

---

# Supplementary material for the paper: ”Adaptive Bandits: Towards the best history-dependent strategy“

---

**Odalric-Ambrym Maillard**  
INRIA Lille Nord-Europe

**Rémi Munos**  
INRIA Lille Nord-Europe

## Abstract

In this document, we detail further some technical proofs not covered in the paper corresponding to this supplementary material.

## 1 Playing against an opponent using a known model

### 1.1 Regret upper bounds against the best history-class-based strategy

**Theorem 1** *In the case of a  $\Phi$ -constrained opponent, using the  $\Phi$ -UCB algorithm with parameter  $\alpha > 1/2$ , we have the distribution-dependent bound:*

$$R_T^\Phi \leq \sum_{c \in \mathcal{H}/\Phi; \mathbb{E}(I_c(T)) > 0} \sum_{a \in \mathcal{A}; \Delta_c(a) > 0} \frac{4\alpha \log(T)}{\Delta_c(a)} + \Delta_c(a)c_\alpha$$

where  $I_c(T) = \sum_{t=1}^T \mathbb{I}_{[h_{<t}] = c}$ , the per-class gaps  $\Delta_c(a) \stackrel{\text{def}}{=} \mu_c(a^*) - \mu_c(a)$ , and the constant  $c_\alpha = 1 + \frac{4}{\log(\alpha+1/2)} \left(\frac{\alpha+1/2}{\alpha-1/2}\right)^2$ . We also have a distribution-free bound (i.e. which does not depend on the gaps):

$$R_T^\Phi \leq \sqrt{TAC\overline{C}(4\alpha \log(T) + c_\alpha)}$$

where  $\overline{C} = |\{c \in \mathcal{H}/\Phi; \mathbb{E}(I_c(T)) > 0\}|$  is the number of classes that may be activated during the run.

Now, in the case of an arbitrary opponent, using  $\Phi$ -Exp3 algorithm, we have:

$$\tilde{R}_T^\Phi \leq \frac{3}{\sqrt{2}} \sqrt{T\overline{C}A \log(A)}.$$

*Proof:*  **$\Phi$ -UCB:** The distribution-dependent bound for  $\Phi$ -UCB is a direct application of the result of [2]

---

Appearing in Proceedings of the 14<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2011, Fort Lauderdale, FL, USA. Volume 15 of JMLR: W&CP 15. Copyright 2011 by the authors.

for the algorithm UCB about  $\tau_a(t) \stackrel{\text{def}}{=} \sum_{s=1}^t \mathbb{I}_{a_s = a}$  where  $a_t$  is played by UCB, that states that  $\mathbb{E}(\tau_a(t)) \leq \frac{4\alpha \log(t)}{\Delta_c^2(a)} + c_\alpha$ . Indeed, we use the fact that  $R_T^\Phi = \sum_{c \in \mathcal{H}/\Phi} R_T(c)$  and thus remark that when a class  $c$  is visited, then we play according to a UCB algorithm for this class.

Thus, for the distribution-free bound, we have:

$$\begin{aligned} R_T^\Phi &= \sum_c \sum_a \Delta_c(a) \mathbb{E}(\tau_a(I_c(T))) \\ &\leq \sum_c \sum_a \sqrt{\mathbb{E}(\tau_a(I_c(T)))} \sqrt{4\alpha \log(T) + c_\alpha} \\ &\leq \sum_c \sqrt{\mathbb{E}(I_c(T))} \sqrt{A} \sqrt{4\alpha \log(T) + c_\alpha} \\ &\leq \sqrt{T\overline{C}A} \sqrt{4\alpha \log(T) + c_\alpha}, \end{aligned}$$

where we used that  $\sum_a \tau_a(s) = s$  for all  $s$ , and  $\sum_c I_c(T) = T$ , and the Cauchy-Schwartz inequality twice.

**$\Phi$ -Exp3:** The bound for  $\Phi$ -Exp3 follows from the bound of the anytime version of the Exp3 algorithm. Indeed we have

$$\tilde{R}_T^\Phi \leq \sum_c \mathbb{E} \left( \frac{A}{2} \sum_{i=1}^{I_c(T)} \eta_i^c + \frac{\log(A)}{\eta_{I_c(T)}^c} \right);$$

we deduce the bound by setting  $\eta_i^c = \sqrt{\frac{2 \log A}{Ai}}$ .  $\square$

### 1.2 Lower bounds on the regret

**Theorem 2** *Let  $\sup$  represents the supremum taken over all  $\Phi$ -constrained opponents and  $\inf$  the infimum over all forecasters, then the stochastic  $\Phi$ -regret is lower-bounded as:*

$$\sup_{\Phi; |\mathcal{H}/\Phi| = C} \inf_{\text{algo}} \sup_{\Phi\text{-opp}} R_T^\Phi \geq \frac{1}{20} \sqrt{TAC}.$$

