

---

# Learning how contextual bandits learn

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In cognitive modeling, understanding how an agent leverages contextual informa-  
2 tion to learn about an adversarial environment and take what it considers good  
3 decisions is a fundamental investigation. By observing the agent’s learning pro-  
4 cess, can we estimate how the agent is using this contextual information? One  
5 way of doing this is to approximate the agent’s learning behavior by contextual  
6 bandit algorithms. The aim of this work is to provide model selection procedures  
7 that will pick the contextual bandit procedure that best fits the agent’s learning  
8 process. We introduce a hold-out estimator and a penalized maximum likelihood  
9 estimator and show that both satisfy oracle inequalities. We give several examples  
10 of bandit algorithms for which the assumptions are satisfied, and assess our results  
11 on both synthetic and experimental learning data in a human categorization task.  
12 We also discuss why bandits with expert advice satisfy the same type of oracle  
13 inequalities and how they can be used to model metalearning in cognition.

## 14 1 Introduction

### 15 1.1 Cognitive models

16 Imagine an agent (human or animal) learning sequentially to make good decisions and having  
17 access at each time step to some contextual information. By looking at the agent’s successive actions,  
18 can we estimate the agent’s learning strategy, that is the way the agent used this contextual infor-  
19 mation to make its decisions? This problem belongs to the more general framework of cognitive  
20 modeling [12]. Cognitive models help to understand the mechanisms that occur while for instance  
21 learning, remembering or predicting tasks. They have been widely studied in the cognition liter-  
22 ature [32, 15] and have a major impact on education for example. Usually in cognitive modeling  
23 [42, 16], maximum likelihood estimation (MLE) is applied and the best cognitive model is selected  
24 by cross-validation or an Akaike information criterion (AIC). One of the main challenges of cog-  
25 nitive modeling on learning data is that, since the agent remembers its past actions to learn, the  
26 data are not stationary and not independent. There are very few theoretical statistical works in this  
27 context: in [5], the properties of the MLE are studied for the Exp3 model on learning data; in [6], a  
28 very general model selection procedure is presented that can be applied to non stationary data, but  
29 nothing in the setup of learning with contextual information. Our present goal is to provide model  
30 selection procedures that are valid for learning data in this contextual setup.

### 31 1.2 Contextual bandits

32 The purpose of a contextual bandit algorithm [28] is to find an optimal policy for selecting actions  
33 based on additional information (the context) given at each time step. In Machine Learning, contex-  
34 tual bandits have many applications [10] such as recommendation, patient follow-up in healthcare,  
35 etc. Here, we use them as learning models. Although not traditionally employed in cognition for  
36 modeling real behavioral data, contextual bandits are gaining popularity in the cognition literature  
37 [27, 41] and most of the cognitive psychology models of learning with contextual data such as Com-  
38 ponent Cue [20] or Alcove [26] can be expressed as contextual bandit algorithms since they treat the  
39 same problem: bandit feedback and choice based on past decisions and present context.

40 Let us formalize the statistical problem we treat. We observe a sequence of contexts and actions  
41  $(X_1, A_1, \dots, X_T, A_T)$  for an integer  $T \geq 1$ , where the contexts  $X_t$  belong to some finite space  $\mathcal{X}$   
42 and the actions  $A_t$  belong to a finite set  $[K] = \{1, \dots, K\}$  for some integer  $K \geq 1$ . Let  $\mathcal{F}_0$  be the  
43 trivial sigma-algebra and for  $t \geq 1$ , let  $\mathcal{F}_t = \sigma(X_1, A_1, \dots, X_t, A_t)$  be the sigma-algebra generated  
44 by observations until time  $t$ . Let  $p^* = (p_t^*)_{t \in [T]}$  be the successive conditional probability distribu-  
45 tions:  $\forall x \in \mathcal{X}, \forall a \in [K], p_t^*(a, x) = \mathbb{P}(A_t = a | X_t = x, \mathcal{F}_{t-1})$ . In reinforcement learning, this  
46 vector is called the *policy* of the agent. Recall that here,  $p^*$  is fixed, but unknown.  
47 Our goal is to select the best model approximating  $p^*$  among a family of models  $(\{p^m =$   
48  $(p_t^m)_{t \in [T]}\})_{m \in \mathcal{M}}$ , where  $\mathcal{M}$  is a countable set. Each  $p_t^m$  is a conditional distribution over  $[K]$   
49 given  $(X_t, \mathcal{F}_{t-1})$  and a candidate at being  $p_t^*$ .

### 50 1.3 Partition-based contextual bandits: an example of parametric models

51 The leading example of contextual bandit algorithm that we use here is *partition-based contextual*  
52 *bandits* [28, Chapter 18]. It consists in assuming that the agent partitions the context space  $\mathcal{X}$  into  
53 disjoint cells  $C$ . This may typically happen if the agent is already familiar with the contexts and  
54 has already built a personal opinion on their meaning. The agent only has to learn the new task  
55 thanks to this fixed view of the space by updating elementary bandit algorithms in each cell  $C$ , that  
56 we denote  $\text{CellBandit}(C)$ , each time the context belongs to the corresponding cell. Our goal is to  
57 estimate the partitioning of the context space that the agent is using, *i.e.* understanding how the agent  
58 uses the contexts for the learning task. As an example, we illustrate numerically our approach on a  
59 categorization task, see Section 5 where contexts are objects to classify. By selecting the partition  
60 that best fits the learning data of a given individual, we have access to the similarity between objects  
61 as perceived by the learner.

62 To formalize *partition-based contextual bandits*, let  $g_t = (g_{1,t}, \dots, g_{K,t}) \in [0, 1]^K$  be the vector of  
63 losses (or rewards) at time  $t$ , which models the feedback of the environment. We make no particular  
64 assumptions on the way losses are generated, except that  $g_t$  needs to be  $\sigma(X_t, \mathcal{F}_{t-1})$ -measurable.  
65 They may be adversarial or stochastic (see Section 4 for some examples). In the same way, the gener-  
66 ation of the contexts  $X_t$  does not need to be specified: they can be independent of past actions or  
67 the result of the past actions. Then, each model  $m \in \mathcal{M}$  corresponds to a partition  $\mathcal{P}_m$  of  $\mathcal{X}$  into  $D_m$   
68 cells. The model  $m$  is parameterized by a vector  $\theta^m = (\theta_C)_{C \in \mathcal{P}_m}$ , where each  $\text{CellBandit}(C)$  is  
69 using a procedure parameterized by a parameter  $\theta_C$  – for instance, the learning rate in Exp3. The re-  
70 sulting candidate for  $p^*$  is therefore  $p_{\theta^m}^m = (p_{\theta^m, t}^m)_{t \in [T]}$ . For a given cell  $C \in \mathcal{P}_m$ ,  $\text{CellBandit}(C)$   
71 is updated each time  $X_t \in C$ , and therefore its decision at time  $t$  only depends on the contexts and  
72 actions happening at times in  $F_t(C) = \{s \in [t], X_s \in C\}$ , which is of cardinality  $T_t^C = |F_t(C)|$ .  
73 We write  $a \mapsto \pi_{C, T_t^C}^{\theta_C}(a)$  the distribution over the set of actions  $[K]$  at time  $t$  for the procedure  
74  $\text{CellBandit}(C)$  with parameter  $\theta_C$  (see Algorithm 1). With this notation,

$$\forall t \in [T], \forall a \in [K], p_{\theta^m, t}^m(a, X_t) = \mathbb{P}_{\theta^m}^m(A_t = a | X_t, \mathcal{F}_{t-1}) = \sum_{C \in \mathcal{P}_m} \pi_{C, T_t^C}^{\theta_C}(a) \mathbf{1}_{X_t \in C}. \quad (1)$$

---

#### Algorithm 1 Partition-based contextual bandit for model $m$ [28]

---

**Inputs:** partition  $\mathcal{P}_m$  of the context space  $\mathcal{X}$ ,

parameters  $\theta^m = (\theta_C)_{C \in \mathcal{P}_m} \in \Theta^m = \bigotimes_{C \in \mathcal{P}_m} \Theta_C$ , with  $\Theta_C$  compact parametric set.

**Initialization:** For all  $C \in \mathcal{P}_m$ , for all  $a \in [K]$ ,  $\pi_{C, 1}^{\theta_C}(a) = 1/K$ .

**for**  $t = 1, 2, \dots$  **do**

Learner observes context  $X_t \in \mathcal{X}$  and finds  $C \in \mathcal{P}_m$  such that  $X_t \in C$ .

Learner plays  $\text{CellBandit}(C)$  with parameter  $\theta_C$  and samples action  $A_t \sim \pi_{C, T_t^C}^{\theta_C}$ .

Learner observes loss  $g_{A_t, t}$  and updates the probability distribution  $\pi_{C, T_t^C}^{\theta_C}$  in  $\text{CellBandit}(C)$ .

---

### 75 1.4 Contributions

76 We provide two model selection procedures for modeling learning with contextual information,  
77 based on the partial log-likelihood  $\ell_T(p^m)$  of the observations  $(X_1, A_1, \dots, X_T, A_T)$ , defined

78 by

$$\ell_T(p^m) = \sum_{t=1}^T \log(p_t^m(A_t, X_t)). \quad (2)$$

79 We prove oracle inequalities for the conditional Kullback-Leibler divergence  $D_{\text{KL}}$  between  $p_t^*(\cdot, X_t)$   
80 and  $p_t^m(\cdot, X_t)$ :

$$D_{\text{KL}}(p_t^*(\cdot, X_t), p_t^m(\cdot, X_t)) = \mathbb{E} \left[ \log \frac{p_t^*(\cdot, X_t)}{p_t^m(\cdot, X_t)} \middle| X_t, \mathcal{F}_{t-1} \right].$$

81 In Section 2 we consider a finite family of general models  $\{p^m = (p_t^m)_{t \in [T]}, m \in \mathcal{M}\}$  and show that  
82 a hold-out estimator satisfies an oracle inequality with an  $\mathcal{O}((\log T + \log |\mathcal{M}|)/T)$  error bound,  
83 regardless of the nature of the models. In Section 3, we focus on the partition-based contextual  
84 bandit models defined in (1) with possibly infinite countable family of partitions and consider a  
85 log-likelihood criterion penalized by  $D_m$  times some logarithmic terms. Under some assumptions  
86 on the CellBandid algorithms that are used, we show an oracle inequality with an  $\mathcal{O}(\log(T)^3/T)$   
87 error bound. In Section 4, we prove that Stochastic Gradient Bandits and Exp3-IX are examples of  
88 CellBandid for which assumptions of Section 3 are satisfied. Section 5 is devoted to numerical  
89 illustrations on both synthetic and experimental learning data in a categorization task. In Section 6,  
90 we discuss how bandits with expert advice can be used to model metalearning [9], which refers  
91 to the processes by which an agent acquires knowledge about its own learning abilities, strategies,  
92 and preferences. In Appendix B, we give the details required to obtain model selection results for  
93 metalearning. The complete proofs of the theoretical results are given in Appendix C.

## 94 1.5 Related work

95 Our objective is not to provide a method that improves the regret [11]. Similarly, our work is not  
96 to be misunderstood with [17, 18, 37] in which authors develop model selection algorithms for  
97 contextual bandits that aim at finding the relation between context and action that best optimizes  
98 rewards. Our goal is to understand how an agent learns, not to tell it how to learn better. Thanks  
99 to the learning data of an agent, we select the contextual bandit algorithm that best fits the learning  
100 curve of the agent – without necessarily assuming that the agent understands the relation between  
101 context and actions. Hence we are not trying to find an optimal model, but the most realistic one  
102 w.r.t. learning data. To our knowledge, this theoretical statistical problem was studied for the first  
103 time in [5]. But in contrast with [5], which assumes Exp3 to be true and studies MLE performances,  
104 we want to perform model selection with contexts.

105 From an Imitation Learning (IL) or Inverse Reinforcement Learning (IRL) point of view, this prob-  
106 lem could be seen as a learner trying to reproduce the learning curve of an expert. Usually in IL [23],  
107 we observe an expert who has already mastered the task, so the input data of a classic IL algorithm  
108 are not learning data. In IRL [4], MLE might be used on data [38] but the IRL learner’s goal is  
109 to infer the underlying reward function that best explains the expert’s observed behavior thanks to  
110 multiple trajectories and then use this inferred reward function to guide its own decision-making. In  
111 our setting, the experimentalist already knows the reward function and the goal is to infer the agent’s  
112 perception of the contexts, thanks to a single learning trajectory.

113 Our goal is close to [24], who estimate how a learner’s behavior evolves over time and how it priori-  
114 tises choices for applications to healthcare, except that [24] is in a Bayesian framework. Similarly,  
115 authors in [40] and [41] try to predict the behavior of participants in contextual multi-armed bandit  
116 tasks. The main difference is that they work in specific stochastic bandit settings with a Bayesian  
117 approach whereas we do model selection in a non-stationary and adversarial framework.

118 On a more technical level, hold-out estimators are often used in cognition for learning data [33, 25].  
119 Hold-out procedures have been studied theoretically in the literature in a stationary and independent  
120 data framework [29, 3, 2]. Few results exist for time dependent data [36] and they are quite far from  
121 our setup. Here, the main issue is that the training set is not independent from the validation set, so  
122 more advanced tools such as  $V$ -fold cross-validation cannot be used.

123 Section 1.3 is very similar in design to the framework of [13] for selecting the best histogram for  
124 density estimation or more generally to non asymptotic model selection [29]. The main difference  
125 is that we are in a non stationary and non independent framework. Therefore, to prove the oracle  
126 inequality of Section 3, we use instead a recent result for penalized log-likelihood estimators which  
127 is valid in this framework [6].

## 128 2 Hold-out estimator

129 In this section, we assume that  $\mathcal{M}$  is finite as it is often the case for hold-out estimators [29, Chapter  
130 8], and  $|\mathcal{M}| \geq 2$ . Let  $T > N \geq 1$  and select  $\hat{m} \in \arg \max_{m \in \mathcal{M}} \sum_{t=N}^T \log p_t^m(A_t, X_t)$ .

131 **Theorem 1.** *Assume that for all  $m \in \mathcal{M}$  and for all  $t \in \{N, \dots, T\}$ ,  $p_t^m$  depends only on  $(X_t, \mathcal{F}_{t-1})$ .  
132 There exists a positive numerical constant  $\diamond$ , such that for any  $\kappa \in (0, 1)$ ,*

$$\begin{aligned} & (1 - \kappa) \mathbb{E} \left[ \frac{1}{T - N + 1} \sum_{s=N}^T \text{D}_{\text{KL}} \left( p_s^*(\cdot, X_s), \frac{p_s^*(\cdot, X_s) + p_s^{\hat{m}}(\cdot, X_s)}{2} \right) \middle| X_N, \mathcal{F}_{N-1} \right] \\ & \leq (1 + \kappa) \inf_{m \in \mathcal{M}} \mathbb{E} \left[ \frac{1}{T - N + 1} \sum_{s=N}^T \text{D}_{\text{KL}} (p_s^*(\cdot, X_s), p_s^m(\cdot, X_s)) \middle| X_N, \mathcal{F}_{N-1} \right] \\ & \quad + \frac{\diamond \log(T - N + 1) + \log |\mathcal{M}|}{\kappa (T - N + 1)}. \end{aligned}$$

133 This result can hold for arbitrary  $p^m$  as long as it is adapted to  $\mathcal{F}_t$ , for  $t \geq N$ . In particular it allows,  
134 as usual for hold-out estimator, to use  $p^m = p_{\hat{\theta}^m}^m$ , where  $\hat{\theta}^m \in \arg \max_{\theta^m \in \Theta^m} \sum_{t=1}^{N-1} \log p_{\theta^m, t}^m(A_t, X_t)$ ,  
135 whatever the parameterization of the model  $m$  – not necessarily partition-based. This result is the  
136 equivalent of Theorem 8.9 in [29] for this learning framework, adding only a multiplicative  $\log T$   
137 factor in the error bound. It justifies the use of hold-out procedures to model learning data in cogni-  
138 tive experiments such as [33, 25], using classical cognitive models as Alcové [26], Component-Cue  
139 [20] or Activity-based Credit Assignment (see [25] and the references therein).

140 *Limitations.* Due to the strongly dependent structure of the data, we perform a single split of the  
141 sample between training and testing data at  $t = N$ , unlike the classical hold-out. As usual, a careful  
142 trade-off has to be performed between  $N$  large enough to properly estimate each model and not too  
143 large, in order to reliably compare them. This also means that this approach is unsuited to situations  
144 where the learner learns differently at the start and at the end of the experiment, for instance by  
145 switching models once it has grasped how the task worked.

## 146 3 Penalized maximum likelihood estimator

147 In this section, we restrict ourselves to partition-based contextual bandit (see (1)). Following [5], we  
148 need to assume that the probabilities do not vanish.

149 **Assumption 1.** There exists  $\varepsilon > 0$  and an integer  $T_\varepsilon \geq 2$ , such that, almost surely,

$$\forall t \leq T_\varepsilon, \forall x \in \mathcal{X}, \forall a \in [K], \quad p_t^*(a, x) \geq \varepsilon \quad (3)$$

150 and that for all  $m \in \mathcal{M}$  and all  $C \in \mathcal{P}_m$ , the `CellBandit`( $C$ ) satisfies, for all parameter  $\theta_C \in \Theta_C$

$$\forall t \leq T_\varepsilon, \forall a \in [K], \quad \pi_{C, T_t^C}^{\theta_C}(a) \geq \varepsilon. \quad (4)$$

151 Let  $\text{pen} : \mathcal{M} \rightarrow \mathbb{R}_+$  be a penalty function. For each  $m \in \mathcal{M}$ , let  $\hat{\theta}^m \in \arg \max_{\theta^m \in \Theta^m} \ell_{T_\varepsilon}(p_{\theta^m}^m)$  be  
152 a MLE of model  $m$ , with  $\ell$  defined as in (2), and select a model  $\hat{m}$  that minimizes the penalized  
153 log-likelihood stopped at  $T_\varepsilon$ :

$$\hat{m} \in \arg \min_{m \in \mathcal{M}} \left( -\frac{\ell_{T_\varepsilon}(p_{\hat{\theta}^m}^m)}{T_\varepsilon} + \text{pen}(m) \right). \quad (5)$$

154 To prove oracle inequalities, we also need a smoothness assumption on the parameterization of  
155 `CellBandit`( $C$ ) which can then be propagated to the  $p^m$  in Proposition 2.

