

---

# Sequential Decision Making with Expert Demonstrations under Unobserved Heterogeneity

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We study the problem of online sequential decision-making given auxiliary  
2 demonstrations from *experts* who made their decisions based on unobserved  
3 contextual information. These demonstrations can be viewed as solving related but  
4 slightly different tasks than what the learner faces. This setting arises in many  
5 application domains, such as self-driving cars, healthcare, and finance, where ex-  
6 pert demonstrations are made using contextual information, which is not recorded  
7 in the data available to the learning agent. We model the problem as a zero-  
8 shot meta-reinforcement learning setting with an unknown task distribution and a  
9 Bayesian regret minimization objective, where the unobserved tasks are encoded  
10 as parameters with an unknown prior. We propose the Experts-as-Priors algo-  
11 rithm (ExPerior), an empirical Bayes approach that utilizes expert data to estab-  
12 lish an informative prior distribution over the learner’s decision-making problem.  
13 This prior enables the application of any Bayesian approach for online decision-  
14 making, such as posterior sampling. We demonstrate that our strategy surpasses  
15 existing behaviour cloning and online algorithms, as well as online-offline base-  
16 lines for multi-armed bandits, Markov decision processes (MDPs), and partially  
17 observable MDPs, showcasing the broad reach and utility of ExPerior in using  
18 expert demonstrations across different decision-making setups.

## 19 1 Introduction

20 Reinforcement learning (RL) has found success in complex decision-making tasks, spanning areas  
21 such as game playing [1, 2, 3], robotics [4, 5], and aligning with human preferences [6]. However,  
22 RL’s considerable sample inefficiency, necessitating millions of training frames for convergence,  
23 remains a significant challenge. A notable body of work within RL has been dedicated to integrating  
24 expert demonstrations to accelerate the learning process, employing strategies like offline pretraining  
25 [7] and the use of combined offline-online datasets [8, 9]. While these approaches are theoretically  
26 sound and empirically validated [10, 11], they typically presume homogeneity between the offline  
27 demonstrations and online RL tasks. A vital question arises regarding the effectiveness of these  
28 methods when expert data embody heterogeneous tasks, indistinguishable by the learner.

29 An important example of such heterogeneity is in situations where experts operate with additional  
30 information not available to the learner, a scenario previously explored in imitation learning with  
31 unobserved contexts [12, 13, 14, 15]. Existing literature either relies on the availability of experts to  
32 query during training [16, 17, 18, 19] or focuses on the assumptions that enable imitation learning  
33 with unobserved contexts, sidestepping online reward-based interactions [20, 21]. Recent contribu-  
34 tions by Hao et al. [22, 23] suggest the utilization of offline expert data for online RL, albeit without  
35 accounting for unobserved contextual variations.

36 Our work addresses the more general challenge of online sequential decision-making given auxiliary  
37 offline expert data with *unobserved* heterogeneity. We view such demonstrations as solving related

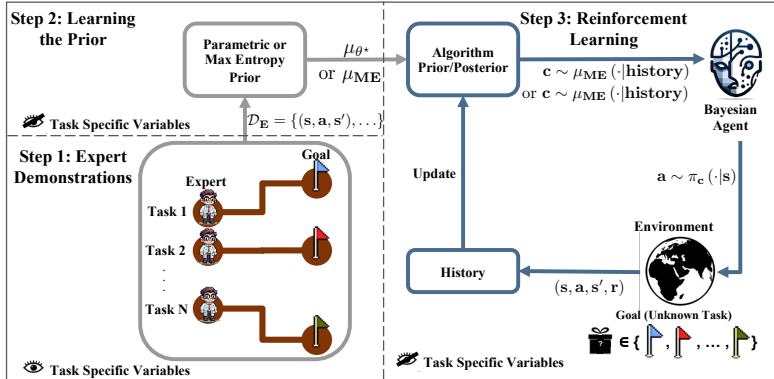


Figure 1: Illustration of ExPerior in a goal-oriented task. Step 1 (Offline): The experts demonstrate their policies for related but different tasks while observing the goal type. Step 2 (Offline): The expert data  $\mathcal{D}_E$  only contains the trajectories states/actions — goal types are not collected. We form a parametric or nonparametric max-entropy prior distribution over tasks using  $\mathcal{D}_E$ . Step 3 (Online): The goal type is unknown but drawn from the same distribution of goals in Step 1. The learner uses the learned prior for posterior sampling.

38 yet distinct tasks from those faced by the learner, where differences remain invisible to the learner.  
 39 For instance, in a personalized education scenario, while a learning agent might access observable  
 40 characteristics like grades or demographics, it might remain oblivious to factors such as learning  
 41 styles, which are visible to an expert teacher and can significantly influence teaching methods. A  
 42 naïve imitation learning algorithm without access to this "private" information will only learn a  
 43 single policy for each observed characteristic [24], leading to sub-optimal decisions. On the other  
 44 hand, a purely online approach will require extensive trial and error to result in meaningful decisions.

45 We integrate offline expert data with online RL, treating the scenario as a zero-shot meta-  
 46 reinforcement learning (meta-RL) problem with an unknown distribution over tasks (unobserved  
 47 factors). Unlike typical meta-RL frameworks where the learner is exposed to multiple tasks during  
 48 training (different students in our example) to learn the underlying task distribution [25, 26].

49 *Contributions:* We define a Bayesian regret minimization objective and consider different tasks as  
 50 parameters under an unknown prior distribution. We use empirical Bayes to derive an informative  
 51 prior over the decision-making task from expert data. We use the learned prior distribution to drive  
 52 exploration in the online RL task, using approaches like posterior sampling [27]. We propose two  
 53 procedures to learn such a prior: (1) a parametric approach that can utilize any existing knowledge  
 54 about the parametric form of the prior distribution, and (2) a nonparametric approach that employs  
 55 the principle of maximum entropy when such prior knowledge does not exist. We call our frame-  
 56 work Experts-as-Priors or ExPerior for short (see Figure 1). ExPerior outperforms existing offline,  
 57 online, and offline-online baselines in multi-armed bandits, Markov decision processes (MDPs),  
 58 and partially observable MDPs. For multi-armed bandits, we find the Bayesian regret incurred by  
 59 ExPerior is proportional to the entropy of the optimal action under the prior distribution, aligning  
 60 with the entropy of expert policy if the experts are optimal. We introduce a frequentist algorithm for  
 61 multi-armed bandits and prove a Bayesian regret bound proportional to a term that closely resembles  
 62 the entropy of the optimal action. Our results suggest using the entropy of expert demonstrations to  
 63 evaluate the impact of unobserved factors.

## 64 2 Related Work

65 Our work is an addition to the recent body of reinforcement learning research that leverages of-  
 66 fline demonstrations to speed up online learning [28, 10, 29, 7, 9]. Classic algorithms such as  
 67 DDPGfD [30] and DQfD [31] achieve this by combining imitation learning and RL. They modify  
 68 DDPG [5] and DQN [1] by warm-starting the algorithms' replay buffers with expert trajectories  
 69 and ensuring that the offline data never gets overridden by online trajectories. Closely related to  
 70 our study is the meta-RL literature, which aims to accelerate learning in a given RL task by using  
 71 prior experience from related tasks [32, 33, 34]. These papers present model-agnostic meta-learning  
 72 training objectives to maximize the expected reward from novel tasks as efficiently as possible.

73 Two unique features distinguish our problem from the settings considered above. First, our setting  
 74 assumes heterogeneity within the offline data and with the online RL task that is unobserved to the

75 learner, while the (optimal) experts are privy to that heterogeneity. Second, we assume the learner  
 76 will only interact with one online task, making our setup similar to zero-shot meta-RL [35, 36, 37].  
 77 Most similar to our work is the ExPLORE algorithm [38], which assigns optimistic rewards to the  
 78 offline data during the online interaction and runs an off-policy algorithm using both online and  
 79 labelled offline data as buffers. For our setting, the algorithm incentivizes the learner to explore the  
 80 expert trajectories, leading to faster convergence. We consider this work one of our baselines.

81 Our methodology utilizes only the state-action trajectory data from expert demonstrations without  
 82 task-specific information or reward labels. Other similar methods require additional offline informa-  
 83 tion. For example, Nair et al. [29] assume that the offline data contains the reward labels and use that  
 84 to pre-train a policy, which is then fine-tuned online. Mendonca et al. [39] require task labelling for  
 85 each trajectory and use the offline data to learn a single meta-learner. Similarly, Zhou et al. [40] and  
 86 Rakelly et al. [41] require the task label and reward labels. They then infer the task during online  
 87 interaction and use the task-specific offline data. Finally, our methodology builds on posterior sam-  
 88 pling [42]. Hao et al. [22, 23] consider a similar problem using posterior sampling to leverage offline  
 89 expert demonstration data to improve online RL. However, they assume homogeneity between the  
 90 expert data and online tasks. In contrast, our setting accounts for heterogeneity.

### 91 3 Problem Setup

92 **Decision Model for Unobserved Heterogeneity of Tasks.** To account for unobserved heterogene-  
 93 ity, we consider a generalization of finite-horizon Markov Decision Processes (MDPs) with a notion  
 94 of probabilistic task variables [43, 13, 21]. The MDP’s underlying model will additionally depend  
 95 on an unobserved task variable that encapsulates some information about the specific task. In a  
 96 personalized education setup where teaching a student corresponds to a task, and the learning agent  
 97 can observe students’ characteristics, like their demographic status and grades. Other factors, such  
 98 as the student’s learning style (e.g., visual learners vs self-study), may not be readily available, even  
 99 though they are important in determining the optimal teaching style.

100 Let  $\mathcal{C}$  be the set of all *unobserved* variables that can describe the heterogeneity of potential tasks  
 101 (e.g., the set of all possible learning styles). A (contextual) MDP  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, R, H, \rho, \mu^*)$  is  
 102 parameterized by states  $\mathcal{S}$ , actions  $\mathcal{A}$ , transition function  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \times \mathcal{C} \rightarrow \Delta(\mathcal{S})$ , reward function  
 103  $R : \mathcal{S} \times \mathcal{A} \times \mathcal{C} \rightarrow \Delta(\mathbb{R})$ , horizon  $H > 0$ , initial state distribution  $\rho \in \Delta(\mathcal{S})$ , and task distribution  
 104  $\mu^*$ . We assume the transition/reward functions and  $\mu^*$  are unknown, and for simplicity,  $\rho$  does not  
 105 depend on the task variable. For each task  $c \sim \mu^*$ , we consider  $T$  episodes, where at the beginning  
 106 of each episode  $t \in [T]$ , an initial state  $s_1 \sim \rho$  is sampled. Then, at each timestep  $h \in [H]$ , the  
 107 learner chooses an action  $a_h \in \mathcal{A}$ , observes a reward  $r_h \sim R(s_h, a_h, c)$  and the next state  $s_{h+1} \sim$   
 108  $\mathcal{T}(s_h, a_h, c)$ . Without loss of generality, we assume the states are partitioned by  $[H]$  to make the  
 109 notation invariant to timestep  $h$ . Let  $\Pi$  be the set of all Markovian policies. For a policy function  
 110  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A}) \in \Pi$  and task variable  $c$ , we define the value function  $V_c(\pi) = \mathbb{E} \left[ \sum_{h=1}^H r_h \mid \pi, c \right]$   
 111 and the Q-function as  $Q_c^\pi(s, a) := \mathbb{E} \left[ \sum_{h'=h}^H r_{h'} \mid s_h = s, a_h = a, \pi, c \right]$  for all  $s \in \mathcal{S}, a \in \mathcal{A}$ .  
 112 Moreover, we define the optimal policy for a task variable  $c \in \mathcal{C}$  as  $\pi_c := \arg \max_{\pi \in \Pi} V_c(\pi)$ . Note  
 113 that since the task variable is unobserved, the learner’s policy will not depend on it. The learning  
 114 agent’s goal is to learn history-dependent distributions  $p^1, \dots, p^T \in \Delta(\Pi)$  over Markovian policies  
 115 to minimize the expected regret, defined as  $\text{Reg} := \mathbb{E}_{c \sim \mu^*} \left[ \sum_{t=1}^T V_c(\pi_c) - \mathbb{E}_{\pi^t \sim p^t} [V_c(\pi^t)] \right]$ .

116 The above setup assumes a fixed distribution  $\mu^*$  over the set of learning styles and aims to minimize  
 117 expected regret over the population of students. Our setup and regret assume the unobserved factors  
 118 remain fixed during training. This captures scenarios wherein the unobserved variables correspond  
 119 to less-variant factors (a student’s learning style is more likely to remain unchanged). No learn-  
 120 ing algorithm can control the regret value if we allow the unobserved factors to change arbitrarily  
 121 throughout  $T$  episodes without access to hidden information. Consider a two-armed bandit with  
 122 a task value drawn with uniform probability from  $\mathcal{C} = \{c_1, c_2\}$  and can change at each episode.  
 123 Assume the expected reward of the first arm under  $c_1$  and  $c_2$  is one and zero, respectively, and it is  
 124 reversed for the other arm. Any algorithm that does not have access to  $c$  would result in linear regret  
 125 since each action is sub-optimal with a probability of 0.5, independent of the algorithm’s choice.

126 **Remark.** Our setup can be formulated as a Bayesian model parameterized by  $\mathcal{C}$ , and our regret  
 127 can be seen as the Bayesian regret of the learner. However, the distribution  $\mu^*$  is not the learner’s  
 128 prior belief about the true model as it is often formulated in Bayesian learning, but a distribution

129 over potential tasks that the learner can encounter. Our setup can thus be seen as a meta-learning  
 130 problem. In fact, it is *zero-shot* meta-learning since we do not assume having access to more than  
 131 one online task during training — we only learn the prior distribution using the offline data.

132 **Expert Demonstrations.** In addition to the online setting described above, we assume the learner  
 133 has access to an offline dataset of expert demonstrations  $\mathcal{D}_E$ , where each demonstration  $\tau_E =$   
 134  $(s_1, a_1, s_2, a_2, \dots, s_H, a_H, s_{H+1})$  refers to an interaction of the expert with a decision-making task  
 135 during a *single* episode, containing the actions made by the expert and the resulting states. We as-  
 136 sume that the unobserved task-specific variables for  $\mathcal{D}_E$  are drawn i.i.d. from distribution  $\mu^*$ , and the  
 137 expert had access to such unobserved variables (private information) during their decision-making.  
 138 Moreover, we assume the expert follows a near-optimal strategy [22, 23].