Now, let  $\sup$  represents the supremum taken over all possible opponents and  $\inf$  the infimum over all forecasters, then the adversarial  $\Phi$ -regret is lower-bounded as:

$$\sup_{\Phi; |\mathcal{H}/\Phi|=C} \inf \sup_{\text{algo opp}} \tilde{R}_T^\Phi \geq \frac{1}{20} \sqrt{TAC}.$$

*Proof:* Let us fix the horizon  $T$  and the number of classes  $C$ . We consider the opponent defined using the specific class-function  $\Phi$  such that each class  $c$  is periodically visited every  $C$  time steps, thus  $T/C$  times. Note that  $T = \frac{T}{C}C$  and that this is intuitively the opponent that makes the algorithm switch between classes the most.

Now we define more precisely the rewards output by the opponent. Let us consider the stochastic bandits such that for each class  $c$ , one arm  $a_c$  is a Bernoulli  $B((1 + \epsilon_c)/2)$ , and all others are  $B((1 - \epsilon_c)/2)$ .

Then by application of Lemma 2.2 in [2], for  $\epsilon_c$  of order  $\sqrt{\frac{A}{s}}$ , we have in the Bandit information setting the following inequality:

$$\sup_{a_c} \sum_{t=1}^s \mu_c(a_c) - \mu_c(a_t) \geq s\epsilon_c \left(1 - \frac{1}{A} - \sqrt{\frac{s\epsilon_c}{2A} \log\left(\frac{1 + \epsilon_c}{1 - \epsilon_c}\right)}\right).$$

Thus with the notations  $I_c(T) = \sum_{t=1}^T \mathbb{1}_{c=[h_{<t}]}$  and  $t_c(i) = \min\{t; I_c(t) \geq i\}$ , we deduce that:

$$\begin{aligned} & \sup_{(a_c)_c} \sum_c \sum_{t=1}^T (\mu_c(a_c) - \mu_c(a_t)) \mathbb{1}_{c=[h_{<t}]} = \\ & \sum_c \sup_{a_c} \sum_{i=1}^{I_c(T)} (\mu_c(a_c) - \mu_c(a_{t_c(i)})) \geq \\ & \sum_c I_c(T) \epsilon_c \left(1 - \frac{1}{A} - \sqrt{\epsilon_c \log\left(\frac{1 + \epsilon_c}{1 - \epsilon_c}\right)} \sqrt{\frac{I_c(T)}{2A}}\right). \end{aligned}$$

Since the  $a_c$  are chosen by the opponent such that each class is visited exactly  $I_c(T) = T/C$  times, then we deduce that the  $\Phi$ -pseudo-regret is lower-bounded as:

$$\begin{aligned} & \sup_{(a_c)_c} \sum_c \sum_{t=1}^T (\mu_c(a_c) - \mu_c(a_t)) \mathbb{1}_{c=[h_{<t}]} \\ & \geq \sum_c \frac{T}{C} \epsilon_c \left(1 - \frac{1}{A} - \sqrt{\epsilon_c \log\left(\frac{1 + \epsilon_c}{1 - \epsilon_c}\right)} \sqrt{\frac{T}{2AC}}\right). \end{aligned}$$

Thus, after some tedious computations to optimize  $\epsilon_c$ , we finally get a lower bound of order:  $\frac{1}{20} \sum_c \sqrt{\frac{T}{C} A} = \frac{1}{20} \sqrt{TAC}$ . Note that this is valid only if  $\epsilon_c \sim \sqrt{\frac{A}{I_c(T)}}$  is small (less than 1), i.e. if the number of classes  $C$  is smaller than a constant times  $\frac{T}{A}$  (and if this not the case, the lower bound becomes obviously of order  $T$ ).

The second part of the Theorem can be proved using the same construction.  $\square$

## 2 Playing against an opponent using a pool of models

We first remind the result of [1] relating the cumulative reward of the Exp4 algorithm to the one of the best expert on top of which it is run. We have:

**Lemma 1** *For any  $\gamma \in (0, 1]$ , for any family of experts which includes the uniform expert, one has*

$$\begin{aligned} & \max_{\theta} \sum_{t=1}^T \mathbb{E}_{a \sim \xi_t^\theta} (r_t(a)) - \mathbb{E}_{a_1, \dots, a_T} \left( \sum_{t=1}^T r_t(a_t) \right) \\ & \leq (e - 1)\gamma T + \frac{A \log(|\Theta|)}{\gamma}. \end{aligned}$$

In our case, since the  $\xi_t^\theta$  are not fixed in advance but are random variables, we can not apply this Lemma directly and need to adapt it. Based on the proof of [1], we can prove the following bound:

**Lemma 2** *For any  $\gamma \in (0, 1]$ , for any family of experts which includes the uniform expert such that all expert advices are adapted to the filtration of the past, one has*

$$\begin{aligned} & \max_{\theta} \sum_{t=1}^T \mathbb{E}_{a_1, \dots, a_{t-1}} (\mathbb{E}_{a \sim \xi_t^\theta} (r_t(a))) - \mathbb{E}_{a_1, \dots, a_T} \left( \sum_{t=1}^T r_t(a_t) \right) \\ & \leq (e - 1)\gamma T + \frac{A \log(|\Theta|)}{\gamma}. \end{aligned}$$

*Proof:* Indeed, by construction of the algorithm, the beginning of the original proof from [1] applies and gives

$$\begin{aligned} & \sum_{t=1}^T r_t(a_t) \geq (1 - \gamma) \sum_{t=1}^T \mathbb{E}_{a \sim \xi_t^\theta} (\tilde{r}_t(a)) - \frac{A \log(|\Theta|)}{\gamma} \\ & \quad - (e - 2) \frac{\gamma}{A} \sum_{t=1}^T \sum_{a \in \mathcal{A}} \tilde{r}_t(a), \end{aligned}$$

where  $\tilde{r}_t(a) = \frac{r_t(a_t)}{q_t(a)} \mathbb{1}_{a_t=a}$ .

Now, we use the fact that  $\xi_t^\theta(a)$  is adapted to the filtration of the past (which we denote  $\mathcal{F}^{t-1}$ ) and the property that  $\mathbb{E}(\tilde{r}_t(a) | \mathcal{F}^{t-1}) = \mathbb{E}(r_t(a) | \mathcal{F}^{t-1})$  to deduce that

$$\begin{aligned} & \mathbb{E}(\mathbb{E}_{a \sim \xi_t^\theta} (\tilde{r}_t(a))) = \mathbb{E} \left( \sum_{a \in \mathcal{A}} \mathbb{E}(\tilde{r}_t(a) \xi_t^\theta(a) | \mathcal{F}^{t-1}) \right) \\ & = \mathbb{E} \left( \sum_{a \in \mathcal{A}} \mathbb{E}(\tilde{r}_t(a) | \mathcal{F}^{t-1}) \xi_t^\theta(a) \right) \\ & = \mathbb{E} \left( \sum_{a \in \mathcal{A}} \mathbb{E}(r_t(a) | \mathcal{F}^{t-1}) \xi_t^\theta(a) \right) \\ & = \mathbb{E} \left( \sum_{a \in \mathcal{A}} \mathbb{E}(r_t(a) \xi_t^\theta(a) | \mathcal{F}^{t-1}) \right) \\ & = \mathbb{E}(\mathbb{E}_{a \sim \xi_t^\theta} (r_t(a))) \end{aligned}$$

On the other hand, since by assumption the uniform expert belongs to the set of considered experts, we also have

$$\begin{aligned} \frac{1}{A} \mathbb{E} \left( \sum_{t=1}^T \sum_{a \in \mathcal{A}} \tilde{r}_t(a) \right) &= \sum_{t=1}^T \mathbb{E} (\mathbb{E}_{a \sim U(\mathcal{A})} (r_t(a))) \\ &\leq \max_{\theta} \sum_{t=1}^T \mathbb{E} (\mathbb{E}_{a \sim \xi_t^\theta} (r_t(a))), \end{aligned}$$

where  $U(\mathcal{A})$  denotes the uniform distribution over the set of actions  $\mathcal{A}$ . This concludes the proof.  $\square$

We now define the best model of the pool  $\theta^*$  to be

$$\theta^* = \operatorname{argmax}_{\theta \in \Theta} \sup_{g: \mathcal{H}/\Phi \rightarrow \mathcal{A}} \mathbb{E} \left( \sum_{t=1}^T \left[ r_t(g([h_{<t}])) - r_t(a_t) \right] \right),$$

We then define for any class  $c \in \mathcal{H}/\Phi_{\theta^*}$ , the action  $a_c^* \stackrel{\text{def}}{=} \operatorname{argmax}_a \mu_c(a)$  that corresponds to the best history-class-based strategy. Thus we can also write  $a_t^* = a_{[h_{<t}]_{\theta^*}}^*$ .