156 **Assumption 2.** With the notation of Assumption 1, there exists  $L_\varepsilon > 0$  such that, almost surely, for  
157 all  $m \in \mathcal{M}$ , all  $C \in \mathcal{P}_m$ ,

$$\forall \delta_C, \theta_C \in \Theta_C, \forall t \leq T_\varepsilon, \quad \sup_{a \in [K]} \left| \log \left( \frac{\pi_{C, T_t^C}^{\delta_C}(a)}{\pi_{C, T_t^C}^{\theta_C}(a)} \right) \right| \leq L_\varepsilon \|\delta_C - \theta_C\|_2. \quad (6)$$

158 **Proposition 2.** Assume that  $p^m$  is a partition-based contextual bandit as in (1) or Algorithm 1 and  
 159 that there exists  $T_\varepsilon$  such that for all  $C \in \mathcal{P}_m$ ,  $\text{CellBandit}(C)$  satisfies (4) and (6). Then, almost  
 160 surely, for all  $\theta^m, \delta^m \in \Theta^m$ , for all  $t \leq T_\varepsilon$ , for all  $x \in \mathcal{X}$ , for all  $a \in [K]$ ,

$$p_{\theta^m, t}^m(a, x) \geq \varepsilon \quad \text{and} \quad \sup_{a \in [K]} \left| \log \left( \frac{p_{\delta^m, t}^m(a, x)}{p_{\theta^m, t}^m(a, x)} \right) \right| \leq L_\varepsilon \sup_{C \in \mathcal{P}_m} \|\delta_C - \theta_C\|_2.$$

161 Assume that the numbers of parameters of all  $\text{CellBandit}$  procedures are uniformly bounded, and  
 162 let  $d = \sup_{m \in \mathcal{M}} \sup_{C \in \mathcal{P}_m} \dim(\Theta_C)$ . Since the models are smooth enough and the probabilities  
 163 are lower bounded, by applying [6], one can prove the following result.

164 **Theorem 3.** Let  $\mathcal{M}$  be a countable set, and for each  $m \in \mathcal{M}$ , consider a partition-based contextual  
 165 bandit model  $\{p_{\theta^m}^m, \theta^m \in \Theta^m\}$  (see Algorithm 1 and (1)). Let  $R$  and  $r$  be such that all coordinates  
 166  $\theta_{i,C}$ 's of  $\theta_C \in \Theta_C$ , for  $C \in \mathcal{P}_m$  and  $m \in \mathcal{M}$ , satisfy  $r \leq \theta_{i,C} \leq R$  and let  $A_\varepsilon = L_\varepsilon \sqrt{d}(R -$   
 167  $r) + 2 \log(\varepsilon^{-1})$ . Let  $\Sigma_\varepsilon = \log(A_\varepsilon) \sum_{m \in \mathcal{M}} e^{-D_m} < +\infty$ . Under Assumptions 1 and 2, there exist  
 168 positive numerical constants  $c$  and  $c'$  such that for all  $\kappa \in (0, 1]$ , the following holds: if for all  
 169  $m \in \mathcal{M}$ ,

$$\text{pen}(m) \geq \frac{c}{\kappa} A_\varepsilon^2 \log(\varepsilon^{-1})^{3/2} \log(T_\varepsilon A_\varepsilon)^2 \frac{D_m}{T_\varepsilon},$$

170 then,

$$\begin{aligned} & \frac{1 - \kappa}{T_\varepsilon} \sum_{t=1}^{T_\varepsilon} \mathbb{E} \left[ \text{D}_{\text{KL}} \left( p_t^*(\cdot, X_t), p_{\theta^m, t}^m(\cdot, X_t) \right) \right] \\ & \leq \inf_{m \in \mathcal{M}} \left( (1 + \kappa) \inf_{\theta^m \in \Theta^m} \frac{1}{T_\varepsilon} \sum_{t=1}^{T_\varepsilon} \mathbb{E} \left[ \text{D}_{\text{KL}} \left( p_t^*(\cdot, X_t), p_{\theta^m, t}^m(\cdot, X_t) \right) \right] + 2 \text{pen}(m) \right) \\ & \quad + \frac{18c'}{\kappa} A_\varepsilon \Sigma_\varepsilon \log(\varepsilon^{-1})^{3/2} \log(T_\varepsilon A_\varepsilon)^2 \frac{\log(T_\varepsilon)}{T_\varepsilon}. \end{aligned}$$

171 This result is very similar to model selection "à la Birgé-Massart" [29, Section 7.4] with a trade-off  
 172 between bias and variance represented by the penalty in  $D_m/T_\varepsilon$ , with additional logarithmic terms  
 173 in  $\log^2 T_\varepsilon$  in the penalty and in  $\log^3 T_\varepsilon$  in the residual error. It is obtained by applying the recent and  
 174 very general result of [6] which holds even for dependent non stationary data. However it is quite  
 175 tedious to validate the assumptions of [6]. The partition-based contextual bandits are an example  
 176 where this holds easily thanks to the partition which involves easier assumptions (namely (4) and  
 177 (6)) to check on the  $\text{CellBandit}$  (see Section 4 for examples of  $\text{CellBandit}$  that satisfy them).  
 178 Another example of contextual bandit where the assumptions of [6] are satisfied is given Section 6  
 179 in metalearning.

180 *Limitations.* Compared to the hold-out procedure, this approach does not require to split the sample.  
 181 While this estimator can still work well when using all data, as shown in Section 5, the oracle  
 182 inequality only holds when using the data from the time interval  $[T_\varepsilon]$ , which can be significantly  
 183 less: in the models of Section 4,  $T_\varepsilon$  is of order  $\sqrt{T}$ . Moreover, this theorem does not cover every  
 184 kind of cognitive models. Finally, while the penalty is in  $c \log^2(T_\varepsilon) D_m/T_\varepsilon$ , the constant  $c$  in this  
 185 theoretical result is not known a priori and one needs to calibrate it by numerical simulations (see  
 186 Section 5). We could use the hold-out procedure described in Section 2 to choose it, with similar  
 187 issues, or other heuristics such as the dimension jump method or slope heuristics [7, 1].

## 188 4 Examples of CellBandit

189 In this section we provide examples of  $\text{CellBandit}$  satisfying (4) and (6). All the algorithms below  
 190 are written for a cell  $C$  and a  $\text{CellBandit}(C)$  parameterized by  $\theta_C \in \Theta_C$  compact subset of  $\mathbb{R}^d$   
 191 such that  $R \geq \sup_{\theta_C \in \Theta_C} \|\theta_C\|_\infty$ .

### 192 4.1 Example 1: Exp3-IX

193 This algorithm is a generalization of Exp3 and was introduced in [35]. Following [5], we write  
 194 Exp3-IX with parameters decreasing as a square root of the sample size to ensure a good MLE  
 195 estimation of the parameters. Note in addition that, for Exp3 and its variants, it is well known that

196 sublinear convergence of the regret occurs when the learning rate  $\eta$  and the exploration term  $\gamma$  are  
 197 decreasing as a square root of the sample size. This renormalization ensures that the learner is able  
 198 to learn at a good pace and at the same time be robust to changes in the environment.

---

**Algorithm 2** Exp3-IX[35] as a CellBandid( $C$ )

---

**Inputs:**  $T$  (Sample size),  $\theta_C = (\eta, \gamma) \in \Theta_C$  (Parameter),  $K$  (Number of actions).

**Initialization:**  $\pi_{C,1}^{\theta_C} = (\frac{1}{K}, \dots, \frac{1}{K})$ .

**for**  $t \in F_T(C)$ , the set of times where  $X_s \in C$ , **do**

Draw an action  $A_t \sim \pi_{C,T_t}^{\theta_C}$  and receive a loss  $g_{A_t,t} \in [0, 1]$ .

Update for all  $a \in [K]$ ,

$$\pi_{C,T_t+1}^{\theta_C}(a) = \frac{\exp\left(-\frac{\eta}{\sqrt{T}} \sum_{s \in F_t(C)} \hat{g}_{a,s}^{\theta_C}\right)}{\sum_{b \in [K]} \exp\left(-\frac{\eta}{\sqrt{T}} \sum_{s \in F_t(C)} \hat{g}_{b,s}^{\theta_C}\right)} \quad \text{where } \hat{g}_{b,s}^{\theta_C} = \frac{g_{b,s}}{\gamma/\sqrt{T} + \pi_{C,T_s}^{\theta_C}(b)} \mathbf{1}_{A_s=b}$$


---

199 In this case,  $\Theta_C \subset \mathbb{R}^2$ . When  $\gamma = 0$ , we recover the classical Exp3 algorithm, studied from the  
 200 MLE point of view in [5]. Note that while  $g_{A_t,t}$  is observed and known, the estimated loss  $\hat{g}_{b,s}^{\theta_C}$   
 201 depends on the parameterization. The following result shows that one can choose Exp3-IX as a  
 202 CellBandid in the partition-based contextual bandits to perform partition selection.

203 **Proposition 4.** *Let  $\varepsilon \in (0, 1/K)$  and let  $\Theta_C \subset [0, R]^2$  with  $R > 0$ . Then Exp3 – IX can be a*  
 204 *CellBandid( $C$ ) with parameterization  $\theta_C \in \Theta_C$  that satisfies (4) and (6), as soon as*

$$T_\varepsilon = \left[ \left( \frac{1}{K} - \varepsilon \right) \frac{\sqrt{T}}{R} \right] \wedge T \quad \text{and} \quad L_\varepsilon = \frac{\sqrt{R^2/T + \varepsilon^2}}{\varepsilon^3 R} e^{1/\varepsilon^2}.$$

205 This shows that one can apply Theorem 3 with Exp3-IX as CellBandid as long as we stop us-  
 206 ing observations after  $\sqrt{T}$  time steps. The dependence in  $\varepsilon$  is not very critical, since it has been  
 207 proved at least for Exp3 in [5], that in practice, we may take  $\varepsilon$  quite large (non-vanishing) with  
 208 almost no impact on  $T_\varepsilon$ . This is a good thing since the theoretical dependency of  $L_\varepsilon$  in  $\varepsilon$  is quite  
 209 pessimistic.

210 *Limitations.* This algorithm considers the horizon  $T$  fixed in order to renormalize the parameteriza-  
 211 tion. From Proposition 4, it follows that Theorem 3 holds when only the first  $\sqrt{T}$  observations are  
 212 used in the MLE, but this in no way means that the estimator will perform poorly when based on  
 213 all data. Taking  $\sqrt{T}$  observations compounds on the usual issue that if the number of cells is large,  
 214 only a small amount of data may be available for each cell, making estimation difficult.

## 215 4.2 Example 2: Gradient Bandit

216 Gradient Bandit is another possible algorithm. We still choose for similar reason a parameteri-  
 217 zation in  $\eta/\sqrt{T}$ , which echoes the Robbins-Monro conditions [39] even if [31] proved convergence  
 218 in a stochastic bandit framework even for non renormalized parameters.

---

**Algorithm 3** Gradient Bandit [31] as a CellBandid

---

**Inputs:**  $T$  (Sample size),  $\theta_C \in [r, R]$  (Parameter),  $K$  (Number of actions).

**Initialization:**  $\pi_{C,1}^{\theta_C} = (\frac{1}{K}, \dots, \frac{1}{K})$ .

**for**  $t \in F_T(C)$  **do**

Draw an action  $A_t \sim \pi_{C,T_t}^{\theta_C}$  and receive a reward  $g_{A_t,t} \in [0, 1]$ .

Update for all  $a \in [K]$ ,

$$\pi_{C,T_t+1}^{\theta_C}(a) = \frac{\exp\left(-\frac{\theta_C}{\sqrt{T}} \sum_{s \in F_t(C)} \hat{g}_{a,s}^{\theta_C}\right)}{\sum_{b \in [K]} \exp\left(-\frac{\theta_C}{\sqrt{T}} \sum_{s \in F_t(C)} \hat{g}_{b,s}^{\theta_C}\right)} \quad \text{where } \hat{g}_{b,s}^{\theta_C} = \left(\mathbf{1}_{A_s=b} - \pi_{C,T_s}^{\theta_C}(b)\right) g_{A_s,s}$$


---

219 **Proposition 5.** Let  $\varepsilon \in (0, 1)$  and let  $\Theta_C \subset [0, R]^2$  with  $R > 0$ . Then, Gradient Bandit can be a  
 220 CellBandit( $C$ ) with parameterization  $\theta_C \in \Theta_C$  that satisfies (4) and (6), as soon as

$$T_\varepsilon := \left\lceil \log \left( \sqrt{\frac{1}{K\varepsilon}} \right) \frac{\sqrt{T}}{R} \right\rceil \wedge T \quad \text{and} \quad L_\varepsilon = \frac{\sqrt{2} \log \left( \sqrt{\frac{1}{K\varepsilon}} \right)}{R\varepsilon \sqrt{K\varepsilon}}.$$

221 This theoretical result has the same interpretation as before: the theoretical guarantees of Theorem  
 222 3 with Gradient Bandit as CellBandit hold when we stop using observations after  $\sqrt{T}$  time  
 223 steps. In practice, we can use the observations up to time  $T$  (see Section 5).

## 224 5 Numerical illustrations

225 We consider an experiment on the following categorization task: learners have to classify nine objects  
 226 in two categories  $A$  and  $B$  in a sequential way. Figure 1 presents the objects and the classifica-  
 227 tion rule the learners have to learn. It is a quite difficult task that has been experimented for instance  
 228 in [34], where the learners needed about 300 trials to learn the classification rule.

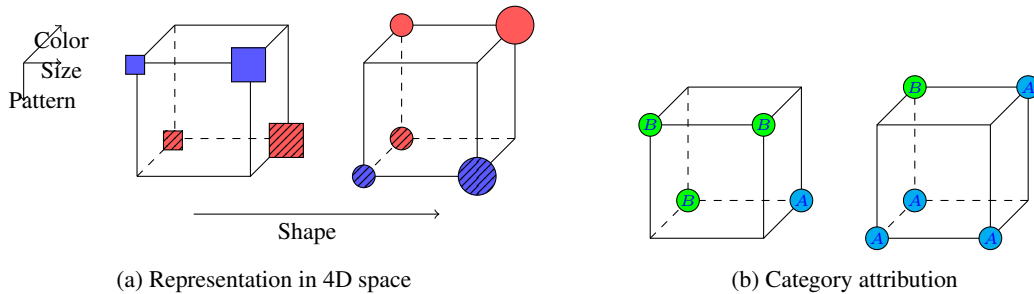


Figure 1: Experiment presentation: classic 5-4 category structure, widely used in cognition [30]. In 1a, the 9 objects to classify represented in a 4D space with respect to their attributes: Color, Size, Filling Pattern, and Shape. In 1b, by position in the 4D space, the category attribution (A or B).

229 The reward is fixed: 1 if the learner finds the good category and 0 in the other case. We focus on 6  
 different models (described in Table 1) that have a good cognitive interpretation.

Table 1: Description of models and their learning abilities

Model	Number of cells	Description of the cells	Learns categorization
OneForAll	1	One giant cell	No
ByShape	2	One for circles, one for squares	Partly
ByPattern	2	One for striped items, one for plain items	Partly
ByShapeExc	4	Cells from ByShape model with exceptions isolated	Yes
ByPatternExc	4	Cells from ByPattern model with exceptions isolated	Yes
OnePerItem	9	One cell for each item	Yes

230

231 **On synthetic data.** We have not been able to run the simulation with Exp3-IX. Indeed, as also  
 232 shown practically in [5] for the simple Exp3 case, the probabilities  $\pi_{C,T}^{\theta_C}$  can go to zero extremely  
 233 fast. When the agent learns over an horizon  $T = 500$ , only  $\sqrt{T} = 22$  observations would be usable  
 234 and the estimations even of just the MLE is unreliable. So all the simulations were performed with

235  $T = 500$  and for Gradient Bandit as CellBandit, for all the 6 models described in Table 1. The  
 236 way synthetic data are generated can be found in Appendix A.1.

237 Figure 2a shows that despite the conservative theoretical bound given in Proposition 5 with  $T_\epsilon$   
 238 of order  $\sqrt{T}$ , Gradient Bandit provides good results when the MLE is applied to all  $T$  data  
 239 points. The truncation at  $\sqrt{T} \simeq 20$  required in the theoretical results does not seem necessary in  
 240 practice, and actually looks suboptimal for Gradient Bandit. Figures 2b and 2c show that the  
 241 hold-out is almost systematically outperformed by the penalized MLE. Both struggle to identify the  
 242 significantly more complex model OnePerItem, preferring simpler alternatives.