139 **Assumption 1** (Noisily Rational Expert). For any  $c \in \mathcal{C}$ , experts select actions based on a dis-  
 140 tribution defined as  $p_E(a | s; c) \propto \exp\{\beta \cdot Q_c^{\pi_c}(s, a)\}$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}$ , and some known  
 141 competence value of  $\beta \in [0, \infty]$ . In particular, the expert follows the optimal policy if  $\beta \rightarrow \infty$ .

142 We assume experts do not provide any rationale for their strategy, nor do we have access to rewards  
 143 in the offline data; this is a combination of imitation and online learning rather than offline RL.

#### 144 4 Experts-as-Priors Framework for Unobserved Heterogeneity

145 Our goal is to leverage offline data to help guide the learner through its interaction with the decision-  
 146 making task. The key idea is to use expert demonstrations to infer a *prior* distribution over  $\mathcal{C}$  and then  
 147 to use a Bayesian approach such as posterior sampling [27] to utilize the inferred prior for a more  
 148 informative exploration. If the current task is from the same distribution of tasks in the offline data,  
 149 we expect that using such priors will lead to faster convergence to optimal trajectories compared to  
 150 the commonly used non-informative priors. Consider the personalized education example. Suppose  
 151 we have gathered offline data on an expert’s teaching strategies for students with similar observed  
 152 information like grade, age, location, etc. The teacher can observe more fine-grained information  
 153 about the students that is generally absent from the collected data (e.g., their learning style). Our  
 154 work relies on the following observation: The space of the optimal strategies for students with  
 155 similar observed information but different learning styles is often much smaller than the space of all  
 156 possible strategies. With the inferred prior distribution, the learner needs only to focus on the span of  
 157 potentially optimal strategies for a new student, allowing for significantly more efficient exploration.

158 We resort to empirical Bayes and use maximum marginal likelihood estimation [44] to construct a  
 159 prior distribution from  $\mathcal{D}_E$ . Given a probability distribution (prior)  $\mu$  on  $\mathcal{C}$ , the marginal likelihood  
 160 of an expert demonstration  $\tau_E = (s_1, a_1, s_2, a_2, \dots, s_H, a_H, s_{H+1}) \in \mathcal{D}_E$  is given by

$$160 \quad P_E(\tau_E; \mu) = \mathbb{E}_{c \sim \mu} \left[ \rho(s_1) \cdot \prod_{h=1}^H p_E(a_h | s_h; c) \mathcal{T}(s_{h+1} | s_h, a_h, c) \right]. \quad (1)$$

161 We aim to find a prior distribution to maximize the log-likelihood of  $\mathcal{D}_E$  under the model described in  
 162 (1). This is equivalent to minimizing the KL divergence between the marginal likelihood  $P_E$  and the  
 163 empirical distribution of expert demonstrations, which we denote by  $\hat{P}_E$ . In particular, we form an  
 164 uncertainty set over the set of plausible priors as  $\mathcal{P}(\epsilon) := \left\{ \mu; D_{\text{KL}}(\hat{P}_E \parallel P_E(\cdot; \mu)) \leq \epsilon \right\}$ , where  
 165 the value of  $\epsilon$  can be chosen based on the number of samples so the uncertainty set contains the  
 166 true prior with high probability [35]. However, the set of plausible priors does not uniquely identify  
 167 the appropriate prior. In fact, even for  $\epsilon = 0$ ,  $\mathcal{P}(\epsilon)$  can have infinite plausible priors. To solve this  
 168 ill-posed problem, we propose two approaches, parametric and nonparametric prior learning.

169 **Parametric Experts-as-Priors.** For settings where we have existing knowledge about the paramet-  
 170 ric form of the prior, we can directly apply maximum marginal likelihood estimation to learn it. In  
 171 particular, we define the parametric expert prior as the following. Note that we can calculate the  
 172 gradients of the marginal likelihood using the score function estimator [45].

173 **Definition 1** (Parametric Expert Prior). Let  $\Theta$  be a set of plausible prior distribution parameters  
 174 (e.g., Beta distribution parameters for a Bernoulli bandit). We call  $\mu_{\theta^*}$  a parametric expert prior, iff  
 175  $\theta^* \in \arg \min_{\theta \in \Theta} \sum_{\tau \in \mathcal{D}_E} -\log P_E(\tau; \mu_{\theta})$ .

176 **Nonparametric Experts-as-Priors.** For settings where there is no existing knowledge on the para-  
 177 metric form of the prior, we can employ the principle of maximum entropy to choose the *least*  
 178 *informative* prior that is compatible with expert data:

179 **Definition 2** (Max-Entropy Expert Prior). Let  $\mu_0$  be a non-informative prior on  $\mathcal{C}$  (e.g., a uniform  
180 distribution). Given some  $\epsilon > 0$ , we define the maximum entropy expert prior  $\mu_{\text{ME}}$  as the solution  
181 to the following optimization problem:

$$\mu_{\text{ME}} = \arg \min_{\mu} \text{D}_{\text{KL}}(\mu \parallel \mu_0) \quad \text{s.t.} \quad \mu \in \mathcal{P}(\epsilon). \quad (2)$$

182 Note that the set of plausible priors  $\mathcal{P}(\epsilon)$  is a convex set, and therefore, (2) is a convex optimization  
183 problem. We derive the solution to problem (2) using Fenchel’s duality theorem [46, 47]:

185 **Proposition 1** (Max-Entropy Expert Prior). Let  $N = |\mathcal{D}_E|$  be the number of demonstrations in  $\mathcal{D}_E$ .  
186 For each  $c \in \mathcal{C}$  and demonstration  $\tau_E = (s_1, a_1, s_2, a_2, \dots, s_H, a_H, s_{H+1}) \in \mathcal{D}_E$ , define  $m_{\tau_E}(c)$  as  
187 the (partial) likelihood of  $\tau_E$  under  $c$ , i.e.,  $m_{\tau_E}(c) = \prod_{h=1}^H p_E(a_h \mid s_h; c) \mathcal{T}(s_{h+1} \mid s_h, a_h, c)$ .

188 Denote  $\mathbf{m}(c) \in \mathbb{R}^N$  as the vector with elements  $m_{\tau_E}(c)$  for  $\tau_E \in \mathcal{D}_E$ . Moreover, let  $\lambda^* \in \mathbb{R}^{\geq 0}$  be  
189 the optimal solution to the Lagrange dual problem of (2). Then, the solution to optimization (2) is:

$$\mu_{\text{ME}}(c) = \lim_{n \rightarrow \infty} \frac{\exp\{\mathbf{m}(c)^\top \boldsymbol{\alpha}_n\}}{\mathbb{E}_{c \sim \mu_0} [\exp\{\mathbf{m}(c)^\top \boldsymbol{\alpha}_n\}]},$$

190 where  $\{\boldsymbol{\alpha}_n\}_{n=1}^\infty$  is a sequence converging to the following supremum:

$$\sup_{\boldsymbol{\alpha} \in \mathbb{R}^N} -\log \mathbb{E}_{c \sim \mu_0} [\exp\{\mathbf{m}(c)^\top \boldsymbol{\alpha}\}] + \frac{\lambda^*}{N} \sum_{i=1}^N \log \left( \frac{N \cdot \alpha_i}{\lambda^*} \right). \quad (3)$$

191 The proof is provided in Appendix A.3. Instead of solving for  $\lambda^*$ , we set it as a hyperparameter  
192 and then solve (3). Even though Proposition 1 requires the correct form of Q-functions for different  
193 values of  $c$ , we will see in the following sections that we can parameterize the Q-functions and treat  
194 those parameters as a proxy for the unobserved factors. Once such a prior is derived, we can employ  
195 any Bayesian approach for the decision-making task. We provide a pseudo-algorithm for ExPerior  
196 in Appendix B. The following sections will detail the algorithm for bandits and MDPs.

## 197 5 Learning in Bandits

198 **K-armed Bandits.** For  $K$ -armed bandits, note that  $\mathcal{S} = \emptyset$ ,  $H = 1$ , and  $\mathcal{A} = \{1, \dots, K\}$ . Each  
199 expert demonstration  $\tau_E = a$  will be the pulled arm by the expert for a particular bandit, and the  
200 (partial) likelihood function in Proposition 1 can be simplified as  $m_{\tau_E}(c) = p_E(a; c)$ . This likeli-  
201 hood function only depends on the task variable  $c$  through the expert policy  $p_E$ , and since  $p_E$  only  
202 depends on  $c$  through the mean reward function (Assumption 1), we can consider the set of mean  
203 reward functions as a proxy for the unobserved task variables  $\mathcal{C}$ . e.g. in a Bernoulli  $K$ -armed bandit  
204 setting, we can define  $\mathcal{C}_{\text{Ber}} = \{a \mapsto \langle \mathbf{e}_a, \boldsymbol{\vartheta} \rangle; \boldsymbol{\vartheta} \in [0, 1]^K\}$ .

205 **Stochastic Contextual Bandits.** In contextual bandits, the state space  $\mathcal{S}$  is the set of contexts and  
206  $H = 1$ . Therefore, the likelihood function for a demonstration  $\tau_E = (s, a)$  will be  $m_{\tau_E}(c) =$   
207  $p_E(a \mid s; c)$ . Like  $K$ -armed bandits, the likelihood function only depends on  $c$  through the expert  
208 policy. Therefore, we can similarly define the set of mean reward functions as the proxy for the  
209 unobserved task variables. For instance, we can consider the task parameters for linear contextual  
210 bandits as  $\mathcal{C}_{\text{Lin}} = \{(s, a) \mapsto \langle \phi(s, a), \boldsymbol{\vartheta} \rangle; \boldsymbol{\vartheta} \in \mathbb{R}^d\}$ , for a known feature function  $\phi: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ .

211 **Posterior Sampling.** With the above parameterizations of  $\mathcal{C}$ , we can use Proposition 1 to derive  
212 the maximum entropy prior distribution over the task parameters. However, we cannot sample from  
213 the exact posterior since the derived prior is not a conjugate prior for standard likelihood functions.  
214 Instead, we resort to approximate posterior sampling via stochastic gradient Langevin dynamics  
215 (SGLD) [48]. We call this method ExPerior-MaxEnt in our experiments. We also employ a  
216 parametric approach as discussed in section 4, which we call ExPerior-Param. In particular, we  
217 use the Beta distribution as our prior model and learn the parametric expert prior in Definition 1.  
218 ExPerior-Param has an advantage over ExPerior-MaxEnt since it provides exact posterior sam-  
219 pling for Bernoulli bandits.

220 We aim to evaluate our approach compared to other baselines, including online methods that do  
221 not use expert data and offline behaviour cloning. We provide an empirical regret analysis for  
222 ExPerior based on the informativeness of expert data, number of actions, and number of training  
223 episodes. We also discuss the robustness of ExPerior to misspecified expert models and the advan-  
224 tage of ExPerior-MaxEnt to ExPerior-Param when the parametric prior model is misspecified.  
225 To characterize the effect of expert data on the learner’s performance, we propose an alternative for  
226  $K$ -armed bandits inspired by the successive elimination and derive a Bayesian regret bound for it.



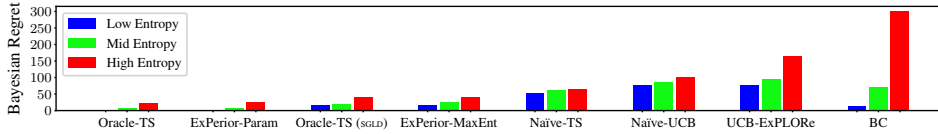


Figure 2: The Bayesian regret of ExPerior and baselines for  $K$ -armed Bernoulli bandits ( $K = 10$ ). We consider three categories of task distributions based on the entropy of the optimal action.

227 **Experiments.** We consider  $K$ -armed Bernoulli bandits for our experimental setup (code: <https://anonymous.4open.science/r/ExPerior-0773>). We evaluate the learning algorithms in  
 228 terms of the Bayesian regret over multiple task distributions  $\mu^*$ . We consider up to  $N_{\mu^*} = 64$   
 229 different beta task distributions, where their parameters are chosen to span a different range of het-  
 230 erogeneity, consisting of tasks with various expert data informativeness. To estimate the Bayesian  
 231 regret, we sample  $N_{\text{task}} = 128$  bandit tasks from each task distribution and calculate the average  
 232 regret. We use  $N_E = 1000$  expert demonstrations for each task distribution in our experiments. We  
 233 compare ExPerior to the following baselines: (1) Behaviour cloning (BC), which learns a policy by  
 234 minimizing the cross-entropy loss between the expert demonstrations and the agent’s policy solely  
 235 based on offline data. (2) Naïve Thompson sampling (Naïve-TS) that chooses the action with the  
 236 highest sampled mean from a posterior distribution under an uninformative prior. (3) Naïve upper  
 237 confidence bound (Naïve-UCB) algorithm that selects the action with the highest upper confidence  
 238 bound. Both Naïve-TS and Naïve-UCB ignore expert demonstrations. (4) UCB-ExPLORe, a variant  
 239 of the algorithm proposed by Li et al. [38] tailored to bandits. It labels the expert data with opti-  
 240 mistic rewards and then uses it alongside online data to compute the upper confidence bounds for  
 241 exploration, and (5) Oracle-TS, which performs exact Thompson sampling having access to the  
 242 true task distribution  $\mu^*$ . For a more fair comparison, we also consider a variant of Oracle-TS,  
 243 which uses SGLD for approximate posterior sampling.  
 244

245 **Comparison to baselines.** Figure 2 demonstrates the average Bayesian regret for various task distri-  
 246 butions over  $T = 1500$  episodes with  $K = 10$  arms. To better understand the effect of expert data,  
 247 we categorize the task distributions by the entropy of their optimal actions into low entropy (less  
 248 than 0.8), high entropy (greater than 1.6), and medium entropy. Oracle-TS and ExPerior-Param  
 249 outperform other baselines, yet the performance of ExPerior is comparable to the SGLD variant  
 250 of Oracle-TS. This indicates that the maximum entropy prior derived from Proposition 1 closely  
 251 approximates the true task distribution,  $\mu^*$ , with the performance difference with Oracle-TS is  
 252 primarily due to approximate posterior sampling. Moreover, the pure online algorithms Naïve-TS  
 253 and Naïve-UCB, which disregard expert data, display similar performance across different entropy  
 254 levels, contrasting with other algorithms that show significantly reduced regret in low-entropy con-  
 255 texts. This underlines the impact of expert data in settings where the unobserved confounding has  
 256 less effect on the optimal actions. Specifically, in the extreme case of no task heterogeneity, BC  
 257 is anticipated to yield optimal performance. Additionally, Naïve-UCB surpasses UCB-ExPLORe in  
 258 medium and high entropy settings, possibly due to the over-optimism of the reward labelling in Li  
 259 et al. [38], which can hurt the performance when the expert demonstrations are uninformative.

