

SOREL AND TOREL: TWO METHODS FOR FULLY OFF-LINE REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Sample efficiency remains a major obstacle for real world adoption of reinforcement learning (RL): success has been limited to settings where simulators provide access to essentially unlimited environment interactions, which in reality are typically costly or dangerous to obtain. Offline RL in principle offers a solution by exploiting offline data to learn a near-optimal policy before deployment. In practice, however, current offline RL methods rely on extensive online interactions for hyperparameter tuning, and have no reliable bound on their initial online performance. To address these two issues, we introduce two algorithms. Firstly, SOReL: an algorithm for **safe offline reinforcement learning**. *Using only offline data* our Bayesian approach infers a posterior over environment dynamics to obtain a reliable estimate of the online performance via the posterior predictive uncertainty. Crucially, all hyperparameters are also tuned fully offline. Secondly, we introduce TOReL: a **tuning for offline reinforcement learning** algorithm that extends our information rate based offline hyperparameter tuning methods to general offline RL approaches. Our empirical evaluation confirms SOReL’s ability to accurately estimate regret in the Bayesian setting whilst TOReL’s offline hyperparameter tuning achieves competitive performance with the *best online hyperparameter tuning* methods *using only offline data*. Thus, SOReL and TOReL make a significant step towards safe and reliable offline RL, unlocking the potential for RL in the real world. Our implementations are publicly available: ANONYMISED.

1 INTRODUCTION

Offline RL (Lange et al., 2012; Levine et al., 2020; Murphy, 2024) promises to unlock the potential for agents to act autonomously, successfully and safely from the moment they are deployed into an environment. However, existing offline RL methods (Tarasov et al., 2023; Kostrikov et al., 2021; Kidambi et al., 2020; Yu et al., 2021) are yet to fulfil this promise due to two key issues that have been largely ignored by the community. **Issue I: there are currently no offline metrics to tune hyperparameters or choose between approaches.** Existing methods rely on *online* samples to carry out the extensive hyperparameter tuning required to achieve high performance (Zhang et al., 2021; Jackson et al., 2025). As we sketch in Fig. 1a, this results in cycles of training offline, deployment, failure online, further hyperparameter tuning and/or model selection, re-training offline and redeployment until the online performance of the agent is acceptable. Without a method for offline tuning, existing methods suffer from *high online sample complexity*, which is precisely the problem offline RL intends to solve.

For many problem settings, we also need reliable online performance guarantees *before* the agent is deployed online. Precisely, **issue II: current methods offer no reliable offline method to approximate true online regret**, i.e. the difference between expected returns of an optimal policy and a policy trained using offline data. This is concerning from an AI safety perspective, as without a reliable regret bound, we cannot deploy agents into the real world where agent failure presents a serious hazard to human life; for these settings it is essential that agents are deployed with near-zero regret. For less sensitive domains, users still need some guarantee of optimality on deployment; there will often be a clear cost associated with deviation from optimal policy in terms of regret, for example in settings where the degree to which a product can fail will result in financial loss that can be determined before deployment using regret and kept in a tolerable range.

To tackle these two essential issues, we develop a Bayesian framework where the posterior (conditioned on the offline data) is used as a *prior* for a Bayesian RL problem (Martin, 1967; Duff, 2002).

