{\gamma}{\gamma - \beta}$$

This gives a computable UCB bound:

$$Q_h^t(s, a) \leq r_h(s, a) + \rho_\gamma^{\hat{p}_h^t(\cdot | s, a)}(V_{h+1}^t) + \frac{r_h^t(s, a)}{\gamma - \beta} \quad \gamma > \beta \quad \kappa = \frac{\gamma}{\gamma - \beta}$$

# Renyi divergence

An optimal upper bound to the state-value function is:

$$Q_h^t(s, a) = r_h(s, a) + \sup_{p \in \mathcal{U}_{h,\kappa}^t(s, a)} \rho_\beta^p(V_{h+1}^t) \quad V_h^t(s) = \max_{a \in \mathcal{A}} Q_h^t(s, a), \quad V_{H+1}^t \equiv 0$$

$$\mathcal{U}_{h,\kappa}^t(s, a) := \left\{ p : R_\kappa(p \| \hat{p}_h^t(\cdot | s, a)) \leq r_h^t(s, a) \right\}$$

And we get by using the variational formula:

$$\sup_{p \in \mathcal{U}_{h,\kappa}^t(s, a)} \rho_\beta^p(V_{h+1}^t) \leq \rho_\gamma^{\hat{p}}(X) + \frac{r}{\gamma - \beta} \quad \gamma > \beta, \quad \kappa = \frac{\gamma}{\gamma - \beta}$$

This gives a computable UCB bound whenever  $p_h(s, a) \in \mathcal{U}_{h,k}^t(s, a)$ :

$$Q_h^t(s, a) \leq r_h(s, a) + \rho_\gamma^{\hat{p}_h^t(\cdot | s, a)}(V_{h+1}^t) + \frac{r_h^t(s, a)}{\gamma - \beta} \quad \gamma > \beta \quad \kappa = \frac{\gamma}{\gamma - \beta}$$



# Problem

We have to control the Renyi divergence

# Renyi divergence is hard to control

## Theorem

*Let  $\alpha > 1$  and let  $\hat{p}_n$  be the empirical distribution of  $n$  i.i.d. samples from an unknown categorical distribution  $p$  on a finite alphabet. There exist constants  $C_\alpha > 0$  and  $\delta_0 \in (0, 1)$  such that for every  $\delta \in (0, \delta_0]$  there exists a (binary) distribution  $p$  with*

$$\mathbb{P}_p \left( \exists n \geq 1 : nD_\alpha(\hat{p}_n \| p) \geq C_\alpha \delta^{-(\alpha-1)} \right) \geq \delta.$$

*Consequently, any distribution-free, time-uniform tail inequality of the form*

$$\sup_p \mathbb{P}_p \left( \exists n \geq 1 : nD_\alpha(\hat{p}_n \| p) \geq b(n, \delta) \right) \leq \delta$$

*must satisfy, for each  $\delta \in (0, \delta_0]$ , that  $b(n, \delta) \geq C_\alpha \delta^{-(\alpha-1)}$  for some  $n$ . In particular, no schedule with  $b(n, \delta) = O(\log(1/\delta))$  works*

# Possible solutions (a burn-in phase)

## Theorem (All-time Rényi bound after Chernoff burn-in)

Let  $p$  be a probability distribution on  $S$ , let  $\hat{p}_n$  be the empirical estimation, and suppose  $b = \min_{i, p_i > 0} p_i$ . Fix  $\delta \in (0, 1)$  and  $\alpha > 1$ . Set the likelihood-ratio cap target to  $L_0 = 1 + \eta$  with  $\eta = \frac{1}{2}$  (so  $L_0 = \frac{3}{2}$ ), and define

$$n_0 = \left\lceil \frac{12}{b} \left( \log \frac{4S}{\delta} + \log \frac{24}{b} \right) \right\rceil$$

For  $\varepsilon \in (0, 1)$  define

$$\beta(n, \delta) = \log \frac{1}{\delta} + (S - 1) \log \left( e \left( 1 + \frac{n}{S-1} \right) \right)$$

Then, with probability at least  $1 - \delta$ , simultaneously for all  $n \geq n_0$ ,

$$D_\alpha(\hat{p}_n \| p) \leq \begin{cases} \frac{3}{n} \beta(n, \frac{\delta}{2}) & 1 < \alpha \leq 2 \\ \frac{\alpha(3/2)^{\alpha-1}}{n} \beta(n, \frac{\delta}{2}) & \alpha > 2 \end{cases}$$

# References

- [1] Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J. Kappen. On the sample complexity of reinforcement learning with a generative model. In *Proceedings of the 29th International Conference on Machine Learning (ICML)*, pages 1707–1714, 2012.
- [2] Vivek S. Borkar. Q-learning for risk-sensitive control. *Mathematics of Operations Research*, 26(2): 294–311, 2001. doi: 10.1287/moor.26.2.294.10559.
- [3] Vivek S. Borkar. A sensitivity formula for risk-sensitive cost and the actor–critic algorithm. *Systems & Control Letters*, 44(5):339–346, 2001. doi: 10.1016/S0167-6911(01)00105-4.
- [4] Vivek S. Borkar. Q-learning for risk-sensitive control. *Mathematics of Operations Research*, 27(2): 294–311, 2002. doi: 10.1287/moor.27.2.294.3238.
- [5] Christoph Dann and Emma Brunskill. Sample complexity of episodic fixed-horizon reinforcement learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2818–2826, 2015.
- [6] Farama Foundation. Cliff walking — gymnasium. URL [https://gymnasium.farama.org/environments/toy\\_text/cliff\\_walking/](https://gymnasium.farama.org/environments/toy_text/cliff_walking/).
- [7] Ronald A. Howard and James E. Matheson. Risk-sensitive markov decision processes. *Management Science*, 18(7):356–369, 1972. doi: 10.1287/mnsc.18.7.356.
- [8] Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119 of *Proceedings of Machine Learning Research*, pages 4870–4879. PMLR, 2020. URL <https://proceedings.mlr.press/v119/jin20d.html>.
- [9] Emilie Kaufmann, Pierre Ménard, Omar Darwiche Domingues, Anders Jonsson, Edouard Leurent, and