260 **Empirical regret analysis for Experts-as-Priors.** We examine how the quality of expert demon-  
 261 strations affects the Bayesian regret achieved by Algorithm 2. Settings with highly informative  
 262 demonstrations, where unobserved factors minimally affect the optimal action, should exhibit near-  
 263 zero regret since there is no diversity in the tasks, and the experts are near-optimal. Conversely,  
 264 in scenarios where unobserved factors significantly influence the optimal actions, we anticipate the  
 265 regret to align with standard online regret bounds, similar to the outcomes of Thompson sampling  
 266 with a non-informative prior. We conduct trials with ExPerior and Oracle-TS across various num-  
 267 bers of arms over  $T = 1500$  episodes, calculating the mean and standard error of Bayesian regret  
 268 across distinct task distributions. As depicted in Figure 3 (a), both ExPerior and Oracle-TS yield  
 269 sub-linear regret relative to  $K$  and  $T$ , comparable to the established regret bound of  $\mathcal{O}(\sqrt{KT})$  for  
 270 Thompson sampling. However, the middle panel indicates that the regret of ExPerior is proportional  
 271 to the entropy of the optimal action, having an almost *linear* relationship. This observation seems to  
 272 be in contrast with the standard Bayesian regret bounds for Thompson sampling under correct prior  
 273 that have shown a sublinear relationship of  $\mathcal{O}(\sqrt{\text{Ent}(\pi_c)})$ , where  $\text{Ent}(\pi_c)$  denotes the entropy of  
 274 the optimal action under  $\mu^*$  [49]. We analyze this observation more concretely below.

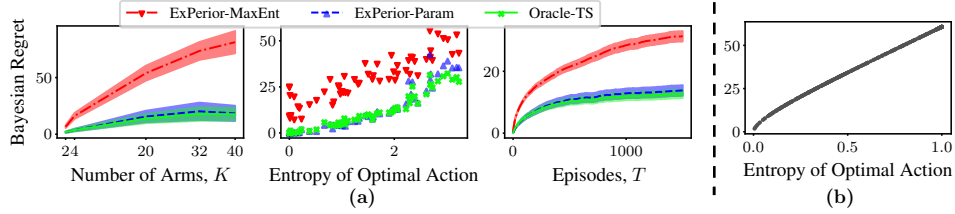


Figure 3: (a) Empirical analysis of ExPerior’s regret in Bernoulli bandits based on the (left) number of arms, (middle) entropy of the optimal action, and (right) number of episodes. (b) The regret bound from Theorem 2 V.S. the entropy of the optimal action. The linear relationship is consistent with the middle panel of (a).

275 **Ablations.** We run additional experiments in Appendix C.1 to assess the robustness of ExPerior to  
 276 misspecified experts. We create expert data from different experts with various competence levels,  
 277 such as optimal, noisily rational, and random-optimal experts, where the latter chooses an action opti-  
 278 mally with a fixed probability and randomly otherwise. Table 2 in the appendix shows ExPerior’s  
 279 robustness to different expert models. With  $\beta = 10$  for training ExPerior-MaxEnt and  $\beta = 1$   
 280 for ExPerior-Param achieves consistent out-performance among different expert types. We evalu-  
 281 ate the advantage of learning nonparametric max-entropy prior over misspecified parametric pri-  
 282 ors in Table 3. Even though ExPerior-Param with Beta model outperforms ExPerior-MaxEnt,  
 283 ExPerior-MaxEnt is superior to ExPerior-Param if the prior is chosen as Gaussian or Gamma.

284 **An Alternative Frequentist Approach for  $K$ -armed Bandits** To analyze the effect of expert  
 285 data on the Bayesian regret, we devise an alternative *frequentist* approach, based on the successive  
 286 elimination algorithm [50], which follows a similar intuition to Experts-as-Priors. In particular, we  
 287 prove a bound on its Bayesian regret and show that the derived bound is proportional to a term that  
 288 closely resembles the entropy of the optimal action, showing that the observation in the middle panel  
 289 of Figure 3 (a) is consistent within different approaches.

290 The idea of successive elimination is to identify suboptimal arms and deactivate them over time. In  
 291 particular, it runs a uniform sampling policy among active arms and builds confidence intervals for  
 292 each. It then deactivates all the arms with an upper confidence bound smaller than at least one arm’s  
 293 lower confidence bound. We modify this algorithm using the policy derived from expert demonstra-  
 294 tions instead of a uniform sampling policy. Recall that in  $K$ -armed bandits, each expert trajectory  
 295  $\tau_E$  represents the pulled arm by the expert. Hence, the empirical distribution of expert demonstra-  
 296 tions can be seen as a sampling policy over different arms. The concrete algorithm is provided in  
 297 Algorithm 1 in Appendix A.4. We now provide a Bayesian regret bound of this algorithm.

298 **Theorem 2.** Consider a stochastic  $K$ -armed bandit and let  $p$  be the empirical expert policy. Assume  
 299 that (i) the mean reward function is bounded in  $[0, 1]$  for all arms, (ii)  $T \geq \frac{1}{\min_{a:p(a) \neq 0} p(a)}$ , (iii) the  
 300 expert is optimal, i.e.,  $\forall a \in \mathcal{A} : p(a) = \mathbb{P}_E(a; \mu^*)$  and  $\beta \rightarrow \infty$ , and (iv) the learner follows  
 301 Algorithm 1. Then, with probability at least  $1 - \delta$ ,

$$Reg \lesssim \sqrt{T \log(TK/\delta)} \sum_{a, a' \in \mathcal{A}, a \neq a'} \sqrt{\frac{p(a)}{p(a) + p(a')}} \left(1 - \frac{p(a)}{p(a) + p(a')}\right) \left[\sqrt{p(a)} + \sqrt{p(a')}\right]. \quad (4)$$

302 See Appendix A.4 for the proof. Two terms in (4) depend on expert data: (1) The relative standard  
 303 deviation between any two pairs of arms and (2) a scaling factor that depends on the magnitude of  
 304 probability that the arms are optimal. For homogeneous demonstrations, where the expert data only  
 305 includes one unique pulled arm, the standard deviation (Term 1) is zero, resulting in zero regret. On  
 306 the other hand, in extreme heterogeneity, where the empirical expert distribution is uniform over the  
 307 arms, we have  $Reg \lesssim \sqrt{KT \log T}$ , a similar bound for standard successive elimination. Finally, to  
 308 assess the relationship between the regret bound and the entropy of the expert data, we fix  $K = 2$ ,  
 309  $T = 100$ , and plot the bound from (4) as a function of the entropy of the optimal action for various  
 310 task distributions. Figure 3 (b) demonstrates a linear relationship, similar to the regret incurred by  
 311 ExPerior in Figure 3 (a). This observation opens up new directions to further analyze the theoretical  
 312 regret for ExPerior and propose similar frequentist approaches for MDPs.

## 313 6 Learning in Markov Decision Processes (MDPs)

314 For MDPs, we need to parameterize both the mean reward and transition functions. However, we  
 315 assume the transition functions are invariant to the task variables to simplify our methodology and

Table 1: The average reward per episode in Frozen Lake (PODMP) after 90,000 training steps.

	Fixed # Hazard = 9				Fixed $\beta = 1$			
	$\beta = 0.1$	$\beta = 1$	$\beta = 2.5$	$\beta = 10$	# Hazard = 2	# Hazard = 5	# Hazard = 7	# Hazard = 9
ExPerior-MaxEnt	-22.58 ± 1.17	<b>6.00 ± 0.00</b>	3.58 ± 0.89	1.62 ± 1.85	11.47 ± 0.52	<b>5.71 ± 0.67</b>	<b>6.00 ± 0.00</b>	<b>6.00 ± 0.00</b>
ExPerior-Param	-23.32 ± 0.69	-4.31 ± 1.80	5.27 ± 0.51	<b>6.00 ± 0.00</b>	<b>12.00 ± 0.37</b>	2.11 ± 1.41	5.42 ± 0.40	-4.31 ± 1.80
Naïve Boot-DQN	-23.32 ± 0.69	-23.32 ± 0.69	-23.32 ± 0.69	-23.32 ± 0.69	-14.36 ± 5.88	-20.57 ± 2.91	-20.39 ± 1.75	-23.32 ± 0.69
ExPLORe	<b>5.99 ± 0.00</b>	<b>6.00 ± 0.00</b>	<b>6.00 ± 0.00</b>	<b>6.00 ± 0.00</b>	-30.68 ± 12.40	-10.64 ± 16.64	-13.00 ± 19.00	<b>6.00 ± 0.00</b>
Optimal	6.00 ± 0.00	6.00 ± 0.00	6.00 ± 0.00	6.00 ± 0.00	12.00 ± 0.37	6.53 ± 0.31	6.00 ± 0.00	6.00 ± 0.00

316 avoid extra modelling assumptions. Under this assumption, it is sufficient to parameterize the *opti-*  
 317 *mal* Q-functions, e.g., using a deep Q-network (DQN) and treat those parameters as a proxy for the  
 318 task variables, i.e.,  $\mathcal{C}_{\text{MDP}} := \{(s, a) \mapsto Q(s, a; \theta); \theta \in \Theta\}$ , where  $\Theta$  is the set of parameters for a  
 319 DQN. We can then derive a closed-form log-pdf of the posterior distribution under the maximum en-  
 320 tropy prior. See Appendix A.5 for details. The derived posterior log-pdf can then be used as the loss  
 321 function for DQN Langevin Monte Carlo [51, 52] as the counterpart for Thompson sampling with  
 322 SGLD. However, running Langevin dynamics can lead to highly unstable policies due to the com-  
 323 plexity of the optimization landscape in DQNs. Instead of sampling from the posterior distribution,  
 324 we use a heuristic that combines the learned prior distribution with bootstrapped DQNs [53].

325 The original method of Bootstrapped DQNs utilizes an ensemble of  $L$  randomly initialized Q-  
 326 networks. It samples a Q-network uniformly at each episode and uses it to collect data. Then,  
 327 each Q-network is trained using the temporal difference loss on parts of or possibly the entire  
 328 collected data. This method and its subsequent iterations [54, 55, 56] achieve deep exploration  
 329 by ensuring diversity among the learned Q-networks. To incorporate Bootstrapped DQN into  
 330 the ExPerior framework and utilize the expert data, we can formulate the ensemble as a discrete  
 331 prior distribution over the Q-networks. Let  $\theta_{\text{ens}} = (\theta_{\text{ens}}^1, \dots, \theta_{\text{ens}}^L)$  be the parameter vector  
 332 for an ensemble of Q-functions. We can define the ensemble prior, parameterized by  $\theta_{\text{ens}}$ , as  
 333  $\mu_{\theta_{\text{ens}}}(\theta) := \frac{1}{L} \sum_{i=1}^L \mathbb{I}(\theta_{\text{ens}}^i = \theta)$  for any  $\theta \in \Theta$ . Based on this prior model, we can learn the  
 334 parametric expert prior using maximum marginal likelihood estimation, as formulated below.

335 **Proposition 3** (Ensemble Marginal Likelihood). *Consider a contextual MDP  $\mathcal{M} =$*   
 336  *$(\mathcal{S}, \mathcal{A}, \mathcal{T}, R, H, \rho, \mu^*)$ . Assume the transition function  $\mathcal{T}$  does not depend on the task variables*  
 337 *and Assumption 1 holds. Then, the negative marginal log-likelihood of expert data  $\mathcal{D}_E$  under the*  
 338 *ensemble prior  $\mu_{\theta_{\text{ens}}}$  is upper bounded by*

$$-\log P_E(\mathcal{D}_E; \mu_{\theta_{\text{ens}}}) \leq \frac{1}{L} \sum_{i=1}^L \sum_{\tau \in \mathcal{D}_E} \sum_{(s,a) \in \tau} \log \left( \sum_{a' \in \mathcal{A}} \exp \left\{ \beta \cdot Q(s, a'; \theta_{\text{ens}}^i) \right\} \right) - \beta \cdot Q(s, a; \theta_{\text{ens}}^i),$$

339 where  $\beta$  is the competence level of the expert in Assumption 1.

340 Proposition 3 is proved in Appendix A.6. We can then initialize the Q-networks in the Bootstrapped  
 341 DQN method using ensemble parameters that minimize the above upper bound. We will refer to this  
 342 method as ExPerior-Param. As an alternative approach, instead of minimizing the above upper  
 343 bound, we can match the discrete prior distribution  $\mu_{\theta_{\text{ens}}}$  to the max-entropy prior by initializing  
 344 the Q-functions in the ensemble with parameters sampled from the max-entropy expert prior. In  
 345 particular, we can apply SGLD on the log-pdf of the max-entropy prior derived in Appendix A.5.  
 346 We will refer to this approach as ExPerior-MaxEnt.