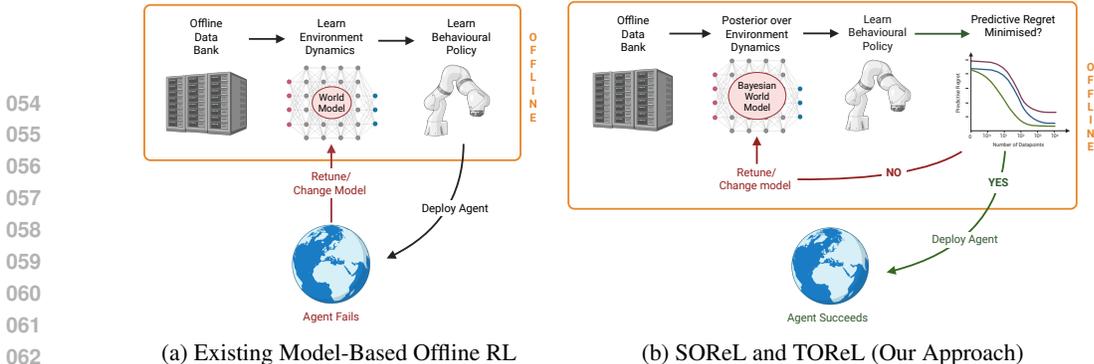


Figure 1: Existing model-based offline approaches rely on online interactions for hyperparameter tuning and verifying accurate model learning before they can achieve good performance, leading to poor online sample efficiency. In contrast, in SOReL/TOReL, model tuning and world model learning is carried out fully offline using the predictive regret as a tuning signal. Only then is the agent trained and deployed.

Our analysis reveals that the regret of the corresponding Bayes-optimal policy is controlled by the posterior information loss (PIL) - that is the expected posterior KL divergence between the model and true dynamic. The change in PIL, known as the *information rate*, measures how much information the model has gained from an incremental amount of offline data. Crucially, the PIL can be estimated and tracked during offline training, allowing us to monitor performance and tune hyperparameters completely offline.

For our first method, we develop SOReL, a theoretically grounded framework for model-based safe offline reinforcement learning which resolves both key issues. By using offline data to infer a posterior over environment dynamics, we approximate regret using the predictive variance and median of policy rollouts *prior to deployment*, directly tackling **issue II**. Moreover, both the predictive regret and the PIL can be used as a signal to tune hyperparameters offline, thereby tackling **issue I** (see Fig. 1b). Only then is the trained agent deployed safely, making SOReL (to the authors’ knowledge) the first fully offline RL approach with reliable performance guarantees once deployed. Our experiments support this claim empirically, showing that in the standard offline RL MuJoCo control tasks (Yu et al., 2020; Kidambi et al., 2020; Ball et al., 2021; Lu et al., 2022b; Sun et al., 2023; Sims et al., 2024), SOReL’s *offline* regret approximation accurately tracks the true regret once deployed *online*.

As SOReL is a general Bayesian approach, existing methods may outperform it in specific domains as many have been designed to exploit heuristics tailored to specific datasets and tasks. To address this, we also extend SOReL’s offline hyperparameter tuning methods to address **issue I** for existing model-free and model-based offline RL approaches (Yu et al., 2020; Kidambi et al., 2020; Ball et al., 2021; Lu et al., 2022b; Sun et al., 2023; Tarasov et al., 2023; Sims et al., 2024) when accurate regret estimation isn’t required, greatly improving their sample efficiency. Using this insight, we develop TOReL, an algorithm for tuning offline reinforcement learning that uses the PIL signal and tracks a *regret metric correlated to the true regret* for tuning. To test our method, we apply TOReL to IQL (Kostrikov et al., 2021), ReBRAC (Tarasov et al., 2023), MOPO (Yu et al., 2020) and MOReL (Kidambi et al., 2020) to carry out hyperparameter tuning in Adroit gymnasium and the standard offline RL MuJoCo control tasks. Using only *offline* data, TOReL achieves similar performance to existing methods that carry out *online* hyperparameter tuning. Notably, when combined with ReBRAC, TOReL consistently finds a hyperparameter combination with near-zero regret, outperforming all hyperparameters for all other algorithms. When comparing TOReL’s offline hyperparameter tuning to a recent online UCB approach (Jackson et al., 2025), we see that UCB typically requires about a dataset’s worth of *online* samples to match TOReL’s performance. We summarise our key contributions:

- I In Section 4, we develop a Bayesian framework for model-based offline RL;
- II In Section 4.3 we carry out a regret analysis for our framework, demonstrating regret is controlled by the PIL and provide a strong frequentist justification for our Bayesian approach;
- III In Section 5.1 we develop SOReL, a method for approximating true regret using predictive uncertainty, which can achieve a desired and safe level of true regret once deployed;
- IV In Section 5.2 we introduce TOReL, adapting SOReL’s offline hyperparameter/model tuning approach to general offline model-based and model-free RL;
- V In Section 6 we empirically confirm that TOReL addresses **issue I** for existing methods and SOReL’s ability to address **issues I and II** for more conservative applications.