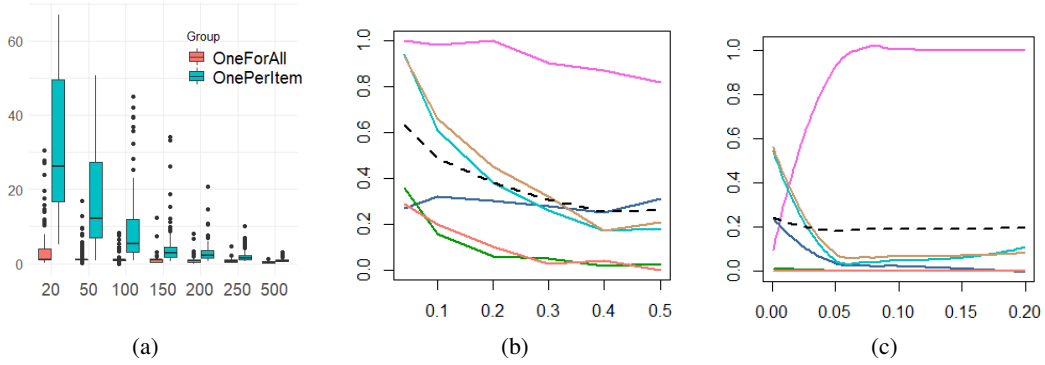


Figure 2: Errors of the procedures as a function of the tuning parameters. In 2a, average of the  $|\hat{\theta}_C - \theta_C|/\theta_C$  over all cells  $C$  in model OneForAll and OnePerItem for the data generated respectively by the same models, where  $\hat{\theta}_C$  is the MLE with likelihood truncated at  $N$  (in abscissa). In 2b and 2c, percentage of mismatch between  $\hat{m}$  and the simulated model over 100 simulations. The colors for each model are the ones given in Figure 3 whereas the average error on the models in the dash line. In 2b, for the hold-out estimator as a function of  $N/T$ . In 2c, for the penalized MLE with  $\text{pen}(m) = c \log(T)^2 D_m/T$ , as a function of  $c$ .

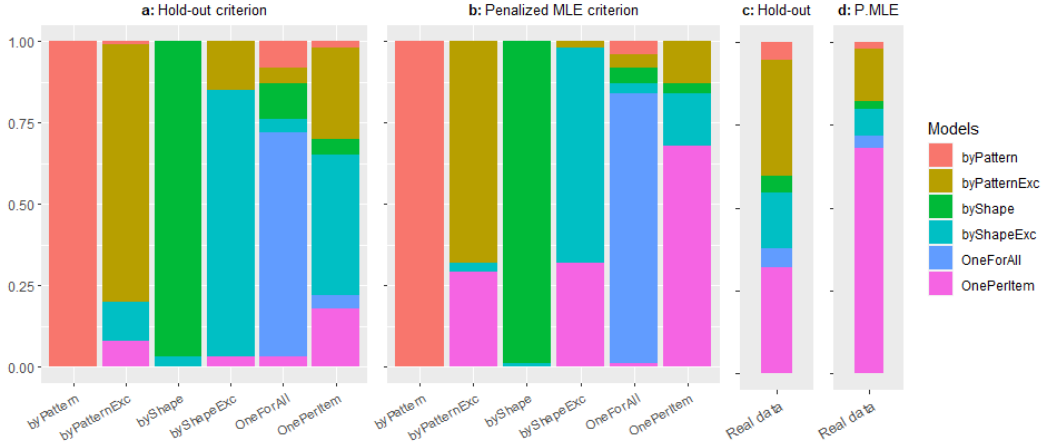


Figure 3: Distribution of the model choices. In **a**, hold-out with  $N = 250$  over 100 simulations. In **b**, penalized MLE with  $c = 0.012$  over 100 simulations. In **c**, hold-out on the data recorded in [34]–176 participants. In **d**, penalized MLE on the same experimental data.

243 Given these results on simulated data, we use  $N = 250$  for the hold-out and  $c = 0.012$  for the  
 244 penalized MLE. The proportion of mismatches for each model are reported in Figure 3a for the  
 245 hold-out and 3b for the penalized MLE. Both methods manage to recover the true model with less  
 246 than 35% of mistakes, except for the model OnePerItem, for which only the penalized MLE is able  
 247 to achieve a successful match more than 60% of the time. The models that are confused the most  
 248 are the ones that are able to correctly learn the categorization, that is ByPatternExc, ByShapeExc  
 249 and OnePerItem.



250 **On real data.** The data have been collected for [34]<sup>1</sup> and we focus only on the learning data. We  
 251 use only the 176 participants that needed at least  $T = 100$  trials. In Figure 3d, we see that most  
 252 of participants are attributed one of the 3 models able to learn. The most frequent is `OnePerItem`  
 253 (about 70% for the penalized MLE) and this percentage is larger than the one obtained on sim-  
 254 ulation, probably meaning that a significant proportion of the participants do not see the division  
 255 along the dimensions `Shape` or `Pattern`. It would be interesting for further study to see if this is  
 256 linked to the presentation order of the objects, as it has been proved for `Alcove` and `Component Cue`  
 257 in [34].

## 258 6 Metalearning

259 By looking at the experiment above, it is hard to believe that learners start directly with a model like  
 260 `ByPatternExc`. It is more likely that they start with a model like `ByPattern` and realize that there  
 261 are too many exceptions, so that they progressively end up with `ByPatternExp`. One way to model  
 262 this progressive switch from one strategy to the other is to use bandits with expert advice. In this  
 263 framework, there is a finite set  $E$  of randomized policies called experts,  $(\xi_{j,t}(\cdot))_{t \in [T]}$ , probabilities  
 264 over the set of actions  $[k]$ , that are modeling the different strategies the learner might have. No  
 265 assumptions are made here on the way experts compute their randomized predictions: they might  
 266 be the result of contextual bandits like `ByPattern` or more generally any kind of computations that  
 267 depend on the learner’s past choices. `Exp4` (see Algorithm 4) is an adaptation of `Exp3` to this case  
 268 (see [28] for regret convergence and variants such as `Exp4.P` [8]).

---

### Algorithm 4 `Exp4` [11]

---

**Inputs:**  $T$  (Sample size),  $\theta \in [r, R]$  (Parameter),  $K$  (Number of actions),  $E$  (Set of experts).

**Initialization:**  $q_{E,1}^\theta$  uniform distribution over the experts  $E$ .

**for**  $t = 1, 2, \dots$  **do**

Receive experts advice  $a \mapsto \xi_{j,t}(a)$  probability distribution over  $[K]$  for all  $j$ .

Draw an action  $A_t \sim \pi_{E,t}^\theta(\cdot) = \sum_{j \in E} q_{E,t}^\theta(j) \xi_{j,t}(\cdot)$  and receive a reward  $g_{A_t,t} \in [0, 1]$ .

Update for all  $j \in E$ ,

$$q_{E,t+1}^\theta(j) = \frac{\exp\left(-\frac{\theta}{\sqrt{T}} \sum_{s=1}^t \hat{y}_{j,s}^\theta\right)}{\sum_{i \in E} \exp\left(-\frac{\theta}{\sqrt{T}} \sum_{s=1}^t \hat{y}_{i,s}^\theta\right)} \quad \text{with} \quad \hat{y}_{i,s}^\theta = \sum_{a \in [K]} \xi_{i,t}(a) \frac{g_{a,s}}{\pi_{E,s}^\theta(a)} \mathbf{1}_{A_s=a}$$


---

269 In this setting, a model  $m$  is defined by a finite set  $E_m$  representing the different experts/strategies  
 270 the learner is learning from. Since there is only one parameter by model (namely  $\theta \in [r, R]$ ), the  
 271 penalty plays no role, nor the calibration of  $c$ . So there is no need for hold-out and one can prove  
 272 that the model with the smallest log-likelihood on the first  $T_\varepsilon \sim \sqrt{T}$  time steps satisfies an oracle  
 273 inequality if  $\mathcal{M}$  is finite, as well as  $|F| := \max_{m \in \mathcal{M}} |E_m|$ . Details are given in the Appendix B.  
 274 This shows that one can select the set  $E_m$  of strategies which is the closest to reality among the sets  
 275 of strategies that are put in competition.

276 *Limitations.* The only limitation with this approach is that we need at first to know the eventual  
 277 parameters of each strategy. Again we could split the data in a hold-out fashion to make the injection  
 278 of estimated parameters possible. However, it would be then nearly impossible to correctly estimate  
 279 the parameters of strategies that are not used at the beginning of the learning. We refer to [9] for  
 280 other methods in meta-learning for cognition.

## 281 References

- 282 [1] S. Arlot. Minimal penalties and the slope heuristics: a survey. *J. SFdS*, 160(3):158–168, 2019.
- 283 [2] S. Arlot and A. Celisse. A survey of cross-validation procedures for model selection. *Statistics*  
 284 *Surveys*, 4:40–79, 2010.

---

<sup>1</sup>We refer the reader to [34] for precise description of the task as well as the ethics agreement.

- 285 [3] S. Arlot and M. Lerasle. Choice of  $V$  for  $V$ -fold cross-validation in least-squares density  
286 estimation. *J. Mach. Learn. Res.*, 17:Paper No. 208, 50, 2016.
- 287 [4] S. Arora and P. Doshi. A survey of inverse reinforcement learning: Challenges, methods and  
288 progress. *Artificial Intelligence*, 297:103500, 2021.
- 289 [5] J. Aubert, L. Lehéricy, and P. Reynaud-Bouret. On the convergence of the MLE as an estimator  
290 of the learning rate in the exp3 algorithm. In *Proceedings of the 40th International Confer-*  
291 *ence on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages  
292 1244–1275. PMLR, 7 2023.
- 293 [6] J. Aubert, L. Lehéricy, and P. Reynaud-Bouret. General oracle inequalities for a penalized  
294 log-likelihood criterion based on non-stationary data. Working paper or preprint, 5 2024.
- 295 [7] J.-P. Baudry, C. Maugis, and B. Michel. Slope heuristics: overview and implementation. *Statis-*  
296 *tics and Computing*, 22:455–470, 2012.
- 297 [8] A. Beygelzimer, J. Langford, L. Li, L. Reyzin, and R. Schapire. Contextual bandit algorithms  
298 with supervised learning guarantees. In *Proceedings of the Fourteenth International Confer-*  
299 *ence on Artificial Intelligence and Statistics*, pages 19–26. JMLR Workshop and Conference  
300 Proceedings, 2011.
- 301 [9] M. Binz, I. Dasgupta, A. K. Jagadish, M. Botvinick, J. X. Wang, and E. Schulz. Meta-learned  
302 models of cognition. *Behavioral and Brain Sciences*, pages 1–38, 2023.
- 303 [10] D. Bouneffouf and I. Rish. A survey on practical applications of multi-armed and contextual  
304 bandits. *arXiv preprint arXiv:1904.10040*, 2019.
- 305 [11] S. Bubeck and N. Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed  
306 bandit problems. *CoRR*, abs/1204.5721, 2012.
- 307 [12] J. R. Busemeyer and A. Diederich. *Cognitive modeling*. Sage, 2010.
- 308 [13] G. Castellán. Density estimation via exponential model selection. *IEEE Trans. Inform. Theory*,  
309 49(8):2052–2060, 2003.
- 310 [14] D. S. Clark. Short proof of a discrete Gronwall inequality. *Discrete Applied Mathematics*,  
311 16(3):279–281, 1987.
- 312 [15] A. G. Collins and A. Shenhav. Advances in modeling learning and decision-making in neuro-  
313 science. *Neuropsychopharmacology*, 47(1):104–118, 2022.
- 314 [16] N. D. Daw. Trial-by-trial data analysis using computational models. *Decision making, affect,*  
315 *and learning: Attention and performance XXIII*, 23(1), 2011.
- 316 [17] M. Dimakopoulou, Z. Zhou, S. Athey, and G. Imbens. Estimation considerations in contextual  
317 bandits. *arXiv preprint arXiv:1711.07077*, 2017.
- 318 [18] D. J. Foster, A. Krishnamurthy, and H. Luo. Model selection for contextual bandits. *Advances*  
319 *in Neural Information Processing Systems*, 32, 2019.
- 320 [19] B. Gao and L. Pavel. On the properties of the softmax function with application in game theory  
321 and reinforcement learning. *arXiv preprint arXiv:1704.00805*, 2017.
- 322 [20] M. A. Gluck and G. H. Bower. Evaluating an adaptive network model of human learning.  
323 *Journal of memory and Language*, 27(2):166–195, 1988.
- 324 [21] N. R. Hansen, P. Reynaud-Bouret, and V. Rivoirard. Lasso and probabilistic inequalities for  
325 multivariate point processes. *Bernoulli*, 21(1), 2 2015.
- 326 [22] C. Houdré and P. Reynaud-Bouret. Exponential inequalities for  $u$ -statistics of order two with  
327 constants. *Stochastic Inequalities and Applications. Progress in Probability*, 56, 01 2002.
- 328 [23] A. Hussein, M. M. Gaber, E. Elyan, and C. Jayne. Imitation learning: A survey of learning  
329 methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.

- 330 [24] A. Hüyük, D. Jarrett, and M. van der Schaar. Inverse contextual bandits: Learning how behavior evolves over time. In *International Conference on Machine Learning*, pages 9506–9524. PMLR, 2022.
- 331  
332
- 333 [25] A. James, P. Reynaud-Bouret, G. Mezzadri, F. Sargolini, I. Bethus, and A. Muzy. Strategy inference during learning via cognitive activity-based credit assignment models. *Scientific reports*, 13(1):9408, 2023.
- 334  
335
- 336 [26] J. K. Kruschke. Alcove: An exemplar-based connectionist model of category learning. In *Connectionist Psychology*, pages 107–138. Psychology Press, 2020.
- 337
- 338 [27] A. S. Lan and R. G. Baraniuk. A contextual bandits framework for personalized learning action selection. In *EDM*, pages 424–429, 2016.
- 339
- 340 [28] T. Lattimore and C. Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- 341 [29] P. Massart. *Concentration inequalities and model selection: Ecole d’Eté de Probabilités de Saint-Flour XXXIII-2003*. Springer, 2007.
- 342
- 343 [30] D. L. Medin and M. M. Schaffer. Context theory of classification learning. *Psychological review*, 85(3):207, 1978.
- 344
- 345 [31] J. Mei, Z. Zhong, B. Dai, A. Agarwal, C. Szepesvari, and D. Schuurmans. Stochastic gradient succeeds for bandits. In *International Conference on Machine Learning*, pages 24325–24360. PMLR, 2023.
- 346  
347
- 348 [32] G. Mezzadri. *Statistical inference for categorization models and presentation order*. Phd thesis, Université Côte d’Azur, LJAD, France, 2020.
- 349
- 350 [33] G. Mezzadri, T. Laloë, F. Mathy, and P. Reynaud-Bouret. Hold-out strategy for selecting learning models: application to categorization subjected to presentation orders. *Journal of Mathematical Psychology*, 109:102691, 2022.
- 351  
352
- 353 [34] G. Mezzadri, P. Reynaud-Bouret, T. Laloë, and F. Mathy. Investigating interactions between types of order in categorization. *Scientific Reports*, 12(1):21625, 2022.
- 354
- 355 [35] G. Neu. Explore no more: Improved high-probability regret bounds for non-stochastic bandits. *Advances in Neural Information Processing Systems*, 28, 2015.
- 356
- 357 [36] J. Opsomer, Y. Wang, and Y. Yang. Nonparametric Regression with Correlated Errors. *Statistical Science*, 16(2):134 – 153, 2001.
- 358
- 359 [37] A. Pacchiano, M. Phan, Y. Abbasi Yadkori, A. Rao, J. Zimmert, T. Lattimore, and C. Szepesvari. Model selection in contextual stochastic bandit problems. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 10328–10337. Curran Associates, Inc., 2020.
- 360  
361  
362
- 363 [38] G. Ramponi, G. Drappo, and M. Restelli. Inverse reinforcement learning from a gradient-based learner. *Advances in Neural Information Processing Systems*, 33:2458–2468, 2020.
- 364
- 365 [39] H. Robbins and S. Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- 366
- 367 [40] E. Schulz, E. Konstantinidis, and M. Speekenbrink. Learning and decisions in contextual multi-armed bandit tasks. In *CogSci*, 2015.
- 368
- 369 [41] E. Schulz, E. Konstantinidis, and M. Speekenbrink. Putting bandits into context: How function learning supports decision making. *Journal of experimental psychology: learning, memory, and cognition*, 44(6):927, 2018.
- 370  
371
- 372 [42] R. C. Wilson and A. G. Collins. Ten simple rules for the computational modeling of behavioral data. *Elife*, 8:e49547, 2019.
- 373

## 374 **NeurIPS Paper Checklist**

### 375 **1. Claims**

376 Question: Do the main claims made in the abstract and introduction accurately reflect the  
377 paper's contributions and scope?

378 Answer: [\[Yes\]](#)

379 Justification: We describe our contributions in the abstract and introduction. They are  
380 mainly theoretical in nature, with oracle inequalities for two model selection procedures  
381 (hold-out and penalized maximum likelihood) in Sections 2 and 3. We also show that the  
382 assumptions of these results are satisfied for some common reinforcement learning models  
383 in Section 4 and 6 and finally illustrate the procedure on synthetic and real data in Section 5.

384 Guidelines:

- 385 • The answer NA means that the abstract and introduction do not include the claims  
386 made in the paper.
- 387 • The abstract and/or introduction should clearly state the claims made, including the  
388 contributions made in the paper and important assumptions and limitations. A No or  
389 NA answer to this question will not be perceived well by the reviewers.
- 390 • The claims made should match theoretical and experimental results, and reflect how  
391 much the results can be expected to generalize to other settings.
- 392 • It is fine to include aspirational goals as motivation as long as it is clear that these  
393 goals are not attained by the paper.