347 **Experimental Setup.** A main challenge in RL is the reward *sparsity*, where the learner needs  
 348 to explore the environment deeply to observe reward states. Utilizing expert demonstrations can  
 349 significantly improve the efficiency of exploration. For this reason, we focus on "Deep Sea," a  
 350 sparse-reward tabular RL environment proposed by Osband et al. [55] to assess deep exploration for  
 351 different RL methods. The environment is an  $M \times M$  grid, where the agent starts at the top-left  
 352 corner of the map, and at each time step, it chooses an action from  $\mathcal{A} = \{\text{left}, \text{right}\}$  to move to  
 353 the left or right column, while going down by one row. In the original version of Deep Sea, the goal  
 354 is always on the bottom-right corner of the map. We introduce unobserved task variables by defining  
 355 a distribution over the goal columns while keeping the goal row the same. We consider four types of  
 356 goal distributions where the goal is situated at (1) the bottom-right corner of the grid, (2) uniformly at  
 357 the bottom of any of the right-most  $\frac{M}{4}$  columns, (3) uniformly at the bottom of any of the right-most  
 358  $\frac{M}{2}$  columns, and (4) uniformly at the bottom of any of the  $M$  columns. We set  $M = 30$  and generate  
 359  $\bar{N} = 1000$  samples from the optimal policies as offline expert demonstrations. To further evaluate  
 360 ExPerior and showcase its applicability to partially-observed MDP, we also consider the "Frozen



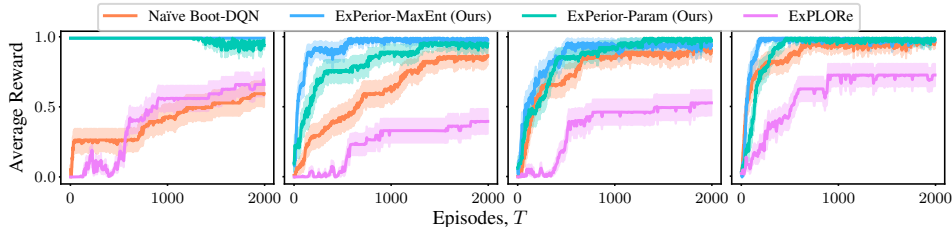


Figure 4: The average reward per episode over 2,000 episodes in "Deep Sea." The goal is located at the right column, uniformly at the right-most quarter of the columns, uniformly at the right-most half, and uniformly at random over all the columns, respectively. ExPerior outperforms the baselines in all instances.

361 Lake" environment, which requires the learner to navigate to a goal while avoiding hazards [17].  
 362 The learner cannot observe the hazard location, while the expert has access to the whole map. Taking  
 363 action, reaching the goal, and hitting the hazard incur rewards of -2, 20, and -100, respectively. The  
 364 frozen lake map is  $5 \times 5$ , where the hazard (weak ice) is randomly located in the interior squares. We  
 365 consider different settings with 2, 5, 7, and 9 potential locations for the hazard. At the start of each  
 366 episode, the hazard will be chosen randomly within the potential locations. We generate  $N = 1000$   
 367 samples from noisily rational experts with different competence levels for this environment. See  
 368 Appendix C.2 for the MDP experiments in Frozen Lake experiments.

369 **Baselines.** We compare ExPerior to the following baselines. (1) ExPLORe, proposed by Li et al. [38]  
 370 to accelerate off-policy reinforcement learning using unlabeled prior data. In this method, the offline  
 371 demonstrations are assigned optimistic reward labels generated using the online data with regular  
 372 updates. This information is then combined with the buffer data to perform off-policy learning.  
 373 (2) Naïve Boot-DQN, which is the original implementation of Bootstrapped DQN with randomly  
 374 initialized Q-networks [53]. The latter baseline is purely online.

375 **Deep Sea Results.** Figure 4 demonstrates the average reward per episode achieved by the baselines  
 376 for  $T = 2000$  episodes. For each goal distribution, we run the baselines with 30 different seeds and  
 377 take the average to estimate the expected reward. ExPerior outperforms the baselines in all instances.  
 378 However, the gap between ExPerior and the fully online Naïve Boot-DQN, which measures the ef-  
 379 fect of using the expert data, decreases as we go from the low-entropy setting (upper left) to the  
 380 high-entropy task distribution (bottom right). This is consistent with the empirical and theoretical  
 381 results discussed in section 5 and confirms our expectation that the expert demonstrations may not  
 382 be helpful under strong unobserved confounding (strong task heterogeneity). The ExPLORe base-  
 383 line substantially underperforms, even compared to the fully online Naïve Boot-DQN (except for the  
 384 first task distribution with zero-entropy). We suspect this is because ExPLORe uses actor-critic  
 385 methods as its backbone model, which are shown to struggle with deep exploration [57].

386 **Frozen Lake Results.** We run all the baselines for 90,000 steps with 30 different seeds. Table 1  
 387 shows the average reward after 500 evaluation steps at the end of the training. ExPerior outperforms  
 388 the baselines in almost all instances except for the case of  $\beta = 0.1$ , which corresponds to a nearly  
 389 random expert. On the other hand, ExPLORe achieves near-optimal results for  $\beta = 0.1$ . We hypoth-  
 390 esize that ExPLORe's performance is mainly due to the superiority of their base actor-critic model  
 391 since it can achieve near-optimal performance even when the expert trajectories are low-quality.

## 392 7 Conclusion

393 We introduce the Experts-as-Priors (ExPerior) framework, a novel empirical Bayes approach, to  
 394 address the problem of sequential decision-making using expert demonstrations with unobserved  
 395 heterogeneity. We ground our methodology in the maximum entropy principle to infer a prior dis-  
 396 tribution from expert data that guides the learning process in both bandit settings and Markov De-  
 397 cision Processes (MDPs). This advantage underscores the utility of our approach in contexts where  
 398 the learner faces uncertainty and variability in task parameters, a common challenge in real-world  
 399 applications from autonomous driving to personalized learning environments. Our work contributes  
 400 to the understanding of leveraging expert demonstrations under unobserved heterogeneity and offers  
 401 a practical framework readily applied to a broad spectrum of decision-making tasks. We provide a  
 402 principled way to incorporate the wealth of information contained in expert behaviours, thus opening  
 403 new avenues for research in meta-reinforcement learning. One limitation of our work is the limited  
 404 set of experiments, especially those with human-in-the-loop. Future directions include extending to  
 405 more complex environments, and further investigating our RL algorithm's theoretical properties.

## References

- 406  
407 [1] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan  
408 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint*  
409 *arXiv:1312.5602*, 2013.
- 410 [2] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van  
411 Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot,  
412 et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529  
413 (7587):484–489, 2016.
- 414 [3] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur  
415 Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, et al. Mastering  
416 chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint*  
417 *arXiv:1712.01815*, 2017.
- 418 [4] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Ried-  
419 miller. Deterministic policy gradient algorithms. In *International conference on machine*  
420 *learning*, pages 387–395. Pmlr, 2014.
- 421 [5] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval  
422 Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning.  
423 *arXiv preprint arXiv:1509.02971*, 2015.
- 424 [6] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,  
425 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to  
426 follow instructions with human feedback. *Advances in Neural Information Processing Systems*,  
427 35:27730–27744, 2022.
- 428 [7] Andrew Wagenmaker and Aldo Pacchiano. Leveraging offline data in online reinforcement  
429 learning. In *International Conference on Machine Learning*, pages 35300–35338. PMLR,  
430 2023.
- 431 [8] Yuda Song, Yifei Zhou, Ayush Sekhari, J Andrew Bagnell, Akshay Krishnamurthy, and Wen  
432 Sun. Hybrid rl: Using both offline and online data can make rl efficient. *arXiv preprint*  
433 *arXiv:2210.06718*, 2022.
- 434 [9] Philip J Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient online reinforcement  
435 learning with offline data. In *International Conference on Machine Learning*, pages 1577–  
436 1594. PMLR, 2023.
- 437 [10] Ashvin Nair, Bob McGrew, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel.  
438 Overcoming exploration in reinforcement learning with demonstrations. In *2018 IEEE inter-*  
439 *national conference on robotics and automation (ICRA)*, pages 6292–6299. IEEE, 2018.
- 440 [11] Serkan Cabi, Sergio Gómez Colmenarejo, Alexander Novikov, Ksenia Konyushkova, Scott  
441 Reed, Rae Jeong, Konrad Zolna, Yusuf Aytar, David Budden, Mel Vecerik, et al. Scaling  
442 data-driven robotics with reward sketching and batch reinforcement learning. *arXiv preprint*  
443 *arXiv:1909.12200*, 2019.
- 444 [12] Richard A. Briesch, Pradeep K. Chintagunta, and Rosa L. Matzkin. Nonparametric Discrete  
445 Choice Models With Unobserved Heterogeneity. *Journal of Business & Economic Statistics*,  
446 28(2):291–307, 2010. ISSN 0735-0015. Publisher: [American Statistical Association, Taylor  
447 & Francis, Ltd.].
- 448 [13] Lantao Yu, Tianhe Yu, Chelsea Finn, and Stefano Ermon. Meta-inverse reinforcement learning  
449 with probabilistic context variables. *Advances in neural information processing systems*, 32,  
450 2019.
- 451 [14] Nathan Kallus and Angela Zhou. Minimax-Optimal Policy Learning Under Unobserved Con-  
452 founding. *Management Science*, 67(5):2870–2890, May 2021. ISSN 0025-1909, 1526-5501.  
453 doi: 10.1287/mnsc.2020.3699.
- 454 [15] Andrew Bennett, Nathan Kallus, Lihong Li, and Ali Mousavi. Off-policy Evaluation in  
455 Infinite-Horizon Reinforcement Learning with Latent Confounders. In *Proceedings of The*  
456 *24th International Conference on Artificial Intelligence and Statistics*, pages 1999–2007.  
457 PMLR, March 2021. ISSN: 2640-3498.

- 458 [16] Sanjiban Choudhury, Mohak Bhardwaj, Sankalp Arora, Ashish Kapoor, Gireeja Ranade, Se-  
459 bastian Scherer, and Debadepta Dey. Data-driven planning via imitation learning. *The Inter-*  
460 *national Journal of Robotics Research*, 37(13-14):1632–1672, 2018.
- 461 [17] Andrew Warrington, Jonathan W Lavington, Adam Scibior, Mark Schmidt, and Frank Wood.  
462 Robust asymmetric learning in pomdps. In *International Conference on Machine Learning*,  
463 pages 11013–11023. PMLR, 2021.
- 464 [18] Aaron Walsman, Muru Zhang, Sanjiban Choudhury, Dieter Fox, and Ali Farhadi. Impossibly  
465 good experts and how to follow them. In *The Eleventh International Conference on Learning*  
466 *Representations*, 2022.
- 467 [19] Idan Shenfeld, Zhang-Wei Hong, Aviv Tamar, and Pulkit Agrawal. TGRL: An algorithm  
468 for teacher guided reinforcement learning. In Andreas Krause, Emma Brunskill, Kyunghyun  
469 Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the*  
470 *40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine*  
471 *Learning Research*, pages 31077–31093. PMLR, 23–29 Jul 2023.
- 472 [20] Junzhe Zhang, Daniel Kumor, and Elias Bareinboim. Causal imitation learning with unob-  
473 served confounders. *Advances in neural information processing systems*, 33:12263–12274,  
474 2020.
- 475 [21] Gokul Swamy, Sanjiban Choudhury, J Bagnell, and Steven Z Wu. Sequence model imitation  
476 learning with unobserved contexts. *Advances in Neural Information Processing Systems*, 35:  
477 17665–17676, 2022.
- 478 [22] Botao Hao, Rahul Jain, Tor Lattimore, Benjamin Van Roy, and Zheng Wen. Leveraging  
479 demonstrations to improve online learning: Quality matters. In *International Conference on*  
480 *Machine Learning*, pages 12527–12545. PMLR, 2023.
- 481 [23] Botao Hao, Rahul Jain, Dengwang Tang, and Zheng Wen. Bridging imitation and online  
482 reinforcement learning: An optimistic tale. *arXiv preprint arXiv:2303.11369*, 2023.
- 483 [24] Luca Weihs, Unnat Jain, Iou-Jen Liu, Jordi Salvador, Svetlana Lazebnik, Aniruddha Kemb-  
484 havi, and Alex Schwing. Bridging the imitation gap by adaptive insubordination. *Advances in*  
485 *Neural Information Processing Systems*, 34:19134–19146, 2021.
- 486 [25] Leonardo Cella, Alessandro Lazaric, and Massimiliano Pontil. Meta-learning with stochastic  
487 linear bandits. In *International Conference on Machine Learning*, pages 1360–1370. PMLR,  
488 2020.
- 489 [26] Leonardo Cella and Massimiliano Pontil. Multi-task and meta-learning with sparse linear  
490 bandits. In *Uncertainty in Artificial Intelligence*, pages 1692–1702. PMLR, 2021.
- 491 [27] Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via  
492 posterior sampling. *Advances in Neural Information Processing Systems*, 26, 2013.
- 493 [28] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman,  
494 Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep  
495 reinforcement learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.
- 496 [29] Ashvin Nair, Abhishek Gupta, Murtaza Dalal, and Sergey Levine. Awac: Accelerating online  
497 reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- 498 [30] Mel Vecerik, Todd Hester, Jonathan Scholz, Fumin Wang, Olivier Pietquin, Bilal Piot, Nicolas  
499 Heess, Thomas Rothörl, Thomas Lampe, and Martin Riedmiller. Leveraging demonstrations  
500 for deep reinforcement learning on robotics problems with sparse rewards. *arXiv preprint*  
501 *arXiv:1707.08817*, 2017.
- 502 [31] Todd Hester, Matej Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Dan Hor-  
503 gan, John Quan, Andrew Sendonaris, Ian Osband, et al. Deep q-learning from demonstrations.  
504 In *Proceedings of the AAAI conference on artificial intelligence*, 2018.
- 505 [32] Abhishek Gupta, Russell Mendonca, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Meta-  
506 reinforcement learning of structured exploration strategies. *Advances in neural information*  
507 *processing systems*, 31, 2018.
- 508 [33] Anusha Nagabandi, Ignasi Clavera, Simin Liu, Ronald S Fearing, Pieter Abbeel, Sergey  
509 Levine, and Chelsea Finn. Learning to adapt in dynamic, real-world environments through  
510 meta-reinforcement learning. *arXiv preprint arXiv:1803.11347*, 2018.