2 PRELIMINARIES

Let X be a $\mathcal{X} \subseteq \mathbb{R}^n$ -valued random variable. We denote a distribution as P_X with density (if it exists) as $p(x)$. We denote the set of all distributions over \mathcal{X} as $\mathcal{P}(\mathcal{X})$. We introduce the notation $\mathcal{G}(p)$ to represent the geometric distribution and $\mathcal{AG}(p)$ to represent the arithmetico-geometric distribution, with probability mass functions: $P_{\mathcal{G}}(i) := (1-p)p^i$ and $P_{\mathcal{AG}}(i) := (1-p)^2 p^i (i+1)$ respectively, for $i \in \mathbb{N}_0^+$ and parameter $p \in [0, 1)$. We denote the uniform distribution over $\{0, 1, \dots, i\}$ as \mathcal{U}_i and the multivariate normal distribution with mean vector μ and covariance matrix Σ as $\mathcal{N}(\mu, \Sigma)$.

2.1 OFFLINE REINFORCEMENT LEARNING

For our offline RL setting, an agent is tasked with solving the learning problem in an infinite-horizon, discounted Markov decision process (Bellman, 1956; 1958; Sutton and Barto, 2018; Puterman, 1994; Szepesvári, 2010): $\mathcal{M}^* := \langle \mathcal{S}, \mathcal{A}, P_0, P_S^*(s, a), P_R^*(s, a), \gamma \rangle$, with state space \mathcal{S} , action space \mathcal{A} and discount factor γ . At time $t = 0$, an agent starts in an initial state allocated according to the initial state distribution: $s_0 \sim P_0$. At every timestep t , an agent in state s_t takes an action according to a policy $a_t \sim \pi(h_t)$ ¹, receives a scalar reward $r_t \sim P_R^*(s_t, a_t)$ and transitions to a new state $s_{t+1} \sim P_S^*(s_t, a_t)$ where $h_t := \{s_0, a_0, r_0, s_1, a_1, r_1, \dots, a_{t-1}, r_{t-1}, s_t\} \in \mathcal{H}_t$ is the observed a history of interactions with the environment. Here $\mathcal{H}_t := \mathcal{S} \times (\mathcal{A} \times \mathbb{R} \times \mathcal{S})^{\times t}$ denotes the corresponding product space. We assume rewards are bounded with $r_t \in [r_{\min}, r_{\max}] \subset \mathbb{R}$ where r_{\min} and r_{\max} denote the minimum and maximum reward values respectively. For convenience, we often write the joint state transition-reward distribution as $P_{R,S}^*(s, a)$. We denote the distribution over history h_t as $P_{t,\pi}^*$. The goal of an agent is to learn an optimal policy $\pi^* \in \Pi^*$ where $\Pi^* := \arg \max_{\pi} J^{\pi}(\mathcal{M}^*)$ is the set of policies that maximise the expected discounted return $J^{\pi}(\mathcal{M}^*) := \mathbb{E}_{h_{\infty} \sim P_{\infty,\pi}^*} [\sum_{i=0}^{\infty} \gamma^i r_i]$. It suffices to consider only optimal policies that condition only on the most recent state (i.e. $\pi^*(s_t)$) as, in a fully observable MDP, any optimal history-conditioned policy will never take an action that cannot be taken by an optimal policy that conditions only on most recent state.