### 394 **2. Limitations**

395 Question: Does the paper discuss the limitations of the work performed by the authors?

396 Answer: [\[Yes\]](#)

397 Justification: After each theoretical result, we wrote a paragraph named limitations to dis-  
398 cuss the issues we thought of with respect to each result. Note that since our focus is on the  
399 theoretical properties of the methods, with the experimental sections serving as illustration  
400 of these properties, we do not care about the computational complexity of the algorithms  
401 considered and do not discuss them beyond their execution time discussed in Appendix A.

402 Guidelines:

- 403 • The answer NA means that the paper has no limitation while the answer No means  
404 that the paper has limitations, but those are not discussed in the paper.
- 405 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 406 • The paper should point out any strong assumptions and how robust the results are to  
407 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
408 model well-specification, asymptotic approximations only holding locally). The au-  
409 thors should reflect on how these assumptions might be violated in practice and what  
410 the implications would be.
- 411 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
412 only tested on a few datasets or with a few runs. In general, empirical results often  
413 depend on implicit assumptions, which should be articulated.
- 414 • The authors should reflect on the factors that influence the performance of the ap-  
415 proach. For example, a facial recognition algorithm may perform poorly when image  
416 resolution is low or images are taken in low lighting. Or a speech-to-text system might  
417 not be used reliably to provide closed captions for online lectures because it fails to  
418 handle technical jargon.
- 419 • The authors should discuss the computational efficiency of the proposed algorithms  
420 and how they scale with dataset size.

- 421 • If applicable, the authors should discuss possible limitations of their approach to ad-  
422 dress problems of privacy and fairness.
- 423 • While the authors might fear that complete honesty about limitations might be used by  
424 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
425 limitations that aren't acknowledged in the paper. The authors should use their best  
426 judgment and recognize that individual actions in favor of transparency play an impor-  
427 tant role in developing norms that preserve the integrity of the community. Reviewers  
428 will be specifically instructed to not penalize honesty concerning limitations.

### 429 3. Theory Assumptions and Proofs

430 Question: For each theoretical result, does the paper provide the full set of assumptions and  
431 a complete (and correct) proof?

432 Answer: [Yes]

433 Justification: We stated precisely all our theoretical results with assumptions that are clearly  
434 referenced. Only the last section about metalearning is written slightly more informally due  
435 to lack of space but formally written in the supplementary material. We give intuition about  
436 how each assumption is used. The proofs are given in the supplementary material.

437 Guidelines:

- 438 • The answer NA means that the paper does not include theoretical results.
- 439 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
440 referenced.
- 441 • All assumptions should be clearly stated or referenced in the statement of any theo-  
442 rems.
- 443 • The proofs can either appear in the main paper or the supplemental material, but if  
444 they appear in the supplemental material, the authors are encouraged to provide a  
445 short proof sketch to provide intuition.
- 446 • Inversely, any informal proof provided in the core of the paper should be comple-  
447 mented by formal proofs provided in appendix or supplemental material.
- 448 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 449 4. Experimental Result Reproducibility

450 Question: Does the paper fully disclose all the information needed to reproduce the main  
451 experimental results of the paper to the extent that it affects the main claims and/or conclu-  
452 sions of the paper (regardless of whether the code and data are provided or not)?

453 Answer: [Yes]

454 Justification: The simulation of data according to each model is precisely explained in  
455 Appendix A, and the code to produce it is given. The codes to compute both estimators  
456 (hold-out and penalized MLE) are given and commented, also in the supplementary mate-  
457 rial. The real data are taken from a published paper [34] and we asked the authors of [34] to  
458 provide us with these data. We don't think it is possible to make these data public because  
459 the ethics agreement that has been signed by the authors of [34], prior to data collection,  
460 might not include this possibility. However, this application to real data is mainly to prove  
461 that this can be done in practice. Since there is not a truth to be compared to in these data  
462 and since even cross validation is hard to perform, this real dataset cannot be used as a  
463 benchmark to compare methods anyway.

464 Guidelines:

- 465 • The answer NA means that the paper does not include experiments.
- 466 • If the paper includes experiments, a No answer to this question will not be perceived  
467 well by the reviewers: Making the paper reproducible is important, regardless of  
468 whether the code and data are provided or not.

- 469 • If the contribution is a dataset and/or model, the authors should describe the steps  
470 taken to make their results reproducible or verifiable.
- 471 • Depending on the contribution, reproducibility can be accomplished in various ways.  
472 For example, if the contribution is a novel architecture, describing the architecture  
473 fully might suffice, or if the contribution is a specific model and empirical evaluation,  
474 it may be necessary to either make it possible for others to replicate the model with  
475 the same dataset, or provide access to the model. In general, releasing code and data  
476 is often one good way to accomplish this, but reproducibility can also be provided via  
477 detailed instructions for how to replicate the results, access to a hosted model (e.g., in  
478 the case of a large language model), releasing of a model checkpoint, or other means  
479 that are appropriate to the research performed.
- 480 • While NeurIPS does not require releasing code, the conference does require all sub-  
481 missions to provide some reasonable avenue for reproducibility, which may depend  
482 on the nature of the contribution. For example
  - 483 (a) If the contribution is primarily a new algorithm, the paper should make it clear  
484 how to reproduce that algorithm.
  - 485 (b) If the contribution is primarily a new model architecture, the paper should describe  
486 the architecture clearly and fully.
  - 487 (c) If the contribution is a new model (e.g., a large language model), then there should  
488 either be a way to access this model for reproducing the results or a way to re-  
489 produce the model (e.g., with an open-source dataset or instructions for how to  
490 construct the dataset).
  - 491 (d) We recognize that reproducibility may be tricky in some cases, in which case au-  
492 thors are welcome to describe the particular way they provide for reproducibility.  
493 In the case of closed-source models, it may be that access to the model is limited in  
494 some way (e.g., to registered users), but it should be possible for other researchers  
495 to have some path to reproducing or verifying the results.

## 496 5. Open access to data and code

497 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
498 tions to faithfully reproduce the main experimental results, as described in supplemental  
499 material?

500 Answer: [Yes]

501 Justification: All codes to generate synthetic data and to perform penalized MLE and hold-  
502 out are provided, so that all the numerical part on the calibration of both methods can  
503 be faithfully reproduced. Only the real dataset, as explained before, cannot be given for  
504 reproduction (Figure 3).

505 Guidelines:

- 506 • The answer NA means that paper does not include experiments requiring code.
- 507 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
508 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 509 • While we encourage the release of code and data, we understand that this might not  
510 be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
511 including code, unless this is central to the contribution (e.g., for a new open-source  
512 benchmark).
- 513 • The instructions should contain the exact command and environment needed to run to  
514 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
515 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 516 • The authors should provide instructions on data access and preparation, including how  
517 to access the raw data, preprocessed data, intermediate data, and generated data, etc.

- 518 • The authors should provide scripts to reproduce all experimental results for the new  
519 proposed method and baselines. If only a subset of experiments are reproducible, they  
520 should state which ones are omitted from the script and why.
- 521 • At submission time, to preserve anonymity, the authors should release anonymized  
522 versions (if applicable).
- 523 • Providing as much information as possible in supplemental material (appended to the  
524 paper) is recommended, but including URLs to data and code is permitted.

## 525 6. Experimental Setting/Details

526 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
527 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
528 results?

529 Answer: [Yes]

530 Justification: The whole purpose of our numerical study in Section 5 is to explain the choice  
531 of hyperparameters (such as the splitting in the hold-out of the calibration of the constant  $c$   
532 in the penalty).

533 Guidelines:

- 534 • The answer NA means that the paper does not include experiments.
- 535 • The experimental setting should be presented in the core of the paper to a level of  
536 detail that is necessary to appreciate the results and make sense of them.
- 537 • The full details can be provided either with the code, in appendix, or as supplemental  
538 material.

## 539 7. Experiment Statistical Significance

540 Question: Does the paper report error bars suitably and correctly defined or other appropri-  
541 ate information about the statistical significance of the experiments?

542 Answer: [Yes]

543 Justification: We have run our simulations on 600 independent simulated learners and we  
544 show with a boxplot (Figure 2a) and mismatch proportion graphs (Figure 2b, 2c and 3)  
545 the proportion of erroneous selections. This cannot be done on real data, since each real  
546 participant to the categorization task is unique.

547 Guidelines:

- 548 • The answer NA means that the paper does not include experiments.
- 549 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
550 dence intervals, or statistical significance tests, at least for the experiments that support  
551 the main claims of the paper.
- 552 • The factors of variability that the error bars are capturing should be clearly stated (for  
553 example, train/test split, initialization, random drawing of some parameter, or overall  
554 run with given experimental conditions).
- 555 • The method for calculating the error bars should be explained (closed form formula,  
556 call to a library function, bootstrap, etc.)
- 557 • The assumptions made should be given (e.g., Normally distributed errors).
- 558 • It should be clear whether the error bar is the standard deviation or the standard error  
559 of the mean.
- 560 • It is OK to report 1-sigma error bars, but one should state it. The authors should prefer-  
561 ably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of  
562 Normality of errors is not verified.

- 563
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 564
- 565
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
- 566
- 567

## 568 8. Experiments Compute Resources

569 Question: For each experiment, does the paper provide sufficient information on the com-  
570 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
571 the experiments?

572 Answer: [Yes]

573 Justification: It is not central in our analysis so it is just mentioned in the supplementary  
574 material in Appendix A.

575 Guidelines:

- The answer NA means that the paper does not include experiments.
  - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
  - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
  - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 576
- 577
- 578
- 579
- 580
- 581
- 582
- 583

## 584 9. Code Of Ethics

585 Question: Does the research conducted in the paper conform, in every respect, with the  
586 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

587 Answer: [Yes]

588 Justification: Up to the real data that are used in this paper, there is absolutely nothing that  
589 would be in violation to the Code of Ethics. For the real data that are used, they are human  
590 categorization data. They have been recorded for another publication and just transmitted to  
591 us. The experimental procedure was approved by the local ethics committee of the authors.  
592 We do not want to share these data publicly since we do not want to breach the Privacy rule  
593 of the Code of Ethics.

594 Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
  - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
  - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 595
- 596
- 597
- 598
- 599

## 600 10. Broader Impacts

601 Question: Does the paper discuss both potential positive societal impacts and negative  
602 societal impacts of the work performed?

603 Answer: [NA]

604 Justification: This work is theoretical. The methods that are validated theoretically here  
605 have already been in use in practice for a long time (see for instance the rules to follow  
606 for cognitive modeling in [42]) and so the expected societal impact of the present work is  
607 negligible.

608 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- 609



- 610 • If the authors answer NA or No, they should explain why their work has no societal  
611 impact or why the paper does not address societal impact.
- 612 • Examples of negative societal impacts include potential malicious or unintended uses  
613 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
614 (e.g., deployment of technologies that could make decisions that unfairly impact spe-  
615 cific groups), privacy considerations, and security considerations.
- 616 • The conference expects that many papers will be foundational research and not tied  
617 to particular applications, let alone deployments. However, if there is a direct path to  
618 any negative applications, the authors should point it out. For example, it is legitimate  
619 to point out that an improvement in the quality of generative models could be used to  
620 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
621 that a generic algorithm for optimizing neural networks could enable people to train  
622 models that generate Deepfakes faster.
- 623 • The authors should consider possible harms that could arise when the technology is  
624 being used as intended and functioning correctly, harms that could arise when the  
625 technology is being used as intended but gives incorrect results, and harms following  
626 from (intentional or unintentional) misuse of the technology.
- 627 • If there are negative societal impacts, the authors could also discuss possible mitiga-  
628 tion strategies (e.g., gated release of models, providing defenses in addition to attacks,  
629 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
630 feedback over time, improving the efficiency and accessibility of ML).

## 631 11. Safeguards

632 Question: Does the paper describe safeguards that have been put in place for responsible  
633 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
634 image generators, or scraped datasets)?

635 Answer: [NA]

636 Justification: We do not think this applies to our research.

637 Guidelines:

- 638 • The answer NA means that the paper poses no such risks.
- 639 • Released models that have a high risk for misuse or dual-use should be released with  
640 necessary safeguards to allow for controlled use of the model, for example by re-  
641 quiring that users adhere to usage guidelines or restrictions to access the model or  
642 implementing safety filters.
- 643 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
644 should describe how they avoided releasing unsafe images.
- 645 • We recognize that providing effective safeguards is challenging, and many papers do  
646 not require this, but we encourage authors to take this into account and make a best  
647 faith effort.

## 648 12. Licenses for existing assets

649 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
650 the paper, properly credited and are the license and terms of use explicitly mentioned and  
651 properly respected?

652 Answer: [Yes]

653 Justification: We clearly stated that the real data come from [34]. The code has been  
654 developed by us solely, using classical packages in R that are clearly mentioned in the code  
655 and supplementary material.

656 Guidelines:

- 657 • The answer NA means that the paper does not use existing assets.
- 658 • The authors should cite the original paper that produced the code package or dataset.

- 659 • The authors should state which version of the asset is used and, if possible, include a  
660 URL.
- 661 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 662 • For scraped data from a particular source (e.g., website), the copyright and terms of  
663 service of that source should be provided.
- 664 • If assets are released, the license, copyright information, and terms of use in the pack-  
665 age should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has  
666 curated licenses for some datasets. Their licensing guide can help determine the li-  
667 cense of a dataset.
- 668 • For existing datasets that are re-packaged, both the original license and the license of  
669 the derived asset (if it has changed) should be provided.
- 670 • If this information is not available online, the authors are encouraged to reach out to  
671 the asset’s creators.

### 672 13. New Assets

673 Question: Are new assets introduced in the paper well documented and is the documenta-  
674 tion provided alongside the assets?

675 Answer: [NA]

676 Justification: We do not provide new packages associated to our results.

677 Guidelines:

- 678 • The answer NA means that the paper does not release new assets.
- 679 • Researchers should communicate the details of the dataset/code/model as part of their  
680 submissions via structured templates. This includes details about training, license,  
681 limitations, etc.
- 682 • The paper should discuss whether and how consent was obtained from people whose  
683 asset is used.
- 684 • At submission time, remember to anonymize your assets (if applicable). You can  
685 either create an anonymized URL or include an anonymized zip file.

### 686 14. Crowdsourcing and Research with Human Subjects

687 Question: For crowdsourcing experiments and research with human subjects, does the pa-  
688 per include the full text of instructions given to participants and screenshots, if applicable,  
689 as well as details about compensation (if any)?

690 Answer: [No]

691 Justification: We did not collect data for the present article but used data collected for [34],  
692 a work that is already published. In this article, all the details are given and we do not think  
693 it makes sense to reproduce it here for our illustration. We only kept the main description  
694 of the task so that the readers can understand what was done.

695 Guidelines:

- 696 • The answer NA means that the paper does not involve crowdsourcing nor research  
697 with human subjects.
- 698 • Including this information in the supplemental material is fine, but if the main contri-  
699 bution of the paper involves human subjects, then as much detail as possible should  
700 be included in the main paper.
- 701 • According to the NeurIPS Code of Ethics, workers involved in data collection, cura-  
702 tion, or other labor should be paid at least the minimum wage in the country of the  
703 data collector.

### 704 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 705 Subjects

706 Question: Does the paper describe potential risks incurred by study participants, whether  
707 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
708 approvals (or an equivalent approval/review based on the requirements of your country or  
709 institution) were obtained?

710 Answer: [No]

711 Justification: The data collection done for [34] had the approval of the local ethic committee  
712 as mentioned in their article. Here we do not feel necessary to reproduce this here but rather  
713 point towards [34] for additional information about the task and its ethic agreement.

714 Guidelines:

- 715 • The answer NA means that the paper does not involve crowdsourcing nor research  
716 with human subjects.
- 717 • Depending on the country in which research is conducted, IRB approval (or equiva-  
718 lent) may be required for any human subjects research. If you obtained IRB approval,  
719 you should clearly state this in the paper.
- 720 • We recognize that the procedures for this may vary significantly between institutions  
721 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
722 guidelines for their institution.
- 723 • For initial submissions, do not include any information that would break anonymity  
724 (if applicable), such as the institution conducting the review.

## 725 A Code and data description

### 726 A.1 Details on numerical illustrations

727 In this section, we give details on the numerical illustrations of Section 5. The images were obtained  
728 using the `ggplot2` package of R. Two types of analyses were conducted, on synthetic data and on  
729 real data.

730 **On synthetic data.** The simulations of the synthetic data helped us calibrate the tuning parameters  
731 choices for the hold-out and the penalized log-likelihood procedure. In Section 2, the parameter  $N$   
732 must be calibrated for choosing the correct training data sample size. In Section 3, as said in the  
733 *Limitations*, since the constant  $c$  in the penalty term is not known a priori, it must be calibrated as  
734 well. To do this, we follow the guidelines of [42]. The procedure is as follows.