- 511 [34] Jacob Beck, Risto Vuorio, Evan Zheran Liu, Zheng Xiong, Luisa Zintgraf, Chelsea  
512 Finn, and Shimon Whiteson. A survey of meta-reinforcement learning. *arXiv preprint*  
513 *arXiv:2301.08028*, 2023.
- 514 [35] Jay Mardia, Jiantao Jiao, Ervin Tánčzos, Robert D Nowak, and Tsachy Weissman. Concentra-  
515 tion inequalities for the empirical distribution of discrete distributions: beyond the method of  
516 types. *Information and Inference: A Journal of the IMA*, 9(4):813–850, 2020.
- 517 [36] Vinay Kumar Verma, Dhanajit Brahma, and Piyush Rai. Meta-learning for generalized zero-  
518 shot learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34,  
519 pages 6062–6069, 2020.
- 520 [37] Julie Jiang, Kristina Lerman, and Emilio Ferrara. Zero-shot meta-learning for small-scale data  
521 from human subjects. In *2023 IEEE 11th International Conference on Healthcare Informatics*  
522 *(ICHI)*, pages 311–320. IEEE, 2023.
- 523 [38] Qiyang Li, Jason Zhang, Dibya Ghosh, Amy Zhang, and Sergey Levine. Accelerating ex-  
524 ploration with unlabeled prior data. *Advances in Neural Information Processing Systems*, 36,  
525 2024.
- 526 [39] Russell Mendonca, Abhishek Gupta, Rosen Kralev, Pieter Abbeel, Sergey Levine, and Chelsea  
527 Finn. Guided meta-policy search. *Advances in Neural Information Processing Systems*, 32,  
528 2019.
- 529 [40] Allan Zhou, Eric Jang, Daniel Kappler, Alex Herzog, Mohi Khansari, Paul Wohlhart, Yunfei  
530 Bai, Mrinal Kalakrishnan, Sergey Levine, and Chelsea Finn. Watch, try, learn: Meta-learning  
531 from demonstrations and reward. *arXiv preprint arXiv:1906.03352*, 2019.
- 532 [41] Kate Rakelly, Aurick Zhou, Chelsea Finn, Sergey Levine, and Deirdre Quillen. Efficient off-  
533 policy meta-reinforcement learning via probabilistic context variables. In *International con-*  
534 *ference on machine learning*, pages 5331–5340. PMLR, 2019.
- 535 [42] Daniel J Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, Zheng Wen, et al. A tutorial  
536 on thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96, 2018.
- 537 [43] Assaf Hallak, Dotan Di Castro, and Shie Mannor. Contextual markov decision processes.  
538 *arXiv preprint arXiv:1502.02259*, 2015.
- 539 [44] Bradley P Carlin and Thomas A Louis. Empirical bayes: Past, present and future. *Journal of*  
540 *the American Statistical Association*, 95(452):1286–1289, 2000.
- 541 [45] Michael C Fu. Chapter 19 gradient estimation. *Simulation*, 13:575–616, 2006.
- 542 [46] R Tyrrell Rockafellar. *Convex analysis*, volume 11. Princeton university press, 1997.
- 543 [47] Miroslav Dudík, Steven J Phillips, and Robert E Schapire. Maximum entropy density esti-  
544 mation with generalized regularization and an application to species distribution modeling.  
545 *Journal of Machine Learning Research*, 2007.
- 546 [48] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics.  
547 In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages  
548 681–688, 2011.
- 549 [49] Daniel Russo and Benjamin Van Roy. An information-theoretic analysis of thompson sam-  
550 pling. *The Journal of Machine Learning Research*, 17(1):2442–2471, 2016.
- 551 [50] Eyal Even-Dar, Shie Mannor, and Yishay Mansour. Pac bounds for multi-armed bandit and  
552 markov decision processes. In *Computational Learning Theory: 15th Annual Conference on*  
553 *Computational Learning Theory, COLT 2002 Sydney, Australia, July 8–10, 2002 Proceedings*  
554 *15*, pages 255–270. Springer, 2002.
- 555 [51] Vikranth Dwaracherla and Benjamin Van Roy. Langevin dqn. *arXiv preprint*  
556 *arXiv:2002.07282*, 2020.
- 557 [52] Haque Ishfaq, Qingfeng Lan, Pan Xu, A. Rupam Mahmood, Doina Precup, Anima Anandku-  
558 mar, and Kamyar Azizzadenesheli. Provable and practical: Efficient exploration in reinforce-  
559 ment learning via langevin monte carlo. In *The Twelfth International Conference on Learning*  
560 *Representations*, 2024.
- 561 [53] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration  
562 via bootstrapped dqn. *Advances in neural information processing systems*, 29, 2016.

- 563 [54] Ian Osband, John Aslanides, and Albin Cassirer. Randomized prior functions for deep rein-  
564 forcement learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- 565 [55] Ian Osband, Benjamin Van Roy, Daniel J Russo, Zheng Wen, et al. Deep exploration via  
566 randomized value functions. *J. Mach. Learn. Res.*, 20(124):1–62, 2019.
- 567 [56] Ian Osband, Zheng Wen, Seyed Mohammad Asghari, Vikranth Dwaracherla, Morteza  
568 Ibrahimi, Xiuyuan Lu, and Benjamin Van Roy. Approximate thompson sampling via epis-  
569 temic neural networks. In *Uncertainty in Artificial Intelligence*, pages 1586–1595. PMLR,  
570 2023.
- 571 [57] Ian Osband, Yotam Doron, Matteo Hessel, John Aslanides, Eren Sezener, Andre Saraiva, Ka-  
572 trina McKinney, Tor Lattimore, Csaba Szepesvari, Satinder Singh, et al. Behaviour suite for  
573 reinforcement learning. *arXiv preprint arXiv:1908.03568*, 2019.
- 574 [58] Jonathan M Borwein and Qiji J Zhu. Techniques of variational analysis, 2004.



## 575 A Proofs

### 576 A.1 Notation

577 We assume  $\mathcal{C}$  is a measurable set with an appropriate  $\sigma$ -algebra and there exists a probability measure  
 578  $\mu_0$  on  $\mathcal{C}$ . We denote  $L^p(\mathcal{C}, \mu_0)$  as the space of all measurable functions  $f : \mathcal{C} \rightarrow \mathbb{R}$  such that  
 579  $\|f\|_p = (\int_{\mathcal{C}} |f|^p d\mu_0)^{1/p} < \infty$ . Moreover, we define  $L^\infty(\mathcal{C}, \mu_0)$  as the space of all essentially  
 580 bounded measurable functions from  $\mathcal{C}$  to  $\mathbb{R}$ . Unless stated otherwise, we assume the probability  
 581 measures are absolutely continuous w.r.t.  $\mu_0$ , and their density functions are in  $L^1(\mathcal{C}, \mu_0)$ . We  
 582 may abuse the notation and use the same symbol for a probability measure and its Radon–Nikodym  
 583 derivative w.r.t.  $\mu_0$ . Finally, we use  $\mathbb{E}[\cdot]$  to denote expectation under the probability measure  $\mu_0$ .

### 584 A.2 Useful Lemmas

585 Here, we state and prove a set of results that will be useful for the rest of this section. The first one  
 586 is Fenchel’s duality theorem:

587 **Lemma 4** (Fenchel’s Duality [58]). *Let  $X$  and  $Y$  be Banach spaces, let  $f : X \rightarrow \mathbb{R} \cup \{+\infty\}$  and  
 588  $g : Y \rightarrow \mathbb{R} \cup \{+\infty\}$  be convex functions and let  $A : X \rightarrow Y$  be a bounded linear map. Define the  
 589 primal and dual values  $p, d \in [-\infty, +\infty]$  by the Fenchel problems*

$$p = \inf_{x \in X} f(x) + g(Ax)$$

$$d = \sup_{y^* \in Y^*} -f^*(A^*y^*) - g^*(-y^*),$$

590 where  $f^*$  and  $g^*$  are the Fenchel conjugates of  $f$  and  $g$  defined as  $f^*(x^*) = \sup_{x \in X} \langle x^*, x \rangle - f(x)$   
 591 (similarly for  $g$ ),  $X^*$  is the dual space of  $X$  and  $\langle \cdot, \cdot \rangle$  is its duality pairing, and  $A^* : Y^* \rightarrow X^*$  is  
 592 the adjoint operator of  $A$ , i.e.,  $\langle A^*y^*, x \rangle = \langle y^*, Ax \rangle$ . Suppose  $A \text{ dom}(f) \cap \text{cont}(g) \neq \emptyset$ , where  
 593  $\text{dom}(f) := \{x \in X; f(x) < \infty\}$  and  $\text{cont}(g)$  are the continuous points of  $g$ . Then, strong duality  
 594 holds, i.e.,  $p = d$ .

595 *Proof.* See the proof of Theorem 4.4.3 in Borwein and Zhu [58]. □

596 We can use Fenchel’s duality to solve generalized maximum entropy problems. In particular, we  
 597 prove a generalization of Theorem 2 in [47] for density functions in  $L^1(\mathcal{C}, \mu_0)$ :

598 **Lemma 5.** *For any function  $\mu \in L^1(\mathcal{C}, \mu_0)$ , define the extended KL divergence as*

$$\psi(\mu) := \begin{cases} D_{\text{KL}}(\mu \parallel \mu_0) & \text{If } \|\mu\|_1 = 1, \\ +\infty & \text{o.w.} \end{cases}$$

599 *Moreover, assume a set of bounded feature functions  $m_1, m_2, \dots, m_N : \mathcal{C} \rightarrow \mathbb{R}$  is given and denote*  
 600  **$\mathbf{m}$**  *as the vector of all  $N$  features. Consider the linear function  $A_{\mathbf{m}} : L^1(\mathcal{C}, \mu_0) \rightarrow \mathbb{R}^N$  defined as*

$$\forall \mu \in L^1(\mathcal{C}, \mu_0) : A_{\mathbf{m}}(\mu) := (\mathbb{E}[m_1 \cdot \mu], \mathbb{E}[m_2 \cdot \mu], \dots, \mathbb{E}[m_N \cdot \mu]).$$

601 *We define the generalized maximum entropy problem as the following:*

$$\inf_{\mu \in L^1(\mathcal{C}, \mu_0)} \psi(\mu) + \zeta(A_{\mathbf{m}}(\mu)), \quad (5)$$

602 *for an arbitrary closed proper convex function  $\zeta : \mathbb{R}^N \rightarrow \mathbb{R}$ . Then the following holds:*

603 1. *The dual optimization of (5) is given by*

$$\sup_{\alpha \in \mathbb{R}^N} -\log \mathbb{E}[\exp\{\mathbf{m}^\top \alpha\}] - \zeta^*(-\alpha), \quad (6)$$

604 *where  $\zeta^*$  is the convex conjugate function of  $\zeta$ .*

605 2. *Denote  $\alpha^1, \alpha^2, \dots$  as a sequence in  $\mathbb{R}^N$  converging to supremum (6), and define the fol-*  
 606 *lowing Gibbs density functions*

$$\mu_{\text{Gibbs}}^\alpha(c) := \frac{\exp\{\mathbf{m}(c)^\top \alpha\}}{\mathbb{E}[\exp\{\mathbf{m}^\top \alpha\}]}.$$

607 *Then,*

$$\inf_{\mu \in L^1(\mathcal{C}, \mu_0)} \psi(\mu) + \zeta(A_{\mathbf{m}}(\mu)) = \lim_{n \rightarrow \infty} \psi(\mu_{\text{Gibbs}}^{\alpha^n}) + \zeta(A_{\mathbf{m}}(\mu_{\text{Gibbs}}^{\alpha^n})).$$

608 *Proof. Part 1:* We first derive the convex conjugate of  $\psi$ . Note that  $(L^1(\mathcal{C}, \mu_0))^* = L^\infty(\mathcal{C}, \mu_0)$   
 609 with the pairing

$$\forall h \in L^\infty(\mathcal{C}, \mu_0), \mu \in L^1(\mathcal{C}, \mu_0) : \langle h, \mu \rangle := \int_{\mathcal{C}} h(c) \cdot \mu(c) \, d\mu_0.$$

610 Hence, by Donsker and Varadhan's variational formula

$$\forall h \in L^\infty(\mathcal{C}, \mu_0) : \psi^*(h) = \sup_{\mu \in L^1(\mathcal{C}, \mu_0)} \langle h, \mu \rangle - \psi(\mu) = \log \mathbb{E} [\exp \{h\}]. \quad (7)$$

611 Moreover, the adjoint operator of  $A_{\mathbf{m}}$  is given by  $A_{\mathbf{m}}^* : \mathbb{R}^N \rightarrow (\mathcal{C} \rightarrow \mathbb{R})$ :

$$\forall \alpha \in \mathbb{R}^N, c \in \mathcal{C} : A_{\mathbf{m}}^*(\alpha)(c) = \mathbf{m}(c)^\top \alpha. \quad (8)$$

612 Using (7) and (8) and Lemma 4 concludes the proof.

613 **Part 2:** Denote the primal and dual objective functions by

$$\begin{aligned} P(\mu) &:= \psi(\mu) + \zeta(A_{\mathbf{m}}(\mu)), \\ D(\alpha) &:= -\log \mathbb{E} [\exp \{\mathbf{m}^\top \alpha\}] - \zeta^*(-\alpha), \end{aligned}$$

614 and their optimal values as  $P^*$  and  $D^*$ . For any  $\nu \in L^1(\mathcal{C}, \mu_0)$ , note that

$$\begin{aligned} \text{D}_{\text{KL}}(\nu \parallel \mu_0) - \text{D}_{\text{KL}}(\nu \parallel \mu_{\text{Gibbs}}^\alpha) &= \int_{\mathcal{C}} \nu \log \nu \, d\mu_0 - \left( \int_{\mathcal{C}} \nu \log \nu \, d\mu_0 - \int_{\mathcal{C}} \nu \log \mu_{\text{Gibbs}}^\alpha \, d\mu_0 \right) \\ &= \int_{\mathcal{C}} (\mathbf{m}(c)^\top \alpha) \nu(c) \, d\mu_0 - \log \mathbb{E} [\exp \{\mathbf{m}^\top \alpha\}] \\ &= A_{\mathbf{m}}(\nu)^\top \alpha - \log \mathbb{E} [\exp \{\mathbf{m}^\top \alpha\}]. \end{aligned} \quad (9)$$