In the learning setting, the true state transition distribution $P_S^*(s, a)$ and reward distribution $P_R^*(s, a)$ are assumed unknown a priori. Once deployed, the agent is faced with the exploration/exploitation dilemma in that it must balance exploring to learn about the unknown environment dynamics with exploiting. In offline RL (Lange et al., 2012; Levine et al., 2020; Murphy, 2024), an agent has access to a dataset of histories of various lengths collected from the true environment. The policies used to collect the data may vary and not be optimal. In the zero-shot model-based offline RL setting (Jackson et al., 2025), the dataset is used to learn the unknown environment dynamics from which a policy is trained prior to any interaction with the environment. The agent is then deployed at test time $t = 0$ and its performance evaluated. The goal of offline RL is to take advantage of offline data so that the deployed policy will be near-optimal from the outset.

3 RELATED WORK

Developing reliable model-based offline RL algorithms remains an open challenge for several reasons; most existing methods are not fully offline, requiring extensive online interactions for tuning. Only limited attention has been given to solving key **issues I and II** introduced in Section 1, which is the focus of this work: Paine et al. (2020) introduce a method for estimating online value and partial hyperparameter tuning of offline model-free algorithms, however their method neither approximates regret, nor is an accurate proxy for true value, resulting in significant overestimation in most domains. As noted by Smith et al. (2023), their approach relies on offline policy evaluation, which is a challenging and provably difficult problem (Wang et al., 2020) whose hyperparameters require tuning online. Moreover, as noted by Jackson et al. (2025), their framework is limited to behavioural cloning and two model-free critic-based methods that have since been outperformed by modern algorithms. Smith et al. (2023) introduce a method for offline hyperparameter tuning, but are limited to the model-free imitation learning setting and offer no regret estimation. Finally, Wang et al. (2022) introduce a method for offline hyperparameter tuning to pre-select hyperparameters for online methods, but do not learn optimal policies offline or provide regret approximation. In contrast to all of these approaches, to the authors’ knowledge, our method is the first offline RL method to reliably approximate regret and carry out *all hyperparameter* tuning for general methods using *only offline data*. Finally, understanding of offline RL from a Bayesian perspective is limited. To the authors’ knowledge, only Chen et al. (2024) have framed solving offline model-based RL as solving a

¹Policies condition on history as we work within a Bayesian paradigm, using methods such as RNN-PPO (Schulman et al., 2017)

BAMDP, however no regret analysis of the Bayes-optimal policy is carried out, a continuous BAMCP (Guez et al., 2014) approximation is used to learn behavioural policies and the algorithm still suffers from a lack of regret approximation, hence relying on online data for tuning.

In addition, the performance of offline RL approaches is particularly dependent on the ability to accurately model transition dynamics as errors in a dynamics model can compound over several timesteps for the long-horizon problems encountered in RL (see also our analysis in Section 4.3) and many datasets used to benchmark methods contain missing datapoints in critical regions of state-action space, which poses a generalisation challenge. We note that both these problems are *orthogonal* to solving key **issues I and II**, which we focus on in this work. Most existing offline RL methods focus on tackling the missing data problem by introducing a form of reward pessimism based on the model uncertainty (Yu et al., 2020; Kidambi et al., 2020; Kumar et al., 2020; Yu et al., 2021; Fujimoto and Gu, 2021; Kostrikov et al., 2021; An et al., 2021; Ball et al., 2021; Lu et al., 2022b; Sun et al., 2023; Tarasov et al., 2023; Sims et al., 2024). Unifloral (Jackson et al., 2025) is a recent framework that unites these offline RL approaches into a single algorithmic space with lightweight and high-performing implementations, as well as providing a clarifying benchmarking protocol. Our implementations and evaluation methods follow this framework.

4 BAYESIAN OFFLINE RL

We now introduce our Bayesian RL framework, which constitutes learning an (approximate) posterior from offline data before solving a Bayesian RL problem with the posterior acting as the prior. We provide an introductory primer on Bayesian RL in Section B and all proofs for all theorems can be found in Section D.