- 735 1) Sample size:  $T = 500$ . It is of the same order of magnitude as real data.
- 736 2) Objects generation: periodic sequence of the nine objects repeated through the  $T$  trials. We  
737 generate a sequence of objects following the same structure as in [34]. Due to the periodic  
738 pattern, each object is therefore seen roughly the same number of times for all time  $t$ .
- 739 3) Actions generation: for each model in Table 1, we generated 100 sequences of actions  
740 called synthetic agents with respect to the procedure given in Algorithm 1 with `Gradient`  
741 `Bandit` as `CellBandit`. The parameters  $\theta_C$  we used were the same for each model and  
742 the same for each cell, equal to  $0.03 \times \sqrt{T}$ , except for the `OnePerItem` model where we  
743 changed slightly the values of the parameter in each cell to make the model identifiable.  
744 For  $m = \text{OnePerItem}$ , we took  $\theta^m = ((0.03/10 + k \times 0.007) \times \sqrt{T})_{k \in \{0, \dots, 8\}}$  following  
745 the same order of presentation of the sequence of objects defined earlier.
- 746 4) Parameters estimation: we then fitted each of the six models on all the synthetic agents  
747 generated data, and we estimated the associated parameters using (MLE) and the pack-  
748 age `DEoptim` in R with range  $(0, 1)$  for the parameters  $\theta_C/\sqrt{T}$  and with the default  
749 parameters and a maxiter value equal to 20. We then computed the log-likelihood as-  
750 sociated to the estimated parameters. We did this for the likelihood stopped at time  
751  $N \in \{25, 50, 100, 150, 200, 250, T\}$ .
  - 752 - With such data, we were able to plot Figure 2a and Figure 2b with the hold-out crite-  
753 rion defined in Section 2. In Figure 2a, we computed the average error made in each  
754 cell by the model fitting of the same model that generated the data. For the Figure 2b,  
755 we simply counted the number of times each model verified the hold-out criterion for  
756 all the synthetic agent and for each model that generated the data.
  - 757 - With the log-likelihood stopped at time  $T$  for the estimated parameters, we were able  
758 to plot Figure 2c according to the penalized log-likelihood criterion defined in (5).  
759 In the same way we counted the number of times each model satisfied the penalized  
760 log-likelihood criterion for all the synthetic agent and for each model that generated  
761 the data.
- 762 5) Choice of the parameters  $N$  and  $c$  for the real data: Given the results of Figure 2b and  
763 Figure 2c, we chose to use  $N$  to be equal to half of the data length and  $c = 0.012$  to  
764 account for a reasonable error for model `OnePerItem`, even if in average  $c = 0.04$  gives  
765 better results. With this data, we were able plot the two first chart of Figure 3.

766 **On real data.** For the real experimental data, here is the process we followed.

- 767 1) Sample size: dependent on each agent, the average data sample size is 300.
- 768 2) Objects and Actions: we collected for each agent their objects sequence and associated  
769 choices.
- 770 3) For each agent, we fitted the 6 models and estimated the parameters associated to each  
771 model. To perform hold-out and penalized log-likelihood model selection, we used the  
772 parameters  $N$  and  $c$  chosen thanks to the synthetic data. With this data, we were able to  
773 plot Figure 3.

## 774 A.2 About the code and the data

775 In this section, we give explanations about the code and data (e.g. computation time, link between  
776 code and data). All the data, code and images used are provided in the zip file associated to sub-  
777 mission, called `ContextualBanditsCode`. We run all the simulations in R and used the following  
778 packages: `DEoptim`, `crayon`, `magrittr`, `dplyr`, `tidyr`, `ggplot2`, `gridExtra`.

779 For the sample size we chose, all the simulations can run on a PC in a reasonable time of execution  
780 (detailed hereafter). Overall, computing the different data and running the code took approximately  
781 6 hours excluding the time needed for the real data. The biggest file is 373 kilobytes. The PC  
782 we used was a Gigabyte - AORUS 15G XC, with processor: Intel(R) Core(TM) i7-10870H CPU  
783 2.20GHz, 2208 MHz, 8 cores, 16 logical processors.

784 **On real data.** As mentioned earlier, we could not provide the experimental data used in [34],  
785 since they have already been published in another paper and we do not want to break the ethic  
786 agreement. We can only provide the results and estimated data resulting from the experimental  
787 data. Note however that the procedures to obtain the following RData files are the same as for  
788 the synthetic data which we detail later. The three RData files on the real data are `realdata mle`,  
789 `realdata_holdout_trainingset`, `realdata_holdout_testingset`.

- 790 • `realdata mle` is a list of estimators and associated log-likelihood for each model and each  
791 agent.
- 792 • `realdata_holdout_trainingset` is a list of estimators and associated log-likelihood on  
793 the first half of the sample for each model and each agent.
- 794 • `realdata_holdout_testingset` is a list of log-likelihood on the testing part  
795 of the sample for each agent and each model with parameters estimated in  
796 `realdata_holdout_trainingset`.

797 **On synthetic data.** All the synthetic data obtained in the other files can be computed by running  
798 the code `ContextualbanditsCodebis`. The code is commented and starts with a list of functions  
799 which are necessary to run the different procedures. In the code, we explain how the different  
800 procedures lead to the following list of files. We have commented with # the parts of the code that  
801 would modify the files so that running the code now would give the same images as the ones used in  
802 the article. If one wants to generate new data, one should uncomment these lines of code. However,  
803 we advise the reader that some of the procedures take a certain time, and would recommend not to  
804 do so. We detail hereafter the content of the different csv and RData files and the time it took to run  
805 them.

- 806 • To begin with, we generate a csv file called `databis_500.csv` of 500 trials and associated  
807 list of objects in the file `synthetic_data`.
- 808 • In the same `synthetic_data` file we create the different model files and within each of  
809 them generate 100 csv files of actions, rewards, and objects according to the procedure  
810 described in A.1. This procedure takes around 5 minutes. Then, we begin to compute the  
811 MLE for each of the synthetic data csv file.
- 812 • `Datalikelihood100agents6modeletabis500horizon` is a nested list of estimators,  
813 associated log-likelihood stopped at time  $T$  for each model fitted to the data of all the  
814 synthetic agents. Computing these data took approximately 2 hours .
- 815 • `holdoutbis100agents6models_horizon_20` is the same nested list of estimators but  
816 computed on a log-likelihood stopped at time  $N = 20$ . Computing these data took approx-  
817 imately 10 minutes .
- 818 • `holdoutbis100agents6models_horizon_50` is the same nested list of estimators but  
819 computed on a log-likelihood stopped at time  $N = 50$ . Computing these data took approx-  
820 imately 20 minutes .
- 821 • `holdoutbis100agents6models_horizon_100` is the same nested list of estimators but  
822 computed on a log-likelihood stopped at time  $N = 100$ . Computing these data took ap-  
823 proximately 30 minutes .

- 824 • holdoutbis100agents6models\_horizon\_150 is the same nested list of estimators but  
825 computed on a log-likelihood stopped at time  $N = 150$ . Computing these data took ap-  
826 proximately 40 minutes .
- 827 • holdoutbis100agents6models\_horizon\_200 is the same nested list of estimators but  
828 computed on a log-likelihood stopped at time  $N = 200$ . Computing these data took ap-  
829 proximately 50 minutes .
- 830 • holdoutbis100agents6models\_horizon\_250 is the same nested list of estimators but  
831 computed on a log-likelihood stopped at time  $N = 250$ . Computing these data took ap-  
832 proximately 1 hour .
- 833 • alldataholdoutbis is a nested list of errors on estimation for the training data and log-  
834 likelihood function on the testing data for all synthetic agents, all models and all training  
835 data sample size  $N \in \{20, 50, 100, 150, 200, 250\}$ . Computing these data took approxi-  
836 mately 10 minutes .

## 837 B Assumptions for metalearning

838 Since we work in a more general setting and not simply with contexts, we assume that we observe  
839 a process  $(A_t)_{1 \leq t \leq T}$  adapted to a general filtration  $(\mathcal{F}_t)_{1 \leq t \leq T}$  where for all  $t \in [T]$ ,  $A_t \in [K]$ . In  
840 particular, for all  $t \in [T]$ ,  $\mathcal{F}_t$  is generated by the past actions  $(A_1, \dots, A_t)$  and any other additional  
841 variable which might be observed or not – such as a context at time  $t + 1$  for instance. We write, for  
842 all  $a \in [K]$ , and all  $t \in [T]$

$$p_t^*(a) = \mathbb{P}(A_t = a | \mathcal{F}_{t-1})$$

843 the true conditional distribution we wish to estimate.

844 Additionally, we consider the family of models  $\{(\pi_{E_m, t}^{\theta^m})_{t \in [T]}, m \in \mathcal{M}\}$  where  $\mathcal{M}$  is a finite set,  
845  $\theta^m \in [r, R]$ , and for all  $m \in \mathcal{M}$ ,  $(\pi_{E_m, t}^{\theta^m})_{t \in [T]}$  is the sequence of mixtures of probability distribu-  
846 tions over actions defined recursively in Algorithm 4 for the finite set  $E_m$ . Each model  $m$  is thus  
847 defined by a set of experts  $(\xi_{j, t}(\cdot))_{j \in E_m, t \in [T]}$  where for all  $m \in \mathcal{M}$ ,  $t \in [T]$ ,  $\xi_{j, t}$  can be any  
848 probability distribution over arms  $[K]$  as long as it is measurable with respect to  $\mathcal{F}_{t-1}$ .

849 Let  $|F| := \max_{m \in \mathcal{M}} |E_m|$ . The goal is to select the set  $E_m$  of policies – that we see as learning  
850 strategies – with which the agent learns to learn. This approach is again based on partial log-  
851 likelihood  $\ell_T(\pi_{E_m}^{\theta^m})$  of the observations  $(A_1, \dots, A_T)$  defined by

$$\ell_T(\pi_{E_m}^{\theta^m}) = \sum_{t=1}^T \log \left( \pi_{E_m}^{\theta^m}(A_t) \right). \quad (7)$$

852 To achieve a model selection result, we need the following assumption on the policies given by the  
853 experts.

854 **Assumption 3.** There exists  $\rho > 0$ , such that almost surely, for all  $m \in \mathcal{M}$ , for all  $t \in [T]$  and all  
855  $i \in [K]$ ,  $\sum_{j \in E_m} \xi_{j, t}(i) \geq \rho$ .

856 Then, with Assumption 3, we can deduce a result similar to Propositions 4 and 5 because of the very  
857 structure of Algorithm 4 which mimics Exp3.

858 **Proposition 6.** Assume Assumption 3 holds. Let  $\rho$  be the associated constant. Let  $\varepsilon \in (0, \rho/|F|)$ ,  
859 and let

$$T_\varepsilon = \left\lfloor \left( \frac{1}{|F|} - \frac{\varepsilon}{\rho} \right) \frac{\sqrt{T}}{R} \right\rfloor \wedge T \quad \text{and} \quad L_\varepsilon = \frac{1}{R\varepsilon^2} \exp \left( \frac{1}{\varepsilon^2} \right).$$

860 Then, for all  $t \in [T_\varepsilon]$ , for all  $m \in \mathcal{M}$ ,  $\theta^m, \delta^m \in [r, R]$ , for all  $k \in [K]$ ,

$$\pi_{E_m, t}^{\theta^m}(k) \geq \varepsilon \quad \text{and} \quad \sup_{k \in [K]} \left| \log \left( \frac{\pi_{E_m, t}^{\theta^m}(k)}{\pi_{E_m, t}^{\delta^m}(k)} \right) \right| \leq L_\varepsilon |\theta^m - \delta^m|.$$

861 Finally, we still assume that the true distribution is bounded away from 0 (as in (3)).

862 **Assumption 4.** Assume that Assumption 3 holds. Let  $\varepsilon$  and  $T_\varepsilon$  be the constants of Proposition 6.  
 863 Assume that

$$\forall t \leq T_\varepsilon, \forall a \in [K], p_t^*(a) \geq \varepsilon.$$

864 Assumptions 3 and 4 allow us to verify Assumptions 1 and 2 of [6]. As for Section 3, it is thus  
 865 possible to put into competition different sets of experts. Let  $A_\varepsilon = L_\varepsilon(R - r) + 2 \log(\varepsilon^{-1})$ . Since  
 866 all the models have the same dimension, there is no penalty term to account for. So the term  $\Sigma_\varepsilon$  in  
 867 Theorem 3 becomes  $\log(A_\varepsilon)|\mathcal{M}|e^{-1}$ . The result of [6, Corollary 2] states that there exist constants  
 868  $c, c'$  such that, for all  $\kappa \in (0, 1]$ ,

$$\begin{aligned} \frac{1-\kappa}{T_\varepsilon} \sum_{t=1}^{T_\varepsilon} \mathbb{E} \left[ D_{\text{KL}} \left( p_t^*, \pi_{E_{\hat{m},t}}^{\hat{\theta}^{\hat{m}}} \right) \right] &\leq \inf_{m \in \mathcal{M}} \left( (1+\kappa) \inf_{\theta \in \Theta^{D_m}} \frac{1}{T_\varepsilon} \sum_{t=1}^{T_\varepsilon} \mathbb{E} \left[ D_{\text{KL}} \left( p_t^*, \pi_{E_{m,t}}^{\theta^m} \right) \right] \right) \\ &+ \frac{c}{\kappa} A_\varepsilon^2 \log(\varepsilon^{-1})^{3/2} \log(T_\varepsilon A_\varepsilon)^2 \frac{1}{T_\varepsilon} \\ &+ \frac{18e^{-1}c'}{\kappa} A_\varepsilon \log(A_\varepsilon) |\mathcal{M}| \log(\varepsilon^{-1})^{3/2} \log(T_\varepsilon A_\varepsilon)^2 \frac{\log(T_\varepsilon)}{T_\varepsilon}. \end{aligned}$$

## 869 C Proofs

### 870 C.1 Proof of Section 2

871 *Proof of Theorem 1.* For any  $m \in \mathcal{M}$ , and  $k \in [K]$ , let

$$\begin{cases} g_m(k, X_t) = -\frac{1}{2} \log \left( \frac{p_t^m(k, X_t)}{p_t^*(k, X_t)} \right) \\ f_m(k, X_t) = -\log \left( \frac{p_t^*(k, X_t) + p_t^m(k, X_t)}{2p_t^*(k, X_t)} \right). \end{cases}$$

872 For any function  $h : [K] \times \mathcal{X} \rightarrow \mathbb{R}$ , let

$$\begin{cases} \tilde{\nu}_T(h) = \frac{1}{T-N+1} \sum_{t=N}^T \sum_{k=1}^K h(k, X_t) (\mathbf{1}_{A_t=k} - p_t^*(k, X_t)), \\ \tilde{P}_T(h) = \frac{1}{T-N+1} \sum_{t=N}^T \sum_{k=1}^K h(k, X_t) \mathbf{1}_{A_t=k}, \\ \tilde{C}_T(h) = \frac{1}{T-N+1} \sum_{t=N}^T \sum_{k=1}^K h(k, X_t) p_t^*(k, X_t). \end{cases}$$

873 From the definition of  $\hat{m}$ ,

$$\tilde{P}_T(g_{\hat{m}}) \leq \tilde{P}_T(g_m).$$

874 Since,  $f_{\hat{m}}(k, X_t) \leq g_{\hat{m}}(k, X_t)$  by concavity of  $\log$ , it holds that

$$\tilde{\nu}_T(f_{\hat{m}}) + \tilde{C}_T(f_{\hat{m}}) = \tilde{P}_T(f_{\hat{m}}) \leq \tilde{P}_T(g_m) = \tilde{\nu}_T(g_m) + \tilde{C}_T(g_m).$$

875 That is

$$\tilde{\nu}_T(f_{\hat{m}} - f_m) + \tilde{C}_T(f_{\hat{m}}) \leq \tilde{\nu}_T(g_m - f_m) + \tilde{C}_T(g_m).$$

876 Let  $U_m = \tilde{\nu}_T(g_m - f_m)$ , then

$$\tilde{C}_T(f_{\hat{m}}) \leq \tilde{C}_T(g_m) - \tilde{\nu}_T(f_{\hat{m}} - f_m) + U_m. \quad (8)$$

877 Note that  $U_m$  is centered. For  $m' \in \mathcal{M}$ , let  $M_N^{m'} = 0$ , and for  $t \geq N+1$ , let

$$M_t^{m'} = - \sum_{s=N}^{t-1} \sum_{k=1}^K (f_{m'}(k, X_s) - f_m(k, X_s)) (\mathbf{1}_{A_s=k} - p_s^*(k, X_s)).$$

878 For all  $t \geq N$ , let  $\mathcal{H}_t = \sigma(X_t, \mathcal{F}_{t-1})$ . Then,  $(M_t^{m'})_{t \geq N}$  is an  $(\mathcal{H}_t)_{t \geq N}$ -martingale. For  $\ell \geq 2$ , let  
879  $B_N^\ell = 0$ , and for  $t \geq N + 1$ , let

$$B_t^\ell := \sum_{s=N}^{t-1} \mathbb{E} \left[ \left( M_{s+1}^{m'} - M_s^{m'} \right)^\ell \middle| \mathcal{H}_s \right].$$

880 For  $t \in \{N, \dots, T-1\}$ , note that

$$|M_{t+1}^{m'} - M_t^{m'}| \leq 2 \sum_{k=1}^K |f_{m'}(k, X_t) - f_m(k, X_t)| \frac{\mathbf{1}_{A_t=k} + p_t^*(k, X_t)}{2},$$

881 so that, by convexity of  $x \mapsto x^\ell$  on  $[0, +\infty)$ ,

$$|M_{t+1}^{m'} - M_t^{m'}|^\ell \leq 2^\ell \sum_{k=1}^K |f_{m'}(k, X_t) - f_m(k, X_t)|^\ell \frac{\mathbf{1}_{A_t=k} + p_t^*(k, X_t)}{2}$$

882 Thus,

$$\begin{aligned} B_t^\ell &= \sum_{s=N}^{t-1} \mathbb{E} \left[ \left( M_{s+1}^{m'} - M_s^{m'} \right)^\ell \middle| \mathcal{H}_s \right] \\ &\leq 2^\ell \sum_{s=N}^{t-1} \sum_{k=1}^K |f_{m'}(k, X_s) - f_m(k, X_s)|^\ell p_s^*(k, X_s) \\ &= 2^\ell \sum_{s=N}^{t-1} \sum_{k=1}^K \left| \log \left( \frac{p_s^*(k, X_s) + p_s^{m'}(k, X_s)}{p_s^*(k, X_s) + p_s^m(k, X_s)} \right) \right|^\ell p_s^*(k, X_s). \end{aligned} \quad (9)$$