615 Using (9), we can re-write the dual objective function as:

$$\forall \alpha \in \mathbb{R}^N, \nu \in L^1(\mathcal{C}, \mu_0) : D(\alpha) = -\text{D}_{\text{KL}}(\nu \parallel \mu_{\text{Gibbs}}^\alpha) + \text{D}_{\text{KL}}(\nu \parallel \mu_0) - A_{\mathbf{m}}(\nu)^\top \alpha - \zeta^*(-\alpha). \quad (10)$$

616 Moreover, note that

$$\begin{aligned} -A_{\mathbf{m}}(\nu)^\top \alpha - \zeta^*(-\alpha) &= -A_{\mathbf{m}}(\nu)^\top \alpha - \left( \sup_x \langle x, -\alpha \rangle - \zeta(x) \right) \\ &\leq -A_{\mathbf{m}}(\nu)^\top \alpha - \left( \langle A_{\mathbf{m}}(\nu), -\alpha \rangle - \zeta(A_{\mathbf{m}}(\nu)) \right) \\ &= \zeta(A_{\mathbf{m}}(\nu)). \end{aligned} \quad (11)$$

617 Combining (10) and (11), we get

$$\forall \alpha \in \mathbb{R}^N, \nu \in L^1(\mathcal{C}, \mu_0) : D(\alpha) \leq -\text{D}_{\text{KL}}(\nu \parallel \mu_{\text{Gibbs}}^\alpha) + \text{D}_{\text{KL}}(\nu \parallel \mu_0) + \zeta(A_{\mathbf{m}}(\nu)) \\ = -\text{D}_{\text{KL}}(\nu \parallel \mu_{\text{Gibbs}}^\alpha) + P(\nu). \quad (12)$$

618 Now, fix an arbitrary  $\epsilon > 0$ , and consider a sequence of  $\mu^1, \mu^2, \dots \in L^1(\mathcal{C}, \mu_0)$  such that for all  
 619  $j \in \mathbb{N}$ :

$$P(\mu^j) - P^* < \frac{\epsilon}{2^j}. \quad (13)$$

620 We can re-write (13) using the fact  $P^* = D^* = \lim_{n \rightarrow \infty} D(\alpha^n)$ :

$$\forall j \in \mathbb{N} : \lim_{n \rightarrow \infty} P(\mu^j) - D(\alpha^n) < \frac{\epsilon}{2^j} \quad (14)$$

621 In particular, by setting  $\nu = \mu^j$  in (12) and combining the result with (14), we get

$$\forall j \in \mathbb{N} : \lim_{n \rightarrow \infty} \text{D}_{\text{KL}}(\mu^j \parallel \mu_{\text{Gibbs}}^{\alpha^n}) < \frac{\epsilon}{2^j}.$$

622 Hence,  $\lim_{j \in \infty} \lim_{n \rightarrow \infty} \text{D}_{\text{KL}}(\mu^j \parallel \mu_{\text{Gibbs}}^{\alpha^n}) = 0$ . From properties of the KL divergence, it follows  
 623 that  $\lim_{j \rightarrow \infty} P(\mu^j) = \lim_{n \rightarrow \infty} P(\mu_{\text{Gibbs}}^{\alpha^n})$ , concluding the proof.  $\square$

624 **A.3 Max-Entropy Prior**

625 **Proposition 1.** Let  $N = |\mathcal{D}_E|$  be the number of demonstrations in  $\mathcal{D}_E$ . For each  $c \in \mathcal{C}$  and demon-  
 626 stration  $\tau_E = (s_1, a_1, s_2, a_2, \dots, s_H, a_H, s_{H+1}) \in \mathcal{D}_E$ , define  $m_{\tau_E}(c)$  as the (partial) likelihood of  
 627  $\tau_E$  under  $c$ :

$$m_{\tau_E}(c) = \prod_{h=1}^H p_E(a_h | s_h; c) \mathcal{T}(s_{h+1} | s_h, a_h, c). \quad (15)$$

628 Denote  $\mathbf{m}(c) \in \mathbb{R}^N$  as the vector with elements  $m_{\tau_E}(c)$  for  $\tau_E \in \mathcal{D}_E$ . Moreover, let  $\lambda^* \in \mathbb{R}^{\geq 0}$  be  
 629 the optimal solution to the Lagrange dual problem of (2). Then, the solution to optimization (2) is  
 630 as follows:

$$\mu_{ME}(c) = \lim_{n \rightarrow \infty} \frac{\exp\{\mathbf{m}(c)^\top \boldsymbol{\alpha}_n\}}{\mathbb{E}_{c \sim \mu_0} [\exp\{\mathbf{m}(c)^\top \boldsymbol{\alpha}_n\}]},$$

631 where  $\{\boldsymbol{\alpha}_n\}_{n=1}^\infty$  is a sequence converging to the following supremum:

$$\sup_{\boldsymbol{\alpha} \in \mathbb{R}^N} -\log \mathbb{E}_{c \sim \mu_0} [\exp\{\mathbf{m}(c)^\top \boldsymbol{\alpha}\}] + \frac{\lambda^*}{N} \sum_{i=1}^N \log \left( \frac{N \cdot \alpha_i}{\lambda^*} \right).$$

632 *Proof.* We first simplify the KL-divergence between the empirical distribution of the expert trajec-  
 633 tories  $\hat{\mathbb{P}}_E$  and the marginal likelihood  $\mathbb{P}_E(\cdot; \mu)$ :

$$\begin{aligned} \text{D}_{\text{KL}}(\hat{\mathbb{P}}_E \parallel \mathbb{P}_E(\cdot; \mu)) &= \sum_{\tau^{(i)} \in \mathcal{D}_E} \hat{\mathbb{P}}_E(\tau^{(i)}) \log \frac{\hat{\mathbb{P}}_E(\tau^{(i)})}{\mathbb{P}_E(\tau^{(i)}; \mu)} \\ &= -\log N - \frac{1}{N} \sum_{\tau^{(i)} \in \mathcal{D}_E} \log \mathbb{P}_E(\tau^{(i)}; \mu) \quad (\hat{\mathbb{P}}_E(\tau^{(i)}) = \frac{1}{N}) \\ &= -\log N - \frac{1}{N} \sum_{\tau^{(i)} \in \mathcal{D}_E} \log \mathbb{E}[m_{\tau^{(i)}} \cdot \mu] - \frac{1}{N} \sum_{s_1^{(i)} \in \mathcal{D}_E} \log \rho(s_1^{(i)}). \end{aligned}$$

By (1) and (15)

634 Using the above equality, we can re-write the definition of uncertainty set  $\mathcal{P}(\epsilon)$  as

$$\mathcal{P}(\epsilon) = \left\{ \mu; -\frac{1}{N} \sum_{\tau \in \mathcal{D}_E} \log \mathbb{E}[m_\tau \cdot \mu] - \epsilon - \log N - \frac{1}{N} \sum_{s_1 \in \mathcal{D}_E} \log \rho(s_1) \leq 0 \right\}.$$

635 Therefore, we can re-write the optimization (2) as

$$\mu_{ME} = \arg \min_{\mu \in L^1(\mathcal{C}, \mu_0)} \psi(\mu) \quad \text{s.t.} \quad -\frac{1}{N} \sum_{\tau \in \mathcal{D}_E} \log \mathbb{E}[m_\tau \cdot \mu] - \epsilon - \log N - \frac{1}{N} \sum_{s_1 \in \mathcal{D}_E} \log \rho(s_1) \leq 0, \quad (16)$$

636 where the extended KL divergence  $\psi(\mu)$  is defined as:

$$\psi(\mu) := \begin{cases} \text{D}_{\text{KL}}(\mu \parallel \mu_0) & \text{If } \|\mu\|_1 = 1, \\ +\infty & \text{o.w.} \end{cases}$$

637 Note that  $\mathcal{P}(\epsilon)$  is a convex set. To see this, consider  $\mu_1, \mu_2 \in \mathcal{P}(\epsilon)$ . Then, for any  $0 \leq \lambda \leq 1$ , we  
 638 have  $\mu = (1 - \lambda)\mu_1 + \lambda\mu_2 \in \mathcal{P}(\epsilon)$  since  $\mathbb{E}[m_\tau \cdot \mu]$  is linear in  $\mu$  and  $-\log$  is convex. Moreover, It  
 639 is easy to see there exists a strictly feasible solution for (16) (e.g., consider the true distribution  $\mu^*$   
 640 over  $\mathcal{C}$ ). Thus, strong duality holds, and we can form the Lagrangian function as

$$L(\mu, \lambda) := \psi(\mu) + \lambda \left( \frac{1}{N} \sum_{\tau \in \mathcal{D}_E} -\log \mathbb{E}[m_\tau \cdot \mu] \right) - \lambda \left( \epsilon + \log N + \frac{1}{N} \sum_{s_1 \in \mathcal{D}_E} \log \rho(s_1) \right).$$

641 Given that  $\lambda^* \in \mathbb{R}^{\geq 0}$  is the optimal solution to the Lagrange dual problem, the maximum entropy  
 642 prior  $\mu_{\text{ME}}$  will be the solution to

$$\inf_{\mu \in L^1(\mathcal{C}, \mu_0)} L(\mu, \lambda^*) = \inf_{\mu \in L^1(\mathcal{C}, \mu_0)} \psi(\mu) + \lambda^* \left( \frac{1}{N} \sum_{\tau \in \mathcal{D}_E} -\log \mathbb{E}[m_\tau \cdot \mu] \right) + \text{constant in } \mu. \quad (17)$$

643 Now, for each  $\mathbf{x} \in \mathbb{R}^N$ , define the convex function  $\zeta(\mathbf{x}) := \frac{\lambda^*}{N} \left( \sum_{i=1}^N -\log x_i \right)$ . Moreover, for  
 644  $\mu \in L^1(\mathcal{C}, \mu_0)$ , define  $A_{\mathbf{m}}(\mu) := (\mathbb{E}[m_{\tau(1)} \cdot \mu], \mathbb{E}[m_{\tau(2)} \cdot \mu], \dots, \mathbb{E}[m_{\tau(N)} \cdot \mu])$ . Then,

$$L(\mu, \lambda^*) = \psi(\mu) + \zeta(A_{\mathbf{m}}(\mu)). \quad (18)$$

645 Combining (17) and (18), the maximum entropy prior  $\mu_{\text{ME}}$  is the solution to

$$\inf_{\mu \in L^1(\mathcal{C}, \mu_0)} \psi(\mu) + \zeta(A_{\mathbf{m}}(\mu)).$$

646 Using Lemma 5 and noting that

$$\zeta^*(x^*) = \frac{\lambda^*}{N} \left( \sum_{i=1}^N -1 - \log \left( -\frac{N}{\lambda^*} \cdot x_i^* \right) \right)$$

647 concludes the proof.  $\square$

#### 648 A.4 $K$ -armed Bandit Frequentist Algorithm & Regret

649 To simplify the analysis, we employ a deterministic sampling approach by pulling each arm a fixed  
 650 number of times based on its probability. To do so, we discretize the expert policy with a step size  
 651  $p_{\min}$ , which leads to a relative frequency of  $\lceil \frac{\hat{P}_E(a)}{p_{\min}} \rceil$  for an arm  $a$ . In particular, we can choose  
 652  $p_{\min} = \min_{a \in \mathcal{A}} \hat{P}_E(a)$ .

---

#### Algorithm 1 Successive Elimination with Expert Sampling

---

- 1: **Input:** Episodes  $T$ , Arms  $\mathcal{A} = [K]$ , expert policy  $\hat{P}_E$ , step size  $p_{\min}$ , an unknown task  $c \sim \mu^*$ , and  $\delta \in (0, 1)$ .
  - 2: **for**  $t = 1 \dots T$  **do**
  - 3: Try an active arm  $a$  with a relative frequency of  $\lceil \frac{\hat{P}_E(a)}{p_{\min}} \rceil$ . // all arms are active at  $t = 0$ .  
 //  $n_t(a)$  is the number of times that an arm  $a$  is pulled by episode  $t$  and  $\bar{V}_c^t(a)$  is its empirical mean reward.
  - 4: Increment  $n_t(a)$  and update  $\bar{V}_c^t(a)$ .
  - 5: Construct  $\text{UCB}_a^t = \bar{V}_c^t(a) + \sqrt{\frac{\log(4T^4 K/\delta)}{2n_t(a)}}$  and  $\text{LCB}_a^t = \bar{V}_c^t(a) - \sqrt{\frac{\log(4T^4 K/\delta)}{2n_t(a)}}$ .
  - 6: De-activate all arms  $a$  s.t.  $\exists a'$  with  $\text{UCB}_a \leq \text{LCB}_{a'}$ , and normalize  $\hat{P}_E$ .
  - 7: **end for**
- 

653 **Theorem 2.** Assume that (i) the mean value of reward function  $R$  is bounded in  $[0, 1]$  for all arms,  
 654 (ii)  $T \cdot p_{\min} \geq 1$ , (iii) the expert is optimal, i.e.,  $\hat{P}_E = P_E(\cdot; \mu^*)$  ( $\beta \rightarrow \infty, |\mathcal{D}_E| \rightarrow \infty$ ), and (iv) the  
 655 learner follows Algorithm 1. Then, with probability at least  $1 - \delta$ ,

$$\text{Reg} \lesssim \sqrt{T \log(TK/\delta)} \sum_{a, a' \in \mathcal{A}; a \neq a'} \sqrt{\frac{\hat{P}_E(a)}{\hat{P}_E(a) + \hat{P}_E(a')} \cdot \left( 1 - \frac{\hat{P}_E(a)}{\hat{P}_E(a) + \hat{P}_E(a')} \right)} \left[ \sqrt{\hat{P}_E(a)} + \sqrt{\hat{P}_E(a')} \right].$$

656 *Proof.* Fix  $\delta \in (0, 1)$  and  $c \in \mathcal{C}$ . Let  $\mathcal{E}$  be the event that  $|\bar{V}_c^t(a) - V_c(a)| \leq \sqrt{\frac{\log(4T^4 K/\delta)}{2n_t(a)}}$  for all  
 657 arms  $a \in \mathcal{A}$ , all  $t \leq T$ , and all  $T \in \mathbb{N}$ , where  $n_t(a)$  is the number of times that arm  $a$  was pulled by  
 658 time  $t$ . Note that since  $T \geq \frac{1}{p_{\min}}$ , each arm will be pulled at least once and  $n_t(a) \geq 1$ .