4.1 LEARNING A POSTERIOR WITH OFFLINE DATA

A Bayesian epistemology characterises the agent’s uncertainty in the MDP through distributions over any unknown variable (Martin, 1967; Duff, 2002). We first specify a parametric model $p(r_t, s_{t+1}|s_t, a_t, \theta)$, $P_{R,S}(s_t, a_t, \theta)$, $\theta \in \Theta \subset \mathbb{R}^d$, over the unknown state transitions and reward distributions, with each $\theta \in \Theta \subseteq \mathbb{R}^d$ representing a hypothesis about the MDP \mathcal{M}^* . As we show in Section 4.3, our results can easily be generalised to non-parametric methods like Gaussian process regression (Rasmussen and Williams, 2006; Wiener, 1923; Krige, 1951). A prior distribution over the parameter space P_Θ is specified, which represents the initial *a priori* belief in the true value of $P_{R,S}^*(s, a)$ before the agent has observed any transitions.

We denote an offline dataset of N state-action-state-reward transition observations as: $\mathcal{D}_N = \{(s_i, a_i, s'_i, r_i)\}_{i=0}^{N-1}$, all collected from a single MDP \mathcal{M}^* . Datapoints may be collected from several policies and non-Markovian sampling. Given the dataset \mathcal{D}_N , the prior P_Θ with density $p(\theta)$ is updated to posterior $P_\Theta(\mathcal{D}_N)$ with density $p(\theta|\mathcal{D}_N)$, using Bayes’ rule:

$$p(\theta|\mathcal{D}_N) = \frac{p(\mathcal{D}_N|\theta)p(\theta)}{p(\mathcal{D}_N)} = \frac{\prod_{i=0}^{N-1} p(r_i, s'_i|s_i, a_i, \theta)p(\theta)}{\int_{\Theta} \prod_{i=0}^{N-1} p(r_i, s'_i|s_i, a_i, \theta)p(\theta)d\theta}. \quad (1)$$

The posterior represents the agent’s belief in the unknown environment dynamics once \mathcal{D}_N has been observed. We now detail how a Bayes-optimal policy is learned using the posterior as the initial belief in the environment dynamics.

4.2 LEARNING A BAYES-OPTIMAL POLICY

It is well known that solving a Bayesian RL problem exactly is intractable for all but the simplest models (Martin, 1967; Duff, 2002; Guez et al., 2012; 2013; Zintgraf et al., 2020; Fellows et al., 2024). Inferring the posterior in Eq. (1) is typically infeasible for dynamics models of interest (for example, nonlinear Gaussian world models). This is because there is no analytic solution for the posterior density and the cost of carrying out integration required to evaluate the evidence $p(\mathcal{D}_N)$ grows exponentially in parameter dimensions d . Fortunately, there exist tractable methods to learn an approximate posterior $\hat{P}_\Theta(\mathcal{D}_N) \approx P_\Theta(\mathcal{D}_N)$; in this paper we use randomised priors (Osband and Van Roy, 2017; Osband et al., 2018; Ciosek et al., 2020) (RP) and provide details in Section F.2. In addition, a planning problem must be solved for every conceivable history that an agent could encounter. In our offline RL setting, we ease intractability by replacing the prior P_Θ with a highly informative posterior $P_\Theta(\mathcal{D}_N)$, significantly reducing the hypothesis space in the Bayesian RL problem.

Let $P_{\infty, \pi}(\theta)$ denote the corresponding model distribution over h_∞ for policy $\pi(h_t)$. To obtain an (approximately) Bayes-optimal policy, we use the Bayesian RL objective in the meta-