883 We now need the following Lemma to continue.

884 **Lemma 7.** [29, Lemma 7.26] For all  $\ell \geq 2$  and all  $x > 0$ ,

$$\frac{|\log(x)|^\ell}{\ell!} \leq \frac{9}{64} \left( x - \frac{1}{x} \right)^2.$$

885 *Proof.* The complete Lemma and proof of the Lemma can be found in [29].  $\square$

886 Applying Lemma 7 to  $x = \sqrt{\frac{p_s^*(k, X_s) + p_s^{m'}(k, X_s)}{p_s^*(k, X_s) + p_s^m(k, X_s)}}$  leads to, for all  $k \in [K]$ ,

$$\begin{aligned} &\left| \log \left( \frac{p_s^*(k, X_s) + p_s^{m'}(k, X_s)}{p_s^*(k, X_s) + p_s^m(k, X_s)} \right) \right|^\ell \\ &\leq \frac{9}{64} 2^\ell \ell! \frac{\left( p_s^m(k, X_s) - p_s^{m'}(k, X_s) \right)^2}{\left( p_s^m(k, X_s) + p_s^*(k, X_s) \right) \left( p_s^*(k, X_s) + p_s^{m'}(k, X_s) \right)}. \end{aligned}$$

887 Plugging this in Equation (9) leads to

$$|B_t^\ell| \leq \frac{9}{4} 2^{2(\ell-2)} \ell! \sum_{s=N}^{t-1} \sum_{k=1}^K \frac{\left( p_s^m(k, X_s) - p_s^{m'}(k, X_s) \right)^2 p_s^*(k, X_s)}{\left( p_s^m(k, X_s) + p_s^*(k, X_s) \right) \left( p_s^*(k, X_s) + p_s^{m'}(k, X_s) \right)}.$$

888 For all  $x, y, z \geq 0$ ,

$$(\sqrt{x} + \sqrt{y})^2 z \leq (z + y)(z + x),$$

889 therefore, with  $z = p_s^*(k, X_s)$ ,  $x = p_s^m(k, X_s)$  and  $y = p_s^{m'}(k, X_s)$ ,

$$|B_t^\ell| \leq \frac{9}{4} 4^{\ell-2} \ell! \sum_{s=N}^{t-1} \sum_{k=1}^K \left( \sqrt{p_s^m(k, X_s)} - \sqrt{p_s^{m'}(k, X_s)} \right)^2 \leq \frac{1}{2} 4^{\ell-2} \ell! V_t^{m'}, \quad (10)$$



890 where

$$\begin{aligned} V_t^{m'} &:= \frac{9}{2} \sum_{s=N}^{t-1} \sum_{k=1}^K \left( \sqrt{p_s^m(k, X_s)} - \sqrt{p_s^{m'}(k, X_s)} \right)^2 \\ &= 9 \sum_{s=N}^{t-1} H \left( p_s^m(\cdot, X_s), p_s^{m'}(\cdot, X_s) \right)^2 \end{aligned} \quad (11)$$

891 where  $H$  is the Hellinger distance between the two probability distributions  $p_s^m(\cdot, X_s)$  and  
892  $p_s^{m'}(\cdot, X_s)$ . Lemma 3.3 of [22] gives that for all  $\lambda > 0$ ,

$$\mathcal{E}_t = \exp \left( \lambda M_t^{m'} - \sum_{\ell \geq 2} \frac{\lambda^\ell}{\ell!} B_t^\ell \right)$$

893 is a supermartingale and that in particular  $\mathbb{E}(\mathcal{E}_{T+1}) \leq 1$ . By Equation (10), for  $\lambda \in (0, 1/4)$ ,

$$\sum_{\ell \geq 2} \frac{\lambda^\ell}{\ell!} B_t^\ell \leq \frac{\lambda^2}{2} \sum_{\ell \geq 2} (4\lambda)^{\ell-2} V_t^{m'} = \frac{\lambda^2}{2(1-4\lambda)} V_t^{m'}.$$

894 Let  $\Psi(\lambda) = \frac{\lambda^2}{2(1-4\lambda)}$  for  $\lambda \in (0, 1/4)$ . Then,

$$\mathbb{E} \left[ e^{\lambda M_{T+1}^{m'} - \Psi(\lambda) V_{T+1}^{m'}} \middle| \mathcal{H}_N \right] \leq 1.$$

895 By Markov's inequality, for all  $x \geq 0$  and  $\lambda \in (0, 1/4)$ ,

$$\mathbb{P} \left( M_{T+1}^{m'} \geq V_{T+1}^{m'} \frac{\Psi(\lambda)}{\lambda} + \frac{x}{\lambda} \middle| \mathcal{H}_N \right) \leq e^{-x}. \quad (12)$$

896 Therefore, for all  $x, u \geq 0$  and  $\lambda \in (0, 1/4)$ ,

$$\mathbb{P} \left( M_{T+1}^{m'} \geq \frac{\Psi(\lambda)}{\lambda} u + \frac{x}{\lambda} \quad \text{and} \quad V_{T+1}^{m'} \leq u \middle| \mathcal{H}_N \right) \leq e^{-x}.$$

897 To choose the optimal  $\lambda$ , we use Lemma 2 from [21].

898 **Lemma 8.** [21, Lemma 2] Let  $a, b$  and  $x$  be positive constants and let us consider on  $(0, 1/b)$ ,

$$g(\xi) = \frac{a\xi}{1-b\xi} + \frac{x}{\xi}.$$

899 Then  $\min_{\xi \in (0, 1/b)} g(\xi) = 2\sqrt{ax} + bx$  and the minimum is achieved in  $\xi(a, b, x) = \frac{\sqrt{x}}{\sqrt{xb} + \sqrt{a}}$ .

900 For  $a = \frac{u}{2}$  and  $b = 4$ , Lemma 8 shows that for all  $x, u \geq 0$ ,

$$\mathbb{P} \left( M_{T+1}^{m'} \geq \sqrt{2ux} + 4x \quad \text{and} \quad V_{T+1}^{m'} \leq u \middle| \mathcal{H}_N \right) \leq e^{-x}.$$

901 Let us use a peeling argument similar to [21]:

902 **Lemma 9.** Let  $X, V$  be real-valued random variables and  $\alpha, b, v, w$  be positive numbers such that  
903  $V \in [w, v]$  a.s. and such that for all  $x \geq 0$  and  $u \in [w, v]$ ,

$$\mathbb{P}(X \geq \sqrt{ux} + bx \quad \text{and} \quad (1+\alpha)^{-1}u \leq V \leq u) \leq e^{-x},$$

904 then for any  $x \geq 0$ ,

$$\mathbb{P}(X \geq \sqrt{(1+\alpha)Vx} + bx) \leq \left( 1 + \frac{\log(v/w)}{\log(1+\alpha)} \right) e^{-x}.$$

905 *Proof.* Let  $v_0 = w$ ,  $v_{d+1} = (1 + \alpha)v_d$ , and  $D$  the smallest integer such that  $v_D \geq v$ . For all  $d \in [D]$   
 906 and  $x \geq 0$ ,

$$\mathbb{P}(X \geq \sqrt{v_d x} + bx \quad \text{and} \quad v_{d-1} \leq V \leq v_d) \leq e^{-x}.$$

907 In particular, since  $V \geq v_{d-1} = (1 + \alpha)^{-1}v_d$  on this event,

$$\mathbb{P}(X \geq \sqrt{(1 + \alpha)Vx} + bx \quad \text{and} \quad v_{d-1} \leq V \leq v_d) \leq e^{-x}.$$

908 Taking the union bound,

$$\mathbb{P}(X \geq \sqrt{(1 + \alpha)Vx} + bx) \leq De^{-x},$$

909 and by definition  $D \leq \frac{\log(v/w)}{\log(1 + \alpha)} + 1$ . □

910 We may apply Lemma 9 to  $X = M_{T+1}^{m'}$  and  $b = 4$ . Since  $V_{T+1}^{m'}$  does not have an obvious lower  
 911 bound, we consider  $V = 2(V_{T+1}^{m'} + \beta)$  for some  $\beta > 0$  to be chosen later. We may therefore take  
 912  $w = 2\beta$ . For the upper bound  $v$  on  $V$ , by (11), since the Hellinger distance is upper bounded by 1,  
 913 we may take  $v = 2(\beta + 9(T - N + 1))$ . With these choices, for any  $\beta, \alpha, x > 0$ ,

$$\mathbb{P}\left(M_{T+1}^{m'} \geq \sqrt{2(1 + \alpha)(V_{T+1}^{m'} + \beta)x} + 4x \mid \mathcal{H}_N\right) \leq \left(\frac{\log\left(\frac{9(T-N+1)}{\beta} + 1\right)}{\log(1 + \alpha)} + 1\right) e^{-x}.$$

914 For  $\alpha = \sqrt{2}$ ,

$$\mathbb{P}\left(M_{T+1}^{m'} \geq \sqrt{5(V_{T+1}^{m'} + \beta)x} + 4x \mid \mathcal{H}_N\right) \leq \left(2 \log\left(\frac{9(T - N + 1)}{\beta} + 1\right) + 1\right) e^{-x}. \quad (13)$$

915 By definition of  $V_{T+1}^{m'}$  and the triangle inequality,

$$\begin{aligned} V_{T+1}^{m'} &= 9 \sum_{s=N}^T H(p_s^m(\cdot, X_s), p_s^{m'}(\cdot, X_s))^2 \\ &\leq 18 \sum_{s=N}^T \left( H(p_s^*(\cdot, X_s), p_s^m(\cdot, X_s))^2 + H(p_s^*(\cdot, X_s), p_s^{m'}(\cdot, X_s))^2 \right). \end{aligned} \quad (14)$$

916 We now use [29, Lemma 7.23] giving a connection between the Kullback-Leibler divergence  $D_{\text{KL}}$   
 917 and the Hellinger distance  $H$ .

918 **Lemma 10.** [29, Lemma 7.23] *Let  $P$  and  $Q$  be some probability measures. Then,*

$$D_{\text{KL}}\left(P, \frac{P + Q}{2}\right) \geq (2 \log 2 - 1) H^2(P, Q).$$

919 *Moreover, whenever  $P \ll Q$ ,*

$$2H^2(P, Q) \leq D_{\text{KL}}(P, Q).$$

920 Since  $\frac{18}{2 \log(2) - 1} \leq 48$ ,

$$\begin{aligned} V_{T+1}^{m'} &\leq 48 \sum_{s=N}^T \left( D_{\text{KL}}\left(p_s^*(\cdot, X_s), \frac{p_s^*(\cdot, X_s) + p_s^{m'}(\cdot, X_s)}{2}\right) + \frac{1}{2} D_{\text{KL}}(p_s^*(\cdot, X_s), p_s^m(\cdot, X_s)) \right) \\ &=: 9W_T^{m'}. \end{aligned} \quad (15)$$

921 Let  $\beta = 9(T - N + 1)y^2$ , where  $y > 0$  is to be chosen later. Replacing  $x$  by  $x + \log(|\mathcal{M}|)$  leads to

$$\begin{aligned} \mathbb{P}\left(\frac{M_{T+1}^{m'}}{T - N + 1} \geq 3\sqrt{5\left(\frac{W_T^{m'}}{T - N + 1} + y^2\right)\frac{x + \log(|\mathcal{M}|)}{T - N + 1}} + 4\frac{x + \log(|\mathcal{M}|)}{T - N + 1} \mid \mathcal{H}_N\right) \\ \leq (2 \log(y^{-2} + 1) + 1) e^{-(x + \log(|\mathcal{M}|))}. \end{aligned}$$

922 Let  $\kappa_1 \in (0, 1/(8\sqrt{5})]$ , then, using  $2\sqrt{ab} \leq \kappa_1 a + \kappa_1^{-1} b$  and taking  $y^2 = \frac{x + \log(|\mathcal{M}|)}{(T-N+1)\log 2} \geq \frac{1}{T-N+1}$   
 923 since  $x \geq 0$  and  $|\mathcal{M}| \geq 2$ ,

$$\begin{aligned} \mathbb{P} \left( \frac{M_{T+1}^{m'}}{T-N+1} \geq \frac{3\sqrt{5}}{2} \kappa_1 \frac{W_T^{m'}}{T-N+1} + \left( 4 + \frac{3\sqrt{5}}{2} \left( \frac{\kappa_1}{\log 2} + \kappa_1^{-1} \right) \right) \frac{x + \log(|\mathcal{M}|)}{T-N+1} \middle| \mathcal{H}_N \right) \\ \leq (2 \log(T-N+2) + 1) e^{-(x + \log(|\mathcal{M}|))}. \end{aligned} \quad (16)$$

924 By the union bound on all  $m' \in \mathcal{M}$ , the previous inequality holds with probability at least  $1 -$   
 925  $(2 \log(T-N+2) + 1)e^{-x}$  for all  $m' \in \mathcal{M}$ . It holds in particular for  $\hat{m}$ . Recall with (8) that,

$$\begin{aligned} \frac{1}{T-N+1} \sum_{s=N}^T \text{D}_{\text{KL}} \left( p_s^*(\cdot, X_s), \frac{p_s^*(\cdot, X_s) + p_s^{\hat{m}}(\cdot, X_s)}{2} \right) - U_m \\ \leq \frac{1}{2(T-N+1)} \sum_{s=N}^T \text{D}_{\text{KL}} (p_s^*(\cdot, X_s), p_s^m(\cdot, X_s)) + \frac{M_{T+1}^{\hat{m}}}{T-N+1}. \end{aligned} \quad (17)$$

926 Plugging (15) and (16) in (17) leads to, conditionally on  $\mathcal{H}_N$ , with probability at least  $1 - (2 \log(T -$   
 927  $N + 2) + 1)e^{-x}$

$$\begin{aligned} \frac{(1 - C_{\kappa_1})}{T-N+1} \sum_{s=N}^T \text{D}_{\text{KL}} \left( p_s^*(\cdot, X_s), \frac{p_s^*(\cdot, X_s) + p_s^{\hat{m}}(\cdot, X_s)}{2} \right) - U_m \\ \leq \frac{(1 + C_{\kappa_1})}{T-N+1} \sum_{s=N}^T \frac{1}{2} \text{D}_{\text{KL}} (p_s^*(\cdot, X_s), p_s^m(\cdot, X_s)) + C'_{\kappa_1} \frac{x + \log(|\mathcal{M}|)}{T-N+1}, \end{aligned}$$

928 where

$$\begin{aligned} 929 \quad & \bullet C_{\kappa_1} = 8\sqrt{5}\kappa_1, \\ 930 \quad & \bullet C'_{\kappa_1} = 4 + \frac{3\sqrt{5}}{2} \left( \frac{\kappa_1}{\log 2} + \kappa_1^{-1} \right) \leq 13\kappa_1^{-1} = \frac{104\sqrt{5}}{C_{\kappa_1}}, \end{aligned}$$

931 Integrating on  $x \geq 0$  and noting that  $\mathbb{E}[U_m | \mathcal{H}_N] = 0$  leads to, for all  $m \in \mathcal{M}$ ,

$$\begin{aligned} \frac{(1 - C_{\kappa_1})}{T-N+1} \mathbb{E} \left[ \sum_{s=N}^T \text{D}_{\text{KL}} \left( p_s^*(\cdot, X_s), \frac{p_s^*(\cdot, X_s) + p_s^{\hat{m}}(\cdot, X_s)}{2} \right) \middle| \mathcal{H}_N \right] \\ \leq \frac{(1 + C_{\kappa_1})}{T-N+1} \mathbb{E} \left[ \frac{1}{2} \sum_{s=N}^T \text{D}_{\text{KL}} (p_s^*(\cdot, X_s), p_s^m(\cdot, X_s)) \middle| \mathcal{H}_N \right] \\ + \frac{104\sqrt{5}}{C_{\kappa_1}} \frac{2 \log(T-N+2) + 1 + \log(|\mathcal{M}|)}{T-N+1}. \end{aligned}$$

932 To conclude  $\kappa = \frac{\kappa_1}{8\sqrt{5}}$ , so that  $C_{\kappa_1} = \kappa$ .  $\square$

### 933 C.2 Proof of Section 3

934 *Proof of Proposition 2.* The proof is straightforward with the definition of  $p_{\theta^m, t}^m$  in (1). Let  
 935  $\theta^m, \delta^m \in \Theta^m$ ,  $t \leq T_\varepsilon$ ,  $x \in \mathcal{X}$ , and  $k \in [K]$ .

$$p_{\theta^m, t}^m(k, x) = \sum_{C \in \mathcal{P}_m} \pi_{C, T_t^C}^{\theta^m}(k) \mathbf{1}_{x \in C} \geq \sum_{C \in \mathcal{P}_m} \varepsilon \mathbf{1}_{x \in C} = \varepsilon.$$

936 For the second part of the proof, it holds that, almost surely, for all  $t \leq T_\varepsilon$

$$\begin{aligned} \left| \log \left( \frac{p_{\delta^m, t}^m(k, x)}{p_{\theta^m, t}^m(k, x)} \right) \right| = \sum_{C \in \mathcal{P}_m} \left| \log \left( \frac{\pi_{C, T_t^C}^{\delta^m}(k)}{\pi_{C, T_t^C}^{\theta^m}(k)} \right) \right| \mathbf{1}_{x \in C} \\ \leq L_\varepsilon \sum_{C \in \mathcal{P}_m} \|\delta_C - \theta_C\|_2 \mathbf{1}_{x \in C} \leq L_\varepsilon \sup_{C \in \mathcal{P}_m} \|\delta_C - \theta_C\|_2. \end{aligned}$$

937  $\square$

938 *Proof of Theorem 3.* Our goal is to apply [6].