659 We first show that  $\mathbb{P}(\mathcal{E}) \geq 1 - \delta$ . Fix  $T$ , arm  $a$ , and  $t \leq T$ . Suppose  $n_t(a) = j$  for  $1 \leq j \leq T$ . By  
660 Hoeffding's inequality, we have

$$\mathbb{P}\left(\left|\overline{V}_c^t(a) - V_c(a)\right| \leq \sqrt{\frac{\log(4T^4K/\delta)}{2j}}\right) \geq 1 - \frac{\delta}{2T^4K}. \quad (19)$$

661 Now, using the union bound over all episodes and all actions, we get

$$\begin{aligned} & \mathbb{P}\left(\exists a \in \mathcal{A}, T \in \mathbb{N}, t \leq T, j \leq t: \left|\overline{V}_c^t(a) - V_c(a)\right| > \sqrt{\frac{\log(2T^4K/\delta)}{2j}}\right) \\ & \leq \sum_{T=1}^{\infty} \sum_{a \in \mathcal{A}} \sum_{t=1}^T \sum_{j=1}^t \mathbb{P}\left(\left|\overline{V}_c^t(a) - V_c(a)\right| > \sqrt{\frac{\log(2T^4K/\delta)}{2j}}\right) \\ & \leq \sum_{T=1}^{\infty} \sum_{a \in \mathcal{A}} \sum_{t=1}^T t \cdot \frac{\delta}{2T^4K} \quad \text{By (19)} \\ & \leq \sum_{T=1}^{\infty} \frac{\delta}{2T^4K} \times T^2 \times K = \sum_{T=1}^{\infty} \frac{\delta}{2T^2} \leq \delta, \end{aligned}$$

662 which concludes that  $\mathbb{P}(\mathcal{E}) \geq 1 - \delta$ .

663 The rest of the proof computes the regret for when  $\mathcal{E}$  holds. For simplicity and without loss of  
664 generality, we assume all expert probabilities are dividable by  $p_{\min}$ . Recall that we follow a de-  
665 terministic sampling approach and choose each arm according to its relative frequency  $\frac{\widehat{P}_E(\cdot)}{p_{\min}}$  for  
666 multiple batches, where each batch loops over all active actions. Let  $t_a$  be the episode in which we  
667 eliminate an arm  $a$  in favour of another arm. Then, it is easy to show that

$$\forall a' \in \text{active arms by } t_a: \widehat{P}_E(a') \cdot t_a \leq n_{t_a}(a'), \quad (20)$$

668 This lower bound corresponds to the case where no other arm is eliminated before eliminating  $a$ .  
669 Moreover, we have an upper bound for  $n_{t_a}(a)$  considering the worst-case scenario in which the only  
670 remaining arms are  $a$  and  $a_c$ , where  $a_c$  is the optimal action for task  $c$ :

$$n_{t_a}(a) \leq \frac{\widehat{P}_E(a)}{\widehat{P}_E(a) + \widehat{P}_E(a_c)} \cdot t_a. \quad (21)$$

671 Now, let  $\text{Reg}_c(a)$  be the total regret contributed by the arm  $a$  for a given task  $c \sim \mathcal{C}$ . We can upper  
672 bound the regret as

$$\begin{aligned} \text{Reg}_c(a) &= n_{t_a}(a) (V_c(a_c) - V_c(a)) \\ & \stackrel{(i)}{\leq} 2n_{t_a}(a) \left( \sqrt{\frac{\log(4T^4K/\delta)}{2n_{t_a}(a)}} + \sqrt{\frac{\log(4T^4K/\delta)}{2n_{t_a}(a_c)}} \right) \\ & \leq 2 \frac{\widehat{P}_E(a)}{\widehat{P}_E(a) + \widehat{P}_E(a_c)} \cdot t_a \left( \sqrt{\frac{\log(4T^4K/\delta)}{2n_{t_a}(a)}} + \sqrt{\frac{\log(4T^4K/\delta)}{2n_{t_a}(a_c)}} \right) \quad \text{By (21)} \\ & = \sqrt{2\log(4T^4K/\delta)} \cdot \frac{\widehat{P}_E(a)}{\widehat{P}_E(a) + \widehat{P}_E(a_c)} \cdot t_a \left( \sqrt{\frac{1}{n_{t_a}(a)}} + \sqrt{\frac{1}{n_{t_a}(a_c)}} \right) \\ & \leq \sqrt{2\log(4T^4K/\delta)} \cdot \frac{\widehat{P}_E(a)}{\widehat{P}_E(a) + \widehat{P}_E(a_c)} \cdot t_a \left( \sqrt{\frac{1}{t_a \widehat{P}_E(a)}} + \sqrt{\frac{1}{t_a \widehat{P}_E(a_c)}} \right) \quad \text{By (20)} \\ & = \sqrt{2t_a \log(4T^4K/\delta)} \cdot \frac{\widehat{P}_E(a)}{\widehat{P}_E(a) + \widehat{P}_E(a_c)} \left( \sqrt{\frac{1}{\widehat{P}_E(a)}} + \sqrt{\frac{1}{\widehat{P}_E(a_c)}} \right) \\ & \stackrel{(ii)}{\leq} \sqrt{2T \log(4T^4K/\delta)} \cdot \frac{\widehat{P}_E(a)}{\widehat{P}_E(a) + \widehat{P}_E(a_c)} \left( \sqrt{\frac{1}{\widehat{P}_E(a)}} + \sqrt{\frac{1}{\widehat{P}_E(a_c)}} \right), \end{aligned}$$



673 where (i) holds since the confidence intervals of arm  $a$  and  $a_c$  overlap at episode  $t_a$  (otherwise,  $a$   
674 would have been eliminated before  $t_a$ ), and (ii) follows from the fact that  $t_a \leq T$ .

675 Finally, we upper bound the Bayesian regret by taking the expectation of  $\sum_{a \neq a_c} \text{Reg}_c(a)$  over  $c \sim$   
676  $\mathcal{C}$ . Note that since the expert is optimal, we have  $\hat{P}_E(a) = \mu^*(a_c = a)$  for all  $k \in \mathcal{A}$ .

$$\begin{aligned}
\text{Reg} &= \mathbb{E}_{c \sim \mu^*} \left[ \sum_{a \neq a_c} \text{Reg}_c(a) \right] \\
&\stackrel{(i)}{\leq} \sum_{a' \in \mathcal{A}} \mu^*(a_c = a') \left( \max_{c; a_c = a'} \sum_{a \neq a'} \text{Reg}_c(a) \right) \\
&= \sum_{a' \in \mathcal{A}} \hat{P}_E(a') \left( \max_{c; a_c = a'} \sum_{a \neq a'} \text{Reg}_c(a) \right) \\
&\leq \sqrt{2T \log(4T^4 K / \delta)} \sum_{a' \in \mathcal{A}} \sum_{a \neq a'} \frac{\hat{P}_E(a') \hat{P}_E(a)}{\hat{P}_E(a) + \hat{P}_E(a')} \left( \sqrt{\frac{1}{\hat{P}_E(a)}} + \sqrt{\frac{1}{\hat{P}_E(a')}} \right) \\
&\stackrel{(ii)}{\leq} \sqrt{8T \log(4TK / \delta)} \sum_{a, a' \in \mathcal{A}; a \neq a'} \frac{\hat{P}_E(a') \hat{P}_E(a)}{\hat{P}_E(a) + \hat{P}_E(a')} \left( \sqrt{\frac{1}{\hat{P}_E(a)}} + \sqrt{\frac{1}{\hat{P}_E(a')}} \right) \\
&= \sqrt{8T \log(4TK / \delta)} \sum_{a, a' \in \mathcal{A}; a \neq a'} \sqrt{\frac{\hat{P}_E(a')}{\hat{P}_E(a) + \hat{P}_E(a')} \cdot \frac{\hat{P}_E(a)}{\hat{P}_E(a) + \hat{P}_E(a')}} \left( \sqrt{\hat{P}_E(a)} + \sqrt{\hat{P}_E(a')} \right)
\end{aligned}$$

677 where (i) follows by partitioning  $\mathcal{C}$  into  $\{c; c \in \mathcal{C}, a_c = a'\}_{a' \in \mathcal{A}}$  and choosing the worst-case task  
678 in each partition, and (ii) holds since  $4K/\delta > 1$ . Replacing  $\frac{\hat{P}_E(a)}{\hat{P}_E(a) + \hat{P}_E(a_c)}$  with  $1 - \frac{\hat{P}_E(a')}{\hat{P}_E(a) + \hat{P}_E(a')}$   
679 concludes the proof.  $\square$

## 680 A.5 Max-Entropy Expert Posterior for MDPs

681 **Proposition 6** (Max-Entropy Expert Posterior for MDPs). *Consider a contextual MDP  $\mathcal{M} =$   
682  $(\mathcal{S}, \mathcal{A}, \mathcal{T}, R, H, \rho, \mu^*)$ . Assume the transition function  $\mathcal{T}$  does not depend on the task variables.  
683 Moreover, assume the reward distribution is Gaussian with unit variance and Assumption 1 holds.  
684 Then, the log-pdf posterior function under the maximum entropy prior is given as:*

$$\begin{aligned}
\forall \theta \in \Theta : \log \mu_{ME}(\theta | \mathcal{H}_T) &= - \sum_{t=1}^T \sum_{h=1}^H \frac{1}{2} \left( r_h^t + \max_{a' \in \mathcal{A}} \mathbb{E}_{s'} [Q(s', a'; \theta)] - Q(s_h^t, a_h^t; \theta) \right)^2 \\
&\quad + \sum_{\tau \in \mathcal{D}_E} \alpha_\tau^* \cdot \prod_{(s, a) \in \tau} \frac{\exp\{\beta \cdot Q(s, a; \theta)\}}{\sum_{a' \in \mathcal{A}} \exp\{\beta \cdot Q(s, a'; \theta)\}} + \text{constant in } \theta,
\end{aligned} \tag{22}$$

685 where  $\mathcal{H}_T = \left\{ \left( (s_h^t, a_h^t, r_h^t, s_{h+1}^t)_{h=1}^H \right)_{t=1}^T \right\}$  is the history of online interactions,  $\mathcal{D}_E$  is the expert  
686 demonstration data,  $\beta$  is the competence level of the expert in Assumption 1, and  $\{\alpha_\tau^*\}_{\tau \in \mathcal{D}_E}$  are  
687 derived from Proposition 1.

688 **Remark.** We note that, in principle, the ExPerior framework allows for task-dependent transition  
689 functions. In this case, the log-pdf in (22) provides an optimistic upper bound on the true posterior  
690 log-pdf function. See Hao et al. [23] for a similar analysis. We leave the general case for future  
691 work. Note that the second term of (22) is simply the log-pdf of the max-entropy prior.

692 *Proof.* Since the transition function is task-independent, the likelihood of an expert trajectory  $\tau_E$  can  
693 be simplified as:

$$\forall c \in \mathcal{C} : \quad m_{\tau_E}(c) = \prod_{h=1}^H p_E(a_h | s_h; c) \cdot \prod_{h=1}^H \mathcal{T}(s_{h+1} | s_h, a_h). \tag{23}$$

694 The second term in (23) is constant in  $c$ . This implies that the likelihood function  $m_{\tau_E}(c)$  will  
695 depend on  $c$  only through the expert policy, which itself is a function of optimal Q-functions by  
696 Assumption 1. Note that the second term in the definition of  $m_{\tau_E}$  can be simply removed since we  
697 can re-weight the parameters  $\alpha$  in the optimization step (3) of Proposition 1. Hence, assuming the  
698 deep Q-network is expressive enough, without loss of generality, we can re-define the likelihood  
699 function of an expert trajectory  $\tau_E = (s_1, a_1, s_2, a_2, \dots, s_H, a_H, s_{H+1})$  as

$$\forall \boldsymbol{\theta} \in \Theta : \quad m_{\tau_E}(\boldsymbol{\theta}) = \prod_{h=1}^H \frac{\exp\{\beta \cdot Q(s_h, a_h; \boldsymbol{\theta})\}}{\sum_{a' \in \mathcal{A}} \exp\{\beta \cdot Q(s_h, a'; \boldsymbol{\theta})\}}.$$

700 We can now write the log-pdf of the posterior distribution of  $\boldsymbol{\theta}$  given  $\mathcal{H}_T$ :

$$\begin{aligned} \forall \boldsymbol{\theta} \in \Theta : \quad & \log \mu_{\text{ME}}(\boldsymbol{\theta} \mid \mathcal{H}_T) \\ &= \log P(\mathcal{H}_T \mid \boldsymbol{\theta}) + \log \mu_{\text{ME}}(\boldsymbol{\theta}) + \text{constant in } \boldsymbol{\theta} \\ &= \sum_{t=1}^L \sum_{h=1}^H \log \rho(s_1^t) + \log R(r_h^t \mid s_h^t, a_h^t; \boldsymbol{\theta}) + \log \mathcal{T}(s_{h+1}^t \mid s_h^t, a_h^t) + \log \mu_{\text{ME}}(\boldsymbol{\theta}) + \text{const.} \\ &= \sum_{t=1}^L \sum_{h=1}^H \log R(r_h^t \mid s_h^t, a_h^t; \boldsymbol{\theta}) + \log \mu_{\text{ME}}(\boldsymbol{\theta}) + \text{const.}, \end{aligned} \quad (24)$$

701 Now, given the Bellman equations, we can write the mean value of the reward function as

$$\forall s \in \mathcal{S}, a \in \mathcal{A} : \quad \mathbb{E}[R(s, a; \boldsymbol{\theta})] = Q(s, a; \boldsymbol{\theta}) - \max_{a' \in \mathcal{A}} \mathbb{E}_{s'}[Q(s', a'; \boldsymbol{\theta})]$$