## 2.1 The rebel bandit setting

We now introduce the setting of Rebel bandits that is interesting by itself. It will be used in order to compute the model-based regret of the Exp4 algorithm. In this setting, we consider that at time  $t$  the player  $\theta$  proposes a distribution of probability  $\xi_t^\theta$  over the arms, but he/she actually receives the reward corresponding to an action drawn with another distribution, say  $q_t$ . This will be the distribution of probability proposed by the meta algorithm. We now analyze the ( $\Phi$ -constrained) Exp3 and UCB algorithms in this setting and bound the corresponding rebel-regret define by:

**Definition 1** (*Rebel regret*) *The Rebel-regret of the algorithm that proposes at time  $t$  the distribution  $\xi_t^\theta$  but in the game where the action  $a_t \sim q_t$  is played instead is:*

$$\mathcal{R}_T^q(\theta) = \sum_{t=1}^T \mathbb{E}_{a_1, \dots, a_{t-1}} \left( r_t(a_{[h_{<t}]_{\theta^*}}^*) - \mathbb{E}_{a \sim \xi_t^\theta} (r_t(a)) \right).$$

## 2.2 $\Phi$ -Exp3 in the Rebel bandit setting

We now consider using Exp4 on top of  $\Phi$ -constrained algorithms. We first use the experts  $\Phi_\theta$ -Exp3 for  $\theta \in \Theta$  with a slight modification on the definition of the function  $\tilde{l}_t^c(a)$ . Indeed since the action  $a_t$  are driven according to the meta algorithm and not  $\Phi_\theta$ -Exp3, we redefine  $\tilde{l}_t^c(a) = \frac{1-r_t(a)}{q_t(a)} \mathbb{I}_{a_t=a} \mathbb{I}_{[h_{<t}]_{\theta}=c}$  so as to get unbiased estimate of  $r_t(a)$  for all  $a$ . We now provide a bound on the Rebel-regret of the  $\Phi^*$ -Exp3 algorithm.

**Theorem 3** *The  $\Phi_{\theta^*}$ -Exp3 algorithm in the Rebel bandit setting where  $q_t(a) \geq \delta$  for all  $a$ , and choosing the parameter  $\eta_{t_c^\theta}^\theta = \sqrt{\frac{\delta \log(A)}{i}}$  satisfies*

$$\mathcal{R}_T^q(\theta^*) \leq 2 \sqrt{\frac{T \bar{C} \log A}{\delta}}.$$

The proof of this Theorem is reported in the main paper. We now combine Lemma 2 and Theorem 3 to get the final bound:

**Theorem 4** *For any opponent, the adversarial  $\Phi_\Theta$ -regret of Exp4/Exp3 is bounded as*

$$\tilde{R}_T^\Theta = O(T^{2/3} (A \bar{C} \log(A))^{1/3} \log(|\Theta|)^{1/2}),$$

where  $\bar{C} = \max_{\theta \in \Theta} |\mathcal{H}/\Phi_\theta|$  is the maximum number of classes for models  $\theta \in \Theta$ .

*Proof:* Indeed we can apply Theorem 3 using Exp4 meta algorithm with  $\delta = \frac{\gamma}{A}$ . Thus we get:

$$\begin{aligned} \tilde{R}_T^\Theta &= \sum_{t=1}^T \mathbb{E}_{a_1, \dots, a_{t-1}} (r_t(a_{[h_{<t}]_{\theta^*}}^*) - \mathbb{E}_{a_t \sim q_t} r_t(a_t)) \\ &\leq \mathcal{R}_T^q(\theta^*) + (e-1)\gamma T + \frac{A \log(|\Theta|)}{\gamma} \\ &\leq 2 \sqrt{\frac{T A \bar{C} \log A}{\gamma}} + 2\gamma T + \frac{A \log(|\Theta|)}{\gamma}. \end{aligned}$$

We thus choose  $\gamma = \frac{(A \bar{C} \log(A))^{1/3} \log(|\Theta|)^{1/2}}{(4T)^{1/3}}$  to conclude.  $\square$

## 2.3 $\Phi$ -UCB in the Rebel-bandit setting

Similarly, a bound on the Rebel-regret of the  $\Phi^*$ -UCB algorithm can be derived assuming that we consider a  $\Phi^*$ -constrained opponent with  $\Phi^* = \Phi^{\theta^*} \in \Phi_\Theta$ .

**Theorem 5** *The  $\Phi_{\theta^*}$ -UCB algorithm in the Rebel bandit setting where  $q_t(a) \geq \delta$  for all  $a$ , and provided  $\alpha > 1/2$ , satisfies*

$$\mathcal{R}_T^q(\theta^*) \leq \sum_{c \in \mathcal{H}/\Phi^* a \neq a_c^*} \Delta_c(a) \left[ \frac{2\alpha \log(T)}{\Delta_c(a)^2 \delta} + \sqrt{\frac{\pi \delta \Delta_c(a)^2}{32\alpha \log T}} + c_\alpha \right]$$

We also have the distribution-free bound:

$$\mathcal{R}_T^q(\theta^*) \leq \sqrt{T C^* A} \sqrt{\frac{4\alpha \log(T)}{\delta}} + c_\alpha + \sqrt{\frac{\pi \delta}{32\alpha \log(T)}}.$$