939 Assumption 1 of [6] is satisfied for  $n = T_\varepsilon$  since with Proposition 2, there exists  $\varepsilon > 0$  such that  
 940 a.s., for all  $t \in [T_\varepsilon]$ , for all  $k \in [K]$ ,  $p_t^*(k, X_t) \in [\varepsilon, 1]$  and for all  $m \in \mathcal{M}$  and all  $\theta^m \in \Theta^{D_m}$ ,  
 941  $p_{\hat{\theta}^m, t}^m(k, X_t) \in [\varepsilon, 1]$ .

942 Assumption 2 of [6] is satisfied since with Proposition 2, there exists a positive constants  $L_\varepsilon$  such  
 943 that a.s., for all  $t \in [T_\varepsilon]$ , for all  $m \in \mathcal{M}$  and all  $\delta^m, \theta^m \in \Theta^{D_m}$ ,

$$\sup_{k \in [K]} \left| \log \left( \frac{p_{\delta^m, t}^m(k, X_t)}{p_{\theta^m, t}^m(k, X_t)} \right) \right| \leq L_\varepsilon \sup_{C \in \mathcal{P}_m} \|\delta_C - \theta_C\|_2$$

944 and by Assumption, for all  $\theta^m, \delta^m \in \Theta^{D_m}$

$$\sup_{C \in \mathcal{P}_m} \|\delta_C - \theta_C\|_2 \leq \sqrt{d}(R - r).$$

945 Note in particular that the Lipschitz constant in Proposition 2 does not depend on  $m$ .

946 Assumption 3 in [6] is always satisfied because the set of actions  $[K]$  is finite.

947 Setting  $A_\varepsilon = L_\varepsilon \sqrt{d}(R - r) + 2 \log(\varepsilon^{-1})$ , Corollary 2 in [6] simply reads as follows. There exist  
 948 positive numerical constants  $C$  and  $C'$  such that the following holds. Assume that

$$\Sigma_\varepsilon = \log(A_\varepsilon) \sum_{m \in \mathcal{M}} e^{-D_m} < +\infty.$$

949 Let  $\kappa \in (0, 1]$ . If for all  $m \in \mathcal{M}$ ,

$$\text{pen}(m) \geq \frac{C}{\kappa} A_\varepsilon^2 \log(\varepsilon^{-1})^{3/2} \log(T_\varepsilon A_\varepsilon)^2 \frac{D_m}{T_\varepsilon},$$

950 then,

$$\begin{aligned} & \frac{1 - \kappa}{T_\varepsilon} \sum_{t=1}^{T_\varepsilon} \mathbb{E} \left[ D_{\text{KL}} \left( p_t^*(\cdot, X_t), p_{\hat{\theta}^m, t}^m(\cdot, X_t) \right) \right] \\ & \leq \inf_{m \in \mathcal{M}} \left( (1 + \kappa) \inf_{\theta \in \Theta^m} \frac{1}{T_\varepsilon} \sum_{t=1}^{T_\varepsilon} \mathbb{E} \left[ D_{\text{KL}} \left( p_t^*(\cdot, X_t), p_{\theta, t}^m(\cdot, X_t) \right) \right] + 2 \text{pen}(m) \right) \\ & \quad + \frac{18C'}{\kappa} A_\varepsilon \Sigma_\varepsilon \log(\varepsilon^{-1})^{3/2} \log(T_\varepsilon A_\varepsilon)^2 \frac{\log(T_\varepsilon)}{T_\varepsilon}. \end{aligned}$$

951

□

### 952 C.3 Proof of Section 4.1

953 *Proof of Proposition 4.* Let  $\theta_C \in \Theta_C$ . Write  $\theta_{C,T} = (\eta_{C,T}, \gamma_{C,T}) = \theta_C / \sqrt{T} \in \Theta$ . To ease the  
 954 notations in the proof, we remove the  $C$  from the notations.  $\theta_C$  becomes  $\theta$  and  $\theta_{C,T}$  becomes  $\theta_T$ .  
 955 In the same way,  $\theta_T = (\eta_T, \gamma_T)$  now.

956 Let  $t \in F_T(C)$ . Then,

$$\begin{aligned} \pi_{C, T_t^C+1}^\theta(k) &= \frac{\pi_{C, T_t^C}^\theta(k) e^{-\eta_T g_{k,t} / (\gamma_T + \pi_{C, T_t^C}^\theta(k))} \mathbf{1}_{A_t=k}}{(1 - \pi_{C, T_t^C}^\theta(k)) + \pi_{C, T_t^C}^\theta(k) e^{-\eta_T g_{k,t} / (\gamma_T + \pi_{C, T_t^C}^\theta(k))}} \\ & \quad + \sum_{\substack{j=1 \\ j \neq k}}^K \frac{\pi_{C, T_t^C}^\theta(j) \mathbf{1}_{A_t=j}}{(1 - \pi_{C, T_t^C}^\theta(j)) + \pi_{C, T_t^C}^\theta(j) e^{-\eta_T g_{j,t} / (\gamma_T + \pi_{C, T_t^C}^\theta(j))}}. \end{aligned}$$

957 For any  $q \in [0, 1]$ , since  $g_{k,t} \in [0, 1]$ ,  $1 - q + qe^{-\eta_T g_{k,t} / (q + \gamma)} \leq 1$ . Therefore,

$$\pi_{C, T_t^C+1}^\theta(k) \geq \pi_{C, T_t^C}^\theta(k) e^{-\eta_T g_{k,t} / (\gamma_T + \pi_{C, T_t^C}^\theta(k))} \mathbf{1}_{A_t=k} + \sum_{\substack{j=1 \\ j \neq k}}^K \pi_{C, T_t^C}^\theta(j) \mathbf{1}_{A_t=j}.$$

958 Since  $e^{-\eta_T g_{k,t}/(\gamma_T + \pi_{C,T_t^C}^\theta(k))} \leq 1$ ,

$$\begin{aligned} & \pi_{C,T_t^C+1}^\theta(k) \\ & \geq \pi_{C,T_t^C}^\theta(k) e^{-\eta_T g_{k,t}/(\gamma_T + \pi_{C,T_t^C}^\theta(k))} \mathbf{1}_{A_t=k} + \sum_{\substack{j=1 \\ j \neq k}}^K \pi_{C,T_t^C}^\theta(k) e^{-\eta_T g_{k,t}/(\gamma_T + \pi_{C,T_t^C}^\theta(k))} \mathbf{1}_{A_t=j} \\ & = \pi_{C,T_t^C}^\theta(k) e^{-\eta_T g_{k,t}/(\gamma_T + \pi_{C,T_t^C}^\theta(k))} \geq \pi_{C,T_t^C}^\theta(k) e^{-\eta_T/(\gamma_T + \pi_{C,T_t^C}^\theta(k))} \end{aligned}$$

959 since  $g_{k,t} \in [0, 1]$ . Then,

$$\pi_{C,T_t^C+1}^\theta(k) \geq \pi_{C,T_t^C}^\theta(k) \left( 1 - \frac{\eta_T g_{k,t}}{\gamma_T + \pi_{C,T_t^C}^\theta(k)} \right) \geq \pi_{C,T_t^C}^\theta(k) - \eta_T.$$

960 Summing for all  $s \in F_t(C)$ , since  $\pi_{k,1}^\theta = \frac{1}{K}$ ,

$$\pi_{C,T_t^C}^\theta(k) \geq \frac{1}{K} - \eta_T T_t^C.$$

961 Note that  $T_t^C \leq t \leq T_\varepsilon$ . Since,  $T_\varepsilon \leq \left[ \left( \frac{1}{K} - \varepsilon \right) \frac{\sqrt{T}}{R} \right]$ , for all  $t \leq T_\varepsilon$  and  $1 \leq k \leq K$ ,

$$\varepsilon \leq \frac{1}{K} - \frac{R}{\sqrt{T}} T_\varepsilon \leq \frac{1}{K} - \eta_T T_\varepsilon \leq \frac{1}{K} - \eta_T t \leq \pi_{C,T_t^C}^\theta(k).$$

962 For the second part of the proof, let  $\theta = (\eta, \gamma)$ ,  $\theta' = (\eta', \gamma') \in \Theta_C$ . For  $t \geq 2$  Let  $h_{j,t}^\theta =$   
963  $\eta_T \sum_{s \in F_t(C)} \hat{g}_{j,s}^\theta$ . Then  $\pi_{C,T_t^C}^\theta = \text{softmax}(h_{\cdot,t}^\theta)$ . The function softmax is 1-Lipschitz with respect  
964 to the  $\|\cdot\|_2$ -norm in  $\mathbb{R}^K$  (see [19] for a proof). Therefore,

$$\|\pi_{C,T_t^C}^\theta - \pi_{C,T_t^C}^{\theta'}\|_2 \leq \|h_{\cdot,t}^\theta - h_{\cdot,t}^{\theta'}\|_2.$$

965 Since  $g_{j,s} \in [0, 1]$ , by the triangle inequality

$$\|\pi_{C,T_t^C}^\theta - \pi_{C,T_t^C}^{\theta'}\|_2 \leq \sum_{s \in F_t(C)} \left\| \left( \frac{\eta_T}{\gamma_T + \pi_{C,T_s^C}^\theta(\cdot)} - \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^{\theta'}(\cdot)} \right) \mathbf{1}_{A_s=\cdot} \right\|_2.$$

966 Again, using the triangle inequality,

$$\begin{aligned} \|\pi_{C,T_t^C}^\theta - \pi_{C,T_t^C}^{\theta'}\|_2 & \leq \sum_{s \in F_t(C)} \left\| \left( \frac{\eta_T}{\gamma_T + \pi_{C,T_s^C}^\theta(\cdot)} - \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^\theta(\cdot)} \right) \mathbf{1}_{A_s=\cdot} \right\|_2 \\ & \quad + \sum_{s \in F_t(C)} \left\| \left( \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^\theta(\cdot)} - \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^{\theta'}(\cdot)} \right) \mathbf{1}_{A_s=\cdot} \right\|_2. \end{aligned} \quad (18)$$

967 For  $1 \geq q \geq \varepsilon$ , let  $f : (x_1, x_2) \in [0, R_T] \times [0, R_T] \mapsto \frac{x_1}{x_2 + q}$  where  $R_T = \frac{R}{\sqrt{T}}$ . The function  $f$  is  
968 continuously differentiable, and

$$\nabla f = \frac{1}{(x_2 + q)^2} \begin{pmatrix} x_2 + q \\ -x_1 \end{pmatrix}.$$

969 The  $\ell^2$ -norm of the gradient can be upper bounded by

$$\|\nabla f\|_2 \leq \frac{1}{\varepsilon^2} \sqrt{R_T^2 + \varepsilon^2} =: c_\varepsilon$$

970 By the mean value theorem, for all  $k \in [K]$

$$\left| \frac{\eta_T}{\gamma_T + \pi_{C,T_s^C}^\theta(k)} - \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^\theta(k)} \right| \leq c_\varepsilon \|\theta_T - \theta'_T\|_2.$$

971 As a result,

$$\begin{aligned} \left\| \left( \frac{\eta_T}{\gamma_T + \pi_{C,T_s^C}^\theta} - \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^{\theta'}} \right) \mathbf{1}_{A_s=} \right\|_2^2 &= \sum_{k=1}^K \left( \frac{\eta_T}{\gamma_T + \pi_{C,T_s^C}^\theta(k)} - \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^{\theta'}(k)} \right)^2 \mathbf{1}_{A_s=k} \\ &\leq c_\varepsilon^2 \|\theta_T - \theta'_T\|_2^2 \sum_{k=1}^K \mathbf{1}_{A_s=k} \\ &= c_\varepsilon^2 \|\theta_T - \theta'_T\|_2^2 \end{aligned}$$

972 Therefore,

$$\begin{aligned} \sum_{s \in F_t(C)} \left\| \left( \frac{\eta_T}{\gamma_T + \pi_{C,T_s^C}^\theta} - \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^{\theta'}} \right) \mathbf{1}_{A_s=} \right\|_2 \\ \leq \sum_{s \in F_t(C)} c_\varepsilon \|\theta_T - \theta'_T\|_2 = T_t^C c_\varepsilon \|\theta_T - \theta'_T\|_2 \leq T_\varepsilon c_\varepsilon \|\theta_T - \theta'_T\|_2. \end{aligned} \quad (19)$$

973 For  $(\eta, \gamma) \in [0, R_T] \times [0, R_T]$ , let  $g : q \in [\varepsilon, 1] \mapsto \frac{\eta}{\gamma + q}$ . The function  $g$  is continuously  
974 differentiable, and

$$0 \leq f'(q) = \frac{\eta}{(\gamma + q)^2} \leq \frac{R_T}{\varepsilon^2}.$$

975 By the mean value theorem, for all  $k \in [K]$ ,

$$\left| \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^{\theta'}(k)} - \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^\theta(k)} \right| \leq \frac{R_T}{\varepsilon^2} \left| \pi_{C,T_s^C}^\theta(k) - \pi_{C,T_s^C}^{\theta'}(k) \right|.$$

976 Therefore,

$$\begin{aligned} \left\| \left( \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^{\theta'}} - \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^\theta} \right) \mathbf{1}_{A_s=} \right\|_2^2 &= \sum_{k=1}^K \left( \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^{\theta'}(k)} - \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^\theta(k)} \right)^2 \mathbf{1}_{A_s=k} \\ &\leq \frac{R_T^2}{\varepsilon^4} \sum_{k=1}^K \left| \pi_{C,T_s^C}^\theta(k) - \pi_{C,T_s^C}^{\theta'}(k) \right|^2 \mathbf{1}_{A_s=k} \\ &\leq \frac{R_T^2}{\varepsilon^4} \left( \left\| \pi_{C,T_s^C}^\theta - \pi_{C,T_s^C}^{\theta'} \right\|_2 \right)^2. \end{aligned}$$

977 Thus,

$$\sum_{s \in F_t(C)} \left\| \left( \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^{\theta'}} - \frac{\eta'_T}{\gamma'_T + \pi_{C,T_s^C}^\theta} \right) \mathbf{1}_{A_s=} \right\|_2 \leq \frac{R_T}{\varepsilon^2} \sum_{s \in F_t(C)} \left\| \pi_{C,T_s^C}^\theta - \pi_{C,T_s^C}^{\theta'} \right\|_2 \quad (20)$$

978 Plugging Equations (19) and (20) in Equation (18)

$$\left\| \pi_{C,T_t^C}^\theta - \pi_{C,T_t^C}^{\theta'} \right\|_2 \leq T_\varepsilon c_\varepsilon \|\theta_T - \theta'_T\|_2 + \frac{R_T}{\varepsilon^2} \sum_{s \in F_t(C)} \left\| \pi_{C,T_s^C}^\theta - \pi_{C,T_s^C}^{\theta'} \right\|_2,$$

979 Using the discrete Gronwall Lemma [14] leads to, for all  $t \leq T_\varepsilon$ ,

$$\begin{aligned} \left\| \pi_{C,T_t^C}^\theta - \pi_{C,T_t^C}^{\theta'} \right\|_2 &\leq T_\varepsilon c_\varepsilon \|\theta_T - \theta'_T\|_2 \prod_{s \in F_t(C)} \left( 1 + \frac{R_T}{\varepsilon^2} \right) \\ &\leq T_\varepsilon c_\varepsilon \|\theta_T - \theta'_T\|_2 \exp \left( \frac{R_T T_\varepsilon}{\varepsilon^2} \right). \end{aligned}$$

980 But, since  $\frac{1}{K} - \varepsilon \leq 1$ ,

$$T_\varepsilon \leq \left( \frac{1}{K} - \varepsilon \right) \frac{\sqrt{T}}{R} \leq \frac{\sqrt{T}}{R}.$$

981 Therefore,  $R_T T_\varepsilon \leq 1$  and for  $t \leq T_\varepsilon$ ,

$$\|\pi_{C, T_t^C}^\theta - \pi_{C, T_t^C}^{\theta'}\|_2 \leq \frac{C_\varepsilon}{R} e^{1/\varepsilon^2} \|\theta - \theta'\|_2.$$

982 To conclude note that  $\log$  is  $1/\varepsilon$ -Lipschitz on  $[\varepsilon, 1]$  and that

$$\sup_{k \in [K]} \left| \log \left( \frac{\pi_{C, T_t^C}^{\delta_C}(k)}{\pi_{C, T_t^C}^\theta(k)} \right) \right| \leq \frac{1}{\varepsilon} \|\pi_{C, T_t^C}^\theta - \pi_{C, T_t^C}^{\theta'}\|_2.$$

983 □

#### 984 C.4 Proof of Section 4.2

985 *Proof of Proposition 5.* We take the same notations as the previous section. The updated probability  
986 can be written

$$\pi_{C, T_t^C+1}^\theta(k) = \frac{\pi_{C, T_t^C}^\theta(k) e^{\theta_T (\mathbf{1}_{A_t=k} - \pi_{C, T_t^C}^\theta(k)) g_{A_t, t}}}{\pi_{C, T_t^C}^\theta(A_t) e^{\theta_T (1 - \pi_{C, T_t^C}^\theta(A_t)) g_{A_t, t}} + \sum_{j \neq A_t} \pi_{C, T_t^C}^\theta(j) e^{-\theta_T \pi_{C, T_t^C}^\theta(j) g_{A_t, t}}}.$$

987 Therefore,

$$\begin{aligned} \pi_{C, T_t^C+1}^\theta(k) &\geq \frac{\pi_{C, T_t^C}^\theta(k) e^{\theta_T (\mathbf{1}_{A_t=k} - \pi_{C, T_t^C}^\theta(k)) g_{A_t, t}}}{\pi_{C, T_t^C}^\theta(A_t) e^{\theta_T (1 - \pi_{C, T_t^C}^\theta(A_t)) g_{A_t, t}} + \sum_{j \neq A_t} \pi_{C, T_t^C}^\theta(j)} \\ &\geq \frac{\pi_{C, T_t^C}^\theta(k) e^{-\theta_T}}{\pi_{C, T_t^C}^\theta(A_t) e^{\theta_T} + 1 - \pi_{C, T_t^C}^\theta(A_t)} \\ &\geq \pi_{C, T_t^C}^\theta(k) e^{-2\theta_T} \end{aligned}$$

988 where

- 989 • the first inequality holds because  $e^{-\theta_T \pi_{C, T_t^C}^\theta(j) g_{A_t, t}} \leq 1$ ,
- 990 • the second inequality holds because  $g_{j, t} \in (0, 1)$ , for  $j \in [K]$ ,
- 991 • the last inequality holds because  $\pi_{C, T_t^C}^\theta(A_t) e^{\theta_T} + 1 - \pi_{C, T_t^C}^\theta(A_t) \leq e^{\theta_T}$ .