702 The reward distribution is Gaussian with unit variance. Therefore,

$$\forall s \in \mathcal{S}, a \in \mathcal{A}, r \in \mathbb{R} : \quad R(r \mid s, a; \boldsymbol{\theta}) = \mathcal{N}\left(Q(s, a; \boldsymbol{\theta}) - \max_{a' \in \mathcal{A}} \mathbb{E}_{s'}[Q(s', a'; \boldsymbol{\theta})], 1\right). \quad (25)$$

703 Moreover, by Proposition 1, the log-pdf of the maximum entropy expert prior is given as

$$\forall \boldsymbol{\theta} \in \Theta : \quad \log \mu_{\text{ME}}(\boldsymbol{\theta}) = \sum_{\tau \in \mathcal{D}_E} \alpha_{\tau}^* \cdot m_{\tau}(\boldsymbol{\theta}) = \sum_{\tau \in \mathcal{D}_E} \alpha_{\tau}^* \cdot \prod_{(s,a) \in \tau} \frac{\exp\{\beta \cdot Q(s, a; \boldsymbol{\theta})\}}{\sum_{a' \in \mathcal{A}} \exp\{\beta \cdot Q(s, a'; \boldsymbol{\theta})\}}. \quad (26)$$

704 Combining (24) to (26), we conclude the proof.  $\square$

## 705 A.6 Ensemble Marginal Likelihood

706 **Proposition 3.** Consider a contextual MDP  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, R, H, \rho, \mu^*)$ . Assume the transition  
707 function  $\mathcal{T}$  does not depend on the task variables and Assumption 1 holds. Then, the negative  
708 marginal log-likelihood of expert data  $\mathcal{D}_E$  under the ensemble prior  $\mu_{\boldsymbol{\theta}_{\text{ens}}}$  is upper bounded by

$$-\log P_E(\mathcal{D}_E; \mu_{\boldsymbol{\theta}_{\text{ens}}}) \leq \frac{1}{L} \sum_{i=1}^L \sum_{\tau \in \mathcal{D}_E} \sum_{(s,a) \in \tau} \log \left( \sum_{a' \in \mathcal{A}} \exp\{\beta \cdot Q(s, a'; \boldsymbol{\theta}_{\text{ens}}^i)\} \right) - \beta \cdot Q(s, a; \boldsymbol{\theta}_{\text{ens}}^i),$$

709 where  $\beta$  is the competence level of the expert in Assumption 1.

710 *Proof.* Recalling (1), the log-likelihood of the expert trajectories  $\mathcal{D}_E$  under  $\mu_{\theta_{\text{ens}}}$  is given by

$$\begin{aligned}
-\log P_E(\mathcal{D}_E; \mu_{\theta_{\text{ens}}}) &= \sum_{\tau^{(i)} \in \mathcal{D}_E} -\log \mathbb{E}_{\theta \sim \mu_{\theta_{\text{ens}}}} \left[ \rho(s_1^{(i)}) \prod_{h=1}^H p_E(a_h^{(i)} | s_h^{(i)}; \theta) \mathcal{T}(s_{h+1}^{(i)} | s_h^{(i)}, a_h^{(i)}) \right] \\
&= \sum_{\tau^{(i)} \in \mathcal{D}_E} -\log \mathbb{E}_{\theta \sim \mu_{\theta_{\text{ens}}}} \left[ \prod_{h=1}^H p_E(a_h^{(i)} | s_h^{(i)}; \theta) \right] + \text{constant in } \theta_{\text{ens}} \\
&\hspace{20em} (\rho, \mathcal{T} \text{ do not depend on } \theta) \\
&= \sum_{\tau^{(i)} \in \mathcal{D}_E} -\log \left( \frac{1}{L} \sum_{j=1}^L \prod_{h=1}^H p_E(a_h^{(i)} | s_h^{(i)}; \theta_{\text{ens}}^j) \right) \\
&\hspace{20em} (\text{By Definition of } \mu_{\theta_{\text{ens}}}) \\
&\leq \sum_{\tau^{(i)} \in \mathcal{D}_E} \frac{1}{L} \sum_{j=1}^L \sum_{h=1}^H -\log p_E(a_h^{(i)} | s_h^{(i)}; \theta_{\text{ens}}^j) \quad \text{By Jensen's inequality} \\
&= \frac{1}{L} \sum_{i=1}^L \sum_{\tau \in \mathcal{D}_E} \sum_{(s,a) \in \tau} \left[ \log \left( \sum_{a' \in \mathcal{A}} \exp \{ \beta \cdot Q(s, a'; \theta_{\text{ens}}^i) \} \right) - \beta \cdot Q(s, a; \theta_{\text{ens}}^i) \right] \\
&\hspace{20em} \text{By Assumption 1}
\end{aligned}$$

711

□

## 712 B High-Level Implementation of ExPerior

---

### Algorithm 2 Max-Entropy Posterior Sampling (ExPerior)

---

```

1: Input: Expert demonstrations  $\mathcal{D}_E$ , Reference distribution  $\mu_0, \lambda^* \geq 0$ , and unknown task  $c \sim \mu^*$ .
2:  $\mu_{\text{ME}} \leftarrow \text{MAXENTROPYEXPERTPRIOR}(\mu_0, \mathcal{D}_E, \lambda^*)$ 
3:  $history \leftarrow \{\}$ 
4: for episode  $t \leftarrow 1, 2, \dots$  do
5:   sample  $c_t \sim \mu_{\text{ME}}(\cdot | history)$  // posterior sampling
6:   for timestep  $h \leftarrow 1, 2, \dots, H$  do
7:     take action  $a_h^t \sim \pi_{c_t}(\cdot | s_h)$ 
8:     observe  $r_h^t \sim R(s_h^t, a_h^t, c)$ ,  $s_{h+1}^t \sim \mathcal{T}(s_h^t, a_h^t, c)$  and append  $(a_h^t, r_h^t, s_{h+1}^t)$  to  $history$ 
9:   end for
10: end for

```

---

713 **C Additional Experiments**

714 **C.1 Ablation Studies for Bernoulli Multi-Armed Bandits**

Table 2: Ablation experiments to assess the robustness of ExPerior to misspecified expert models. Random-optimal experts choose the optimal action with probability  $\gamma$  and choose random actions with probability  $1 - \gamma$ . ExPerior-MaxEnt achieves consistent out-performance by setting the hyperparameter  $\beta = 10$ , while ExPerior-Param get almost similar results for  $\beta = 1$  and  $\beta = 2.5$ .

	Optimal	Noisily-Rational				Random-Optimal			
		$\beta = 0.1$	$\beta = 1$	$\beta = 2.5$	$\beta = 10$	$\gamma = 0.0$	$\gamma = 0.25$	$\gamma = 0.5$	$\gamma = 0.75$
$\beta = 0.1$									
ExPerior-MaxEnt	51.7 $\pm$ 5.1	52.3 $\pm$ 5.3	52.3 $\pm$ 5.3	52.0 $\pm$ 5.1	51.7 $\pm$ 5.0	52.3 $\pm$ 5.3	52.1 $\pm$ 5.1	52.0 $\pm$ 5.1	51.8 $\pm$ 5.0
ExPerior-Param	11.1 $\pm$ 4.3	33.1 $\pm$ 7.3	<b>12.6 <math>\pm</math> 3.5</b>	11.7 $\pm$ 3.8	10.9 $\pm$ 4.2	40.1 $\pm$ 9.6	12.3 $\pm$ 4.7	11.4 $\pm$ 4.0	10.7 $\pm$ 4.2
$\beta = 1$									
ExPerior-MaxEnt	45.7 $\pm$ 3.4	52.2 $\pm$ 5.3	51.6 $\pm$ 5.1	50.0 $\pm$ 4.8	47.3 $\pm$ 3.8	52.5 $\pm$ 5.3	51.0 $\pm$ 4.8	49.1 $\pm$ 4.2	48.0 $\pm$ 3.6
ExPerior-Param	9.1 $\pm$ 3.0	<b>21.3 <math>\pm</math> 1.3</b>	13.4 $\pm$ 2.9	<b>10.1 <math>\pm</math> 3.0</b>	9.4 $\pm$ 3.1	<b>22.8 <math>\pm</math> 1.3</b>	<b>9.8 <math>\pm</math> 3.0</b>	<b>8.6 <math>\pm</math> 2.7</b>	<b>8.8 <math>\pm</math> 2.9</b>
$\beta = 2.5$									
ExPerior-MaxEnt	<b>37.0 <math>\pm</math> 1.9</b>	52.1 $\pm$ 5.3	51.0 $\pm$ 4.9	47.1 $\pm$ 4.5	38.3 $\pm$ 2.0	<b>52.1 <math>\pm</math> 5.1</b>	48.9 $\pm$ 4.1	44.8 $\pm$ 3.2	40.5 $\pm$ 2.1
ExPerior-Param	<b>8.5 <math>\pm</math> 2.8</b>	24.3 $\pm$ 1.2	19.0 $\pm$ 2.1	12.8 $\pm$ 2.9	<b>9.2 <math>\pm</math> 3.1</b>	24.6 $\pm$ 1.2	15.9 $\pm$ 3.0	10.9 $\pm$ 3.2	<b>8.8 <math>\pm</math> 2.9</b>
$\beta = 10$									
ExPerior-MaxEnt	38.5 $\pm$ 9.4	<b>52.0 <math>\pm</math> 5.2</b>	<b>47.6 <math>\pm</math> 4.4</b>	<b>39.7 <math>\pm</math> 2.9</b>	<b>29.7 <math>\pm</math> 3.6</b>	52.5 $\pm$ 5.3	<b>41.9 <math>\pm</math> 2.6</b>	<b>37.7 <math>\pm</math> 2.8</b>	<b>31.9 <math>\pm</math> 3.0</b>
ExPerior-Param	11.2 $\pm$ 4.8	26.9 $\pm$ 1.2	25.0 $\pm$ 1.5	21.0 $\pm$ 2.1	11.8 $\pm$ 3.3	26.8 $\pm$ 1.1	23.2 $\pm$ 1.8	20.1 $\pm$ 2.5	16.1 $\pm$ 3.0
Oracle-TS	8.5 $\pm$ 2.7	8.5 $\pm$ 2.7	8.5 $\pm$ 2.7	8.5 $\pm$ 2.7	8.5 $\pm$ 2.7	8.5 $\pm$ 2.7	8.5 $\pm$ 2.7	8.5 $\pm$ 2.7	8.5 $\pm$ 2.7
Oracle-TS (SGLD)	24.2 $\pm$ 3.9	24.2 $\pm$ 3.9	24.2 $\pm$ 3.9	24.2 $\pm$ 3.9	24.2 $\pm$ 3.9	24.2 $\pm$ 3.9	24.2 $\pm$ 3.9	24.2 $\pm$ 3.9	24.2 $\pm$ 3.9

Table 3: Superiority of ExPerior-MaxEnt compared to ExPerior-Param with misspecified parametric prior.

	Low Entropy	Mid-Entropy	High-Entropy
ExPerior-Param	0.7 $\pm$ 0.3	6.8 $\pm$ 0.8	24.5 $\pm$ 2.8
ExPerior-MaxEnt	11.6 $\pm$ 1.3	25.7 $\pm$ 1.2	41.3 $\pm$ 2.2
ExPerior-Param (Gamma)	39.3 $\pm$ 2.2	36.8 $\pm$ 0.9	51.8 $\pm$ 3.6
ExPerior-Param (Beta-SGLD)	60.2 $\pm$ 6.3	40.4 $\pm$ 2.0	45.6 $\pm$ 2.0
ExPerior-Param (Normal)	546.5 $\pm$ 153.4	492.5 $\pm$ 185.6	461.8 $\pm$ 104.8
Oracle-TS	0.9 $\pm$ 0.4	7.3 $\pm$ 0.8	21.5 $\pm$ 2.2
Oracle-TS (SGLD)	11.0 $\pm$ 1.6	21.2 $\pm$ 1.0	39.9 $\pm$ 3.2

715 **C.2 Frozen Lake**

Table 4: The average reward per episode in Frozen Lake (MDP) after 90,000 training steps.

	Fixed # Hazard = 9				Fixed $\beta = 1$			
	$\beta = 0.1$	$\beta = 1$	$\beta = 2.5$	$\beta = 10$	# Hazard = 2	# Hazard = 5	# Hazard = 7	# Hazard = 9
(MDP)								
ExPerior-MaxEnt	-23.36 $\pm$ 1.26	12.26 $\pm$ 0.29	12.68 $\pm$ 0.03	<b>12.71 <math>\pm</math> 0.03</b>	<b>13.02 <math>\pm</math> 0.18</b>	<b>12.78 <math>\pm</math> 0.11</b>	<b>12.78 <math>\pm</math> 0.06</b>	12.26 $\pm$ 0.29
ExPerior-Param	-25.53 $\pm$ 2.35	<b>12.64 <math>\pm</math> 0.08</b>	<b>12.70 <math>\pm</math> 0.03</b>	12.68 $\pm$ 0.03	13.00 $\pm$ 0.18	<b>12.78 <math>\pm</math> 0.12</b>	12.73 $\pm$ 0.07	<b>12.64 <math>\pm</math> 0.08</b>
Naïve Boot-DQN	-23.32 $\pm$ 0.69	-23.32 $\pm$ 0.69	-23.32 $\pm$ 0.69	-23.32 $\pm$ 0.69	-14.39 $\pm$ 5.22	-20.99 $\pm$ 2.86	-20.39 $\pm$ 1.75	-23.32 $\pm$ 0.69
ExPLOre	<b>11.74 <math>\pm</math> 0.41</b>	11.75 $\pm$ 0.63	11.96 $\pm$ 0.28	12.3 $\pm$ 0.22	-113.84 $\pm$ 17.50	-54.89 $\pm$ 13.75	-10.00 $\pm$ 7.60	11.75 $\pm$ 0.63
Optimal	12.71 $\pm$ 0.03	12.71 $\pm$ 0.03	12.71 $\pm$ 0.03	12.71 $\pm$ 0.03	13.02 $\pm$ 0.18	12.78 $\pm$ 0.11	12.76 $\pm$ 0.06	12.64 $\pm$ 0.03