*Proof:* We write  $b_t$  the action proposed by the  $\Phi$ -UCB algorithm at time  $t$ , and  $a_t$  the action effectively played according to distribution  $q_t$ . We introduce the notations:  $I_c(T) = \sum_{t=1}^T \mathbb{I}_{[h_{<t}]_{\theta^*}=c}$ , then  $t_c(i) = \min\{t; I_c(t) = i\}$  and for all  $a \in \mathcal{A}$ ,  $I_c(T, a) =$

$\sum_{t=1}^T \mathbb{I}_{[h_{<t}] = c} \mathbb{I}_{a_t = a}$ . The proof mainly follows the lines of [2]. Note that by definition, we want to bound the following term:

$$\mathcal{R}_T^q(\theta^*) = \sum_c \sum_a \Delta_c(a) \mathbb{E} \left( \sum_{i=1}^{I_c^\theta(T)} \mathbb{I}_{b_i = a} \right) \quad (1)$$

**Step one.** Decompose the event  $b_t = a$ . Let us consider a time  $t$  for which  $[h_{<t}] = c$ . Then let us consider a sub-optimal arm  $a$  such that  $\Delta_c(a) > 0$ . Thus it appears that  $b_t = a$  if one of the following conditions holds:

$$(1) \quad \tilde{\mu}_{t,c}(a_c) \leq \mu_c(a_c)$$

$$(2) \quad \tilde{\mu}_{t,c}(a) > \mu_c(a)$$

$$(3) \quad \Delta_c(a) < 2\sqrt{\frac{\alpha \log T}{I_c^\theta(t-1, a)}}$$

Indeed, otherwise we would have

$$\begin{aligned} \tilde{\mu}_c(a_c) &> \mu_c(a_c) = \mu_c(a) + \Delta_c(a) \\ &\geq \mu_c(a) + 2\sqrt{\frac{\alpha \log I_c(t)}{I_c(t-1, a)}} \geq \tilde{\mu}_c(a). \end{aligned}$$

Thus we introduce the quantity  $u_c(a) = \frac{4\alpha \log T}{\Delta_c(a)^2}$ , and deduce that:

$$\mathbb{E} \left( \sum_{i=1}^{I_c^\theta(T)} \mathbb{I}_{b_i = a} \right) \leq \mathbb{E} \left( \sum_{i=1}^{I_c^\theta(T)} \mathbb{I}_{(1) \text{ or } (2) \text{ or } I_c^\theta(t-1, a) < u_c(a)} \right).$$

**Step 2.** Now since  $I_c^\theta(\cdot, a)$  is an increasing function of time (note though, that it does not increase by one each time  $b_t$  is proposed...), we can define the stopping time  $\tau_c(a) = \min\{t; I_c^\theta(t, a) \geq u_c(a)\}$ , or equivalently the stopping instant  $i_c(a) = \min\{i; I_c^\theta(t_c^\theta(i), a) \geq u_c(a)\}$ . Thus we deduce that:

$$\mathbb{E} \left( \sum_{i=1}^{I_c^\theta(T)} \mathbb{I}_{b_i = a} \right) \leq \mathbb{E}(i_c(a)) + \mathbb{E} \left( \sum_{i=i_c(a)+1}^{I_c^\theta(T)} \mathbb{I}_{(1) \text{ or } (2)} \right) \quad (2)$$

Now we can bound the second term of (2) by a constant depending only on  $\alpha$ , by an easy peeling argument (we refer to Section 2.2 of [2]):

$$\mathbb{E} \left( \sum_{i=i_c(a)+1}^{I_c^\theta(T)} \mathbb{I}_{(1) \text{ or } (2)} \right) \leq 2\mathbb{E} \left( \sum_{i=i_c(a)+1}^{I_c^\theta(T)} \left( \frac{\log i}{\log 1/\beta} + 1 \right) \frac{1}{i^{2\beta\alpha}} \right) \quad (3)$$

where  $\beta = \frac{1}{\alpha+1/2}$ .

Then, we also have, by integration by parts:

$$\begin{aligned} 2\mathbb{E} \left( \sum_{i=i_c(a)+1}^{I_c^\theta(T)} \left( \frac{\log i}{\log 1/\beta} + 1 \right) \frac{1}{i^{2\beta\alpha}} \right) &\leq 2 \int_1^\infty \left( \frac{\log t}{\log 1/\beta} + 1 \right) \frac{1}{t^{2\beta\alpha}} dt \\ &\leq \frac{4}{\log(1/\beta)(2\beta\alpha-1)^2}. \end{aligned}$$

**Step 3.** Thus we focus on the first term  $\mathbb{E}(i_c(a))$  of (2). Since we know that  $q_t(a) \geq \delta$  for all  $a, t$ , we thus deduce that:

$$\begin{aligned} \mathbb{E}(i_c(a)) &= \sum_{l=0}^\infty \mathbb{P}(i_c(a) > l) \leq \\ &l_0 + \sum_{l=l_0}^\infty \mathbb{P}(\forall j \leq l \ I_c^\theta(t_c^\theta(j), a) < u_c(a)) \leq \\ &l_0 + \sum_{l=l_0}^\infty \mathbb{P}(\forall j \leq l \ \sum_{i=1}^j \mathbb{I}_{a_{t_c^\theta(i)} = a} - q_{t_c^\theta(i)}(a) < u_c(a) - \delta j). \end{aligned}$$

Now by property of martingale difference sequences, we deduce by setting  $l_0 = \lceil \frac{u_c(a)}{\delta} \rceil$ , that:

$$\begin{aligned} \mathbb{E}(i_c(a)) &\leq l_0 + \sum_{l=l_0}^\infty \exp(-2(l-l_0)^2 \delta^2 l) \\ &\leq l_0 + \sum_{l=l_0}^\infty \exp\left(-\frac{(l-l_0)^2}{2\sigma^2}\right), \end{aligned}$$

where we introduced the quantity  $\sigma^2 = \frac{1}{4\delta^2 l_0}$ . Thus we deduce that:

$$\mathbb{E}(i_c(a)) \leq \lceil \frac{u_c(a)}{\delta} \rceil + \sqrt{\frac{\pi}{8}} \sqrt{\frac{\delta}{u_c(a)}}. \quad (4)$$

**Step 4.** Finally, by combining (3), (4) with (2) and (1), we deduce the following distribution-dependent bound on the rebel regret:

$$\mathcal{R}_T^q(\theta^*) \leq \sum_{c \in \mathcal{H}/\Phi^*} \sum_{a \neq a^*} \Delta_c(a) \left[ \frac{2\alpha \log(T)}{\Delta_c(a)^2 \delta} + \sqrt{\frac{\pi \delta \Delta_c(a)^2}{32\alpha \log T} + c_\alpha} \right],$$

where  $c_\alpha = 1 + \frac{4}{\log(\alpha+1/2)} \left( \frac{\alpha+1/2}{\alpha-1/2} \right)^2$ . We deduce the distribution-free bound by the same argument as for Theorem 1, remarking that  $\sqrt{\frac{\pi}{8}} \sqrt{\frac{\delta \Delta_c(a)^2}{4\alpha \log T}} \leq \sqrt{\frac{\pi}{32\alpha \log(T)}} = c'_\alpha$ .  $\square$

This enables us to deduce the following Theorem, that we prove using the same method as that of Theorem 4 but for the stochastic  $\Phi_\Theta$ -regret of Exp4/UCB.

**Theorem 6** *Assume that we consider a  $\Phi^*$ -constrained opponent with  $\Phi^* \in \Phi_\Theta$ , then the stochastic  $\Phi_\Theta$ -regret of Exp4/UCB is bounded as:*

$$R_T^\Theta = O\left((TA)^{2/3} (\bar{C} \log(T))^{1/3} \log(|\Theta|)^{1/2}\right),$$

where  $\bar{C} = |\mathcal{H}/\Phi^*|$  is the number of classes of the model  $\Phi^*$  of the opponent.

### 3 Approximation error of the models

The following result sheds light on a specific term that appears to be an approximation term of the true model  $\theta^*$  by other models  $\theta$ .

**Theorem 7** For any  $(p_t(\theta))_{t,\theta} \in [0,1]$ , thus for any meta algorithm run on top of Exp3 algorithm and defined with  $q_t(a) = \sum_{\theta} p_t(\theta) \xi_t^\theta(a)$  and decreasing coefficient  $\eta_t^\theta$ , the following holds true:

$$\begin{aligned} \tilde{R}_T^\theta &\leq \mathbb{E} \left( \sum_{\theta} \sum_{c \in \theta} \sum_{i=1}^{I_c^\theta(T)} \frac{\eta_{t_c^\theta(i)}^\theta A}{2} p_{t_c^\theta(i)}(\theta) + \sum_{\theta} \sum_{c \in \theta} \frac{\log A}{\eta_{t_c^\theta(I_c^\theta(T))}^\theta} \right. \\ &\quad \left. + \sum_{\theta} \sum_{c \in \theta} \inf_{a_c^\theta} \sum_{i=1}^{I_c^\theta(T)} (r_{t_c^\theta(i)}(a_{[h_{<t}^\theta(i)]}^*) - r_{t_c^\theta(i)}(a_c^\theta)) p_{t_c^\theta(i)}(\theta) \right). \end{aligned}$$

The term on the second line is actually a mixture of approximation errors of each model, and it seems it can not be reduced without further assumption on the quality of the models.

*Proof:* The proof is in four steps.

**Step 1.** Rewrite the regret to make appear the probabilities  $\xi_t^\theta(a)$ . We first introduce:

$$\begin{aligned} R_T &= \sum_{t=1}^T r_t(a_{[h_{<t}]^*}) - r_t(a_t) \\ &= \sum_{\theta} \sum_{t=1}^T \mathbb{E}_{a_t \sim q_t} \left( \frac{r_t(a_t) p_t(\theta)}{q_t(a_t)} \mathbb{I}_{a_t = a_{[h_{<t}]^*}} \right) \\ &\quad - \frac{r_t(a_t) \xi_t^\theta(a_t) p_t(\theta)}{q_t(a_t)}. \end{aligned}$$

Now we have:  $\tilde{l}_{t,c,\theta}^\theta(a) = (1 - r_t(a)) \frac{p_t(\theta)}{q_t(a)} \mathbb{I}_{a_t = a} \mathbb{I}_{c_\theta = [h_{<t}]_\theta}$ , thus taking the expectation over  $a_t$  for each time  $t$ , we have:

$$\begin{aligned} \tilde{R}_T &= \sum_{t=1}^T \mathbb{E}_{a_t} (r_t(a_{[h_{<t}]^*}) - r_t(a_t)) \\ &= \sum_{\theta} \sum_{t=1}^T \mathbb{E}_{a_t} \left( \frac{p_t(\theta)}{q_t(a_t)} \mathbb{I}_{a_t = a_{[h_{<t}]^*}} - \tilde{l}_{t,[h_{<t}]_\theta}^\theta(a_{[h_{<t}]^*}) \right) \\ &\quad + \sum_{\theta} \sum_{t=1}^T \mathbb{E}_{a_t} \left( \mathbb{E}_{a_t \sim \xi_t^\theta} (\tilde{l}_{t,[h_{<t}]_\theta}^\theta(a)) - \frac{p_t(\theta) \xi_t^\theta(a_t)}{q_t(a_t)} \right). \end{aligned}$$

We can simplify the above expression since  $\mathbb{E}_{a_t} \left( \frac{p_t(\theta)}{q_t(a_t)} \mathbb{I}_{a_t = a_{[h_{<t}]^*}} \right) = \mathbb{E}_{a_t} \left( \frac{p_t(\theta) \xi_t^\theta(a_t)}{q_t(a_t)} \right) = p_t(\theta)$ .

**Step 2.** Decompose the term  $\mathbb{E}_{a_t \sim \xi_t^\theta} (\tilde{l}_{t,[h_{<t}]_\theta}^\theta(a))$  in order to use the definition of  $\xi_t^\theta$ . Indeed, one can upper bound this term by

$$\begin{aligned} \mathbb{E}_{a_t \sim \xi_t^\theta} (\tilde{l}_{t,[h_{<t}]_\theta}^\theta(a)) &\leq \frac{\eta_t^\theta}{2} \mathbb{E}_{a_t \sim \xi_t^\theta} (\tilde{l}_{t,[h_{<t}]_\theta}^\theta(a)^2) \\ &\quad - \frac{1}{\eta_t^\theta} \log \left( \sum_a \exp(-\eta_t^\theta \tilde{l}_{t,[h_{<t}]_\theta}^\theta(a) \xi_t^\theta(a)) \right). \end{aligned}$$

Thus, since by definition we have that  $\xi_t^\theta(a) = \frac{\exp(-\eta_t^\theta \sum_{s=1}^{t-1} \tilde{l}_{s,[h_{<t}]_\theta}^\theta(a))}{\sum_a \exp(-\eta_t^\theta \sum_{s=1}^{t-1} \tilde{l}_{s,[h_{<t}]_\theta}^\theta(a))}$ , we can introduce the quantity  $\Psi_t^\theta(\eta, c) = \frac{1}{\eta} \log \left( \frac{1}{A} \sum_a \exp(-\eta \sum_{s=1}^t \tilde{l}_{s,c}^\theta(a)) \right)$  so that the previous regret term writes:

$$\begin{aligned} \tilde{R}_T &\leq \sum_{\theta} \sum_{t=1}^T \mathbb{E}_{a_t} \left( \frac{\eta_t^\theta}{2} (1 - r_t(a_t))^2 \frac{p_t^2(\theta) \xi_t^\theta(a_t)}{q_t^2(a_t)} \right) \\ &\quad + \sum_{\theta} \left( \sum_{t=1}^T \mathbb{E}_{a_t} (\Psi_{t-1}^\theta(\eta_t^\theta, [h_{<t}]_\theta) - \Psi_t^\theta(\eta_t^\theta, [h_{<t}]_\theta)) \right. \\ &\quad \left. - \mathbb{E}_{a_t} (\tilde{l}_{t,[h_{<t}]_\theta}^\theta(a_{[h_{<t}]^*})) \right). \end{aligned}$$

**Step 3.** Introduce the equivalence classes. We now consider the term in the right hand side of the above equation defined with  $\Psi$  functions. Note that we do not change the bound on the term  $\tilde{R}_T$  by considering the sum over the  $\theta$  such that  $p_t(\theta) > 0$ .

Let us introduce the following notations  $I_c^\theta(t) = \sum_{s=1}^t \mathbb{I}_{c=[h_{<s}]_\theta} \mathbb{I}_{p_s(\theta) > 0}$  and  $t_c^\theta(i) = \min\{t; I_c^\theta(t) = i\}$ . Thus we can write:

$$\begin{aligned} &\sum_{\theta} \sum_{t=1}^T (\Psi_{t-1}^\theta(\eta_t^\theta, [h_{<t}]_\theta) - \Psi_t^\theta(\eta_t^\theta, [h_{<t}]_\theta)) \mathbb{I}_{p_t(\theta) > 0} \\ &= \sum_{\theta} \sum_{c \in \theta} \sum_{i=1}^{I_c^\theta(T)} \Psi_{t_c^\theta(i)-1}^\theta(\eta_{t_c^\theta(i)}^\theta, c) - \Psi_{t_c^\theta(i)}^\theta(\eta_{t_c^\theta(i)}^\theta, c) \\ &= \sum_{\theta} \sum_{c \in \theta} \sum_{i=1}^{I_c^\theta(T)-1} (\Psi_{t_c^\theta(i)}^\theta(\eta_{t_c^\theta(i)+1}^\theta, c) - \Psi_{t_c^\theta(i)}^\theta(\eta_{t_c^\theta(i)}^\theta, c) \\ &\quad - \Psi_{t_c^\theta(I_c^\theta(T))}^\theta(\eta_{t_c^\theta(I_c^\theta(T))}^\theta, c)). \end{aligned}$$

Now, by definition of  $\Psi_t^\theta$ , the last term of this sum,  $-\Psi_{t_c^\theta(I_c^\theta(T))}^\theta(\eta_{t_c^\theta(I_c^\theta(T))}^\theta, c)$ , is equal to the following quantity

$$\frac{\log A}{\eta} - \frac{1}{\eta} \log \left( \frac{1}{A} \sum_a \exp(-\eta \sum_{s=1}^{t_c^\theta(I_c^\theta(T))} \tilde{l}_{s,c}^\theta(a)) \right),$$

where  $\eta = \eta_{t_c^\theta(I_c^\theta(T))}^\theta$ , which can be upper bounded by  $\frac{\log A}{\eta} + \sum_{s=1}^{t_c^\theta(I_c^\theta(T))} \tilde{l}_{s,c}^\theta(a)$  for any given  $a = a_c^\theta$ .

**Step 4.** Now since  $\eta_{t_c^\theta(i)}^\theta \leq \eta_{t_c^\theta(i)+1}^\theta$  and  $\Psi_{t_c^\theta(i)}^\theta(\cdot, c)$  is increasing for all  $\theta, c$ , we deduce from the previous equations that:

$$\begin{aligned} \tilde{R}_T &\leq \sum_{\theta} \sum_{c \in \theta} \sum_{i=1}^{I_c^\theta(T)} \sum_a \frac{\eta_{t_c^\theta(i)}^\theta p_{t_c^\theta(i)}^2(\theta) \xi_{t_c^\theta(i)}^\theta(a)}{2 q_{t_c^\theta(i)}(a)} \\ &\quad + \sum_{\theta} \sum_{c \in \theta} \frac{\log A}{\eta_{t_c^\theta(I_c^\theta(T))}^\theta} \\ &\quad + \sum_{\theta} \sum_{c \in \theta} \inf_{a_c^\theta} \sum_{t=1}^{t_c^\theta(I_c^\theta(T))} \mathbb{E}_{a_t} (\tilde{l}_{t,c}^\theta(a_c^\theta) - \tilde{l}_{t,c}^\theta(a_{[h_{<t}]^*})). \end{aligned}$$

Now we conclude by taking the expectation, seeing that  $p_{t_c^\theta(i)}(\theta) \xi_{t_c^\theta(i)}^\theta \leq q_{t_c^\theta(i)}(a)$ , and that by definition  $\mathbb{E}_{a_t} (\tilde{l}_{t,c}^\theta(a_c^\theta)) = (1 - r_t(a_c^\theta)) p_t(\theta) \mathbb{I}_{c=[h_{<t}]_\theta}$ .

□

## References

- [1] Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32:48–77, January 2003.
- [2] S. Bubeck. *Bandits Games and Clustering Foundations*. PhD thesis, Université Lille 1, 2010.