learning form (Zintgraf et al., 2020; Beck et al., 2024) (i.e. as an expectation using $\hat{P}_\Theta(\mathcal{D}_N)$) so that a simple RL²(Duan et al., 1987) style algorithm can be applied: $J_{\text{Bayes}}^\pi(\hat{P}_\Theta(\mathcal{D}_N)) := \mathbb{E}_{\theta \sim \hat{P}_\Theta(\mathcal{D}_N)} [\mathbb{E}_{h_\infty \sim P_{\infty, \pi}(\theta)} [\sum_{i=0}^{\infty} \gamma^i r_i]]$. Solving the Bayesian RL objective is known as solving a Bayes-adaptive MDP (BAMDP) (Duff, 2002). We optimise the objective $J_{\text{Bayes}}^\pi(\hat{P}_\Theta(\mathcal{D}_N))$ by sampling a hypothesis environment from the approximate posterior $\theta \sim \hat{P}_\Theta(\mathcal{D}_N)$ then rolling out the policy in the sampled environment dynamics. The Bayes-optimal policy $\pi_{\text{Bayes}}^* \in \arg \max_{\pi} J_{\text{Bayes}}^\pi(\hat{P}_\Theta(\mathcal{D}_N))$ is learned using RNN-PPO (Schulman et al., 2017) as a BAMDP solver on the rollouts. Complete implementation details can be found in Section F.

In addition to having excellent exploration/exploitation properties, a Bayesian approach affords access to epistemic uncertainty in the returns via the variance of predictive rollouts. Uncertainty estimation is essential for tackling our two keys **issues I and II**; firstly, as we show in Section 6.2, the predictive variance and predictive median of policy returns can be used to estimate the true regret at test time. Secondly, monitoring the decay of predictive variance and regret is a powerful tool for diagnosing issues with the BAMPD planner offline. Finally, we remark that a Bayesian approach is relatively simple compared to existing model-based approaches in Section 3 as it does not rely on hand-crafted heuristics tailored to specific problem settings and *scales as well as existing state-of-the-art offline RL methods* (Jackson et al., 2025) which also rely on ensembling methods for uncertainty quantification.

4.3 FREQUENTIST JUSTIFICATION OF BAYESIAN OFFLINE RL

We now carry out a frequentist asymptotic regret analysis for Bayesian offline RL. For finite N , we can measure how far the performance of the Bayes optimal policy is from an optimal policy using the *true regret*, which is the difference between the expected return $J^{\pi_{\text{Bayes}}^*}(\mathcal{M}^*, \mathcal{D}_N)$ of the Bayes-optimal policy π_{Bayes}^* given a posterior $P_\Theta(\mathcal{D}_N)$, all in the true MDP \mathcal{M}^* : $\text{Regret}(\mathcal{M}^*, \mathcal{D}_N) := J^{\pi^*}(\mathcal{M}^*) - J^{\pi_{\text{Bayes}}^*}(\mathcal{M}^*, \mathcal{D}_N)$. The goal of this analysis twofold; firstly, we show that the rate of regret decreases for a Bayes-optimal policy is characterised by an easy to estimate quantity known as the *posterior information loss* (PIL). Secondly, our goal is show (for the first time) that the asymptotic regret of offline Bayesian RL converges at the $\mathcal{O}(1/\sqrt{N})$ rate found in prior work (Yu et al., 2020; Kidambi et al., 2020), which is known to be the optimal asymptotic convergence rate for nonlinear models in general MDPs for offline RL (Agarwal et al., 2022). This provides a strong frequentist justification for our Bayesian offline RL framework. Our key result shows that regret can be bounded using the PIL, defined as: $\mathcal{I}_N^\pi := \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [\mathbb{E}_{s, a \sim \rho_\pi^*} [\text{KL}(P_{R,S}^*(s, a) \| P_{R,S}(s, a, \theta))]]$. Here $\rho_\pi^* := \mathbb{E}_{i \sim \mathcal{AG}(\gamma)} [\mathbb{E}_{j \sim \mathcal{U}_i} [P_{j, \pi}^*]]$ is the arithmetico-geometric ergodic state-action distribution, which places mass over state-action pairs according to how much errors in the model influence the regret at each state. Regions of state-action space that require more timesteps to reach from initial states are weighted significantly less than those that are encountered earlier and more frequently, as state errors encountered early accumulate in each prediction from that timestep onwards. The PIL has an intuitive information-geometric interpretation which we discuss further in Section D.2.

Theorem 1. Let $\mathcal{R}_{\max} := \frac{(r_{\max} - r_{\min})}{1 - \gamma