992 Thus, for all  $t \leq T_\varepsilon$ ,

$$\pi_{C, T_t^C}^\theta(k) \geq \frac{1}{K} e^{-2\theta_T T_t^C}.$$

993 Since  $T_t^C \leq t$  and since by definition,  $T_\varepsilon \leq \log\left(\sqrt{\frac{1}{K\varepsilon}}\right) \frac{\sqrt{T}}{R}$ , it holds that for all  $t \leq T_\varepsilon$ ,

$$\pi_{C, T_t^C}^\theta(k) \geq \frac{1}{K} e^{-2\theta_T t} \geq \frac{1}{K} e^{-2R_T t} \geq \frac{1}{K} e^{-2\frac{R}{\sqrt{T}} \log\left(\sqrt{\frac{1}{K\varepsilon}}\right) \frac{\sqrt{T}}{R}} \geq \frac{1}{K} e^{-\log\left(\frac{1}{K\varepsilon}\right)} \geq \varepsilon.$$

994 For the second part of the proof, for  $t \geq 2$  and  $j \in [K]$ , let  $h_{j, T_t^C}^\theta = \theta_T \sum_{s \in F_t(C)} \hat{g}_{j, s}^\theta$ . Then  
995  $\pi_{C, T_t^C}^\theta = \text{softmax}(h_{\cdot, t}^\theta)$ . The function  $\text{softmax}$  is 1-Lipschitz with respect to the  $\|\cdot\|_2$ -norm in  $\mathbb{R}^K$   
996 (see [19] for a proof). Therefore,

$$\|\pi_{C, T_t^C}^\delta - \pi_{C, T_t^C}^\theta\|_2 \leq \|h_{\cdot, t}^\delta - h_{\cdot, t}^\theta\|_2.$$

997 Then,

$$\begin{aligned} \|h_{\cdot, t}^\delta - h_{\cdot, t}^\theta\|_2 &\leq |\delta_T - \theta_T| \sum_{s \in F_t(C)} \|\hat{g}_{j, s}^\delta\|_2 + \theta_T \sum_{s \in F_t(C)} g_{A_s, s} \|\pi_{C, T_s^C}^\delta - \pi_{C, T_s^C}^\theta\|_2 \\ &\leq \sqrt{2} T_t^C |\delta_T - \theta_T| + \theta_T \sum_{s \in F_t(C)} \|\pi_{C, T_s^C}^\delta - \pi_{C, T_s^C}^\theta\|_2 \\ &\leq \sqrt{2} T_\varepsilon |\delta_T - \theta_T| + \theta_T \sum_{s \in F_t(C)} \|\pi_{C, T_s^C}^\delta - \pi_{C, T_s^C}^\theta\|_2. \end{aligned}$$

998 where

- 999 • the first inequality holds because of the triangle inequality,  
 1000 • the second inequality holds because for all  $j \in [K]$ ,  $g_{j,s} \in [0, 1]$  and

$$\|\hat{g}_{j,s}^\delta\|_2^2 = (1 - \pi_{C,T_s^C}^\delta(A_s))^2 + \sum_{j \neq A_s} (\pi_{C,T_s^C}^\delta(j))^2 \leq 2,$$

- 1001 • the last inequality holds because  $T_t^C \leq T_\varepsilon$ .

1002 By the discrete Gronwall Lemma [14], for all  $t \leq T_\varepsilon$

$$\|\pi_{C,T_t^C}^\delta - \pi_{C,T_t^C}^\theta\|_2 \leq \sqrt{2}|\delta_T - \theta_T|T_\varepsilon \prod_{s \in F_t(C)} (1 + \theta_T) \leq \sqrt{2}|\delta_T - \theta_T|T_\varepsilon e^{\theta_T T_\varepsilon}.$$

1003 Since  $T_\varepsilon \leq \left\lceil \log \left( \sqrt{\frac{1}{K\varepsilon}} \right) \frac{\sqrt{T}}{R} \right\rceil$ ,  $\theta_T T_\varepsilon \leq R T T_\varepsilon \leq \log \left( \sqrt{\frac{1}{K\varepsilon}} \right)$ , therefore,

$$\|\pi_{C,T_t^C}^\delta - \pi_{C,T_t^C}^\theta\|_2 \leq \frac{\sqrt{2} \log \left( \sqrt{\frac{1}{K\varepsilon}} \right)}{R \sqrt{K\varepsilon}} |\delta - \theta|.$$

1004 Finally,  $\log$  is  $\frac{1}{\varepsilon}$ -Lipschitz on  $[\varepsilon, 1]$ . Thus, for all  $k \in [K]$ ,  $t \leq T_\varepsilon$ ,

$$\left| \log \left( \frac{\pi_{C,T_t^C}^\delta(k)}{\pi_{C,T_t^C}^\theta(k)} \right) \right| \leq \frac{1}{\varepsilon} |\pi_{C,T_t^C}^\delta(k) - \pi_{C,T_t^C}^\theta(k)| \leq \frac{\sqrt{2} \log \left( \sqrt{\frac{1}{K\varepsilon}} \right)}{R\varepsilon \sqrt{K\varepsilon}} |\delta - \theta|.$$

1005

□

## 1006 C.5 Proof of Section B

1007 Let us recall that  $|F| = \max_{m \in \mathcal{M}} |E_m|$ . For this Section, we drop the dependence  $m$  of the model  
 1008 and simply write  $E$  and  $\theta$  generic set of policies and parameter in  $[r, R]$ .

1009 *Proof of Proposition 6.* For any  $\theta \in [r, R]$ , write  $\theta_T = \theta/\sqrt{T}$ . Assume that Assumption 3 holds.

1010 Let's write  $q_{E,t+1}^\theta$  as

$$q_{E,t+1}^\theta(j) = \frac{q_{E,t}^\theta(j) e^{-\theta_T \hat{y}_{j,t}^\theta}}{\sum_{i \in E} q_{E,t}^\theta(i) e^{-\theta_T \hat{y}_{i,t}^\theta}}$$

1011 Since  $q_{E,t}^\theta$  is a probability distribution over the experts,

$$\sum_{i \in E} q_{E,t}^\theta(i) e^{-\theta_T \hat{y}_{i,t}^\theta} \leq \sum_{i \in E} q_{E,t}^\theta(i) = 1.$$

1012 Therefore,  $q_{E,t+1}^\theta(j) \geq q_{E,t}^\theta(j) e^{-\theta_T \hat{y}_{j,t}^\theta}$ . By definition,

$$\hat{y}_{j,t}^\theta = \sum_{k=1}^K \xi_{j,t}(k) \frac{g_{k,t}}{\pi_{E,t}^\theta(k)} \mathbf{1}_{A_t=k} = \xi_{j,t}(A_t) \frac{g_{A_t,t}}{\pi_{E,t}^\theta(A_t)}.$$

1013 Using that  $e^{-x} \geq 1 - x$  for any  $x \geq 0$ , leads to

$$q_{E,t+1}^\theta(j) \geq q_{E,t}^\theta(j) \left( 1 - \theta_T \xi_{j,t}(A_t) \frac{g_{A_t,t}}{\pi_{E,t}^\theta(A_t)} \right).$$

1014 Since  $g_{A_t,t} \in [0, 1]$  and  $q_{E,t}^\theta(j) \xi_{j,t}(A_t) \leq \pi_{E,t}^\theta(A_t)$ ,

$$q_{E,t+1}^\theta(j) \geq q_{E,t}^\theta(j) - \theta_T.$$

1015 Summing for all  $s$  from 1 to  $t$ ,

$$q_{E,t}^\theta(j) \geq \frac{1}{|E|} - \theta_T t.$$



1016 Since

$$T_\varepsilon = \left\lfloor \left( \frac{1}{|F|} - \frac{\varepsilon}{\rho} \right) \frac{\sqrt{T}}{R} \right\rfloor \wedge T \quad \text{and} \quad |F| \geq |E|,$$

1017 it holds that,

$$T_\varepsilon \leq \left\lfloor \left( \frac{1}{|E|} - \frac{\varepsilon}{\rho} \right) \frac{\sqrt{T}}{R} \right\rfloor \wedge T.$$

1018 Therefore, for all  $t \leq T_\varepsilon$ ,

$$q_{E,t}^\theta(j) \geq \frac{1}{|E|} - \frac{R}{\sqrt{T}} t \geq \frac{1}{|E|} - \frac{R}{\sqrt{T}} \left( \frac{1}{|E|} - \frac{\varepsilon}{\rho} \right) \frac{\sqrt{T}}{R} = \frac{\varepsilon}{\rho}.$$

1019 Finally, for all  $t \leq T_\varepsilon$ , for all  $k \in [K]$ ,

$$\pi_{E,t}^\theta(k) = \sum_{j \in E} q_{E,t}^\theta(j) \xi_{j,t}(k) \geq \frac{\varepsilon}{\rho} \sum_{j \in E} \xi_{j,t}(k) = \varepsilon.$$

1020 For the second part of the proof, let  $\eta, \delta \in [r, R]$ , and write  $\eta_T = \eta/\sqrt{T}$  and likewise for  $\delta_T$ .

1021 For  $t \geq 2$ , let  $g_{j,t}^\eta = \eta_T \sum_{s=1}^{t-1} \hat{y}_{j,s}^\eta$ . Then,  $q_{E,t}^\eta = \text{softmax}(g_t^\eta)$  where  $g_t^\eta = (g_{j,t}^\eta)_{j \in E}$ . Since the

1022 function softmax is 1-Lipschitz with respect to the  $\|\cdot\|_2$ -norm in  $\mathbb{R}^{|E|}$ ,

$$\|q_{E,t}^\eta - q_{E,t}^\delta\|_2 \leq \|g_t^\eta - g_t^\delta\|_2.$$

1023 Therefore, by the triangle inequality,

$$\begin{aligned} \|q_{E,t}^\eta - q_{E,t}^\delta\|_2 &\leq \sum_{s=1}^{t-1} \|\eta_T \hat{y}_{j,s}^\eta - \delta_T \hat{y}_{j,s}^\delta\|_2 \\ &\leq \sum_{s=1}^{t-1} (\|\eta_T - \delta_T\| \|\hat{y}_{j,s}^\eta\|_2 + \delta_T \|\hat{y}_{j,s}^\eta - \hat{y}_{j,s}^\delta\|_2). \end{aligned} \quad (21)$$

1024 Since  $\xi_{j,t}$  is a probability distribution,

$$\begin{aligned} \|\hat{y}_{j,t}^\eta\|_2^2 &= \sum_{j \in E} (\hat{y}_{j,t}^\eta)^2 = \sum_{j \in E} \left( \sum_{k=1}^K \xi_{j,t}(k) \frac{g_{k,t}}{\pi_{E,t}^\eta(k)} \mathbf{1}_{A_t=k} \right)^2 \\ &= \sum_{j \in E} \left( \xi_{j,t}(A_t) \frac{g_{A_t,t}}{\pi_{E,t}^\eta(A_t)} \right)^2 \leq |E| \left( \frac{g_{A_t,t}}{\pi_{E,t}^\eta(A_t)} \right)^2. \end{aligned}$$

1025 Since  $g_{A_t,t}^\eta \in [0, 1]$  and  $\pi_{E,t}^\eta(A_t) \geq \varepsilon$  for all  $t \leq T_\varepsilon$ ,

$$\|\hat{y}_{j,t}^\eta\|_2 \leq \frac{\sqrt{|E|}}{\varepsilon}. \quad (22)$$

1026 Similarly,

$$\begin{aligned} \|\hat{y}_{j,s}^\eta - \hat{y}_{j,s}^\delta\|_2^2 &= \sum_{j \in E} \left( \sum_{k=1}^K \xi_{j,t}(k) g_{k,t} \left( \frac{1}{\pi_{E,t}^\eta(k)} - \frac{1}{\pi_{E,t}^\delta(k)} \right) \mathbf{1}_{A_t=k} \right)^2 \\ &\leq \sum_{j \in E} \left( \sum_{k=1}^K \xi_{j,t}(k) \left( \frac{1}{\pi_{E,t}^\eta(k)} - \frac{1}{\pi_{E,t}^\delta(k)} \right) \mathbf{1}_{A_t=k} \right)^2. \end{aligned}$$

1027 Thus, for all  $t \leq T_\varepsilon$ ,

$$\|\hat{y}_{j,t}^\eta - \hat{y}_{j,t}^\delta\|_2^2 \leq \frac{1}{\varepsilon^4} \sum_{j \in E} \left( \sum_{k=1}^K \xi_{j,t}(k) (\pi_{E,t}^\eta(k) - \pi_{E,t}^\delta(k)) \mathbf{1}_{A_t=k} \right)^2.$$

1028 Since  $\xi_{j,t}(k) \leq 1$ ,

$$\|\hat{y}_{j,t}^\eta - \hat{y}_{j,t}^\delta\|_2 \leq \frac{\sqrt{|E|}}{\varepsilon^2} \|\pi_{E,t}^\eta - \pi_{E,t}^\delta\|_2. \quad (23)$$

1029 Injecting (22) and (23) in (21) leads to

$$\|q_{E,t}^\eta - q_{E,t}^\delta\|_2 \leq \frac{\sqrt{|E|}}{\varepsilon} |\eta_T - \delta_T| (t-1) + \delta_T \frac{\sqrt{|E|}}{\varepsilon^2} \sum_{s=1}^{t-1} \|\pi_{E,s}^\eta - \pi_{E,s}^\delta\|_2.$$

1030 On the other hand,

$$\begin{aligned} \|\pi_{E,t}^\eta - \pi_{E,t}^\delta\|_2^2 &= \sum_{k=1}^K (\pi_{k,t}^\eta - \pi_{k,t}^\delta)^2 = \sum_{k=1}^K \left( \sum_{j \in E} \xi_{j,t}(k) (q_{E,t}^\eta(j) - q_{E,t}^\delta(j)) \right)^2 \\ &\leq \sum_{k=1}^K \sum_{j \in E} \xi_{j,t}(k)^2 \sum_{j \in E} (q_{E,t}^\eta(j) - q_{E,t}^\delta(j))^2 \\ &\leq |E| \|q_{E,t}^\eta - q_{E,t}^\delta\|_2^2 \\ &\leq |F| \|q_{E,t}^\eta - q_{E,t}^\delta\|_2^2 \end{aligned}$$

1031 where

- 1032 • the first inequality holds by the Cauchy-Schwarz inequality,
- 1033 • the second inequality holds because  $\xi_{j,t}$  is a probability distribution over the actions set
- 1034  $[K]$ ,
- 1035 • the last inequality holds because  $|E| \leq |F|$

1036 Therefore,

$$\|\pi_{E,t}^\eta - \pi_{E,t}^\delta\|_2 \leq \frac{|F|}{\varepsilon^2} \left( \varepsilon |\eta_T - \delta_T| (t-1) + \delta_T \sum_{s=1}^{t-1} \|\pi_{E,s}^\eta - \pi_{E,s}^\delta\|_2 \right).$$

1037 Using the discrete Gronwall Lemma, for all  $t \leq T_\varepsilon$ ,

$$\|\pi_{E,t}^\eta - \pi_{E,t}^\delta\|_2 \leq \frac{|F|}{\varepsilon} |\eta_T - \delta_T| T_\varepsilon \prod_{s=1}^{t-1} \left( 1 + \frac{|F|}{\varepsilon^2} \delta_T \right) \leq \frac{|F|}{\varepsilon} |\eta_T - \delta_T| T_\varepsilon \exp \left( \frac{|F|}{\varepsilon^2} \delta_T T_\varepsilon \right).$$

1038 If Assumption 3 is satisfied, then since  $\delta_T \leq R_T = \frac{R}{\sqrt{T}}$  and  $T_\varepsilon \leq \left( \frac{1}{|F|} - \frac{\varepsilon}{\rho} \right) \frac{\sqrt{T}}{R}$ ,

$$\|\pi_{E,t}^\eta - \pi_{E,t}^\delta\|_2 \leq \frac{|F|}{R\varepsilon} \left( \frac{1}{|F|} - \frac{\varepsilon}{\rho} \right) \exp \left( \frac{|F|}{\varepsilon^2} \left( \frac{1}{|F|} - \frac{\varepsilon}{\rho} \right) \right) |\eta - \delta|.$$

1039 To conclude note that  $x \rightarrow \ln(x)$  is  $1/\varepsilon$ -Lipschitz on  $[\varepsilon, +\infty)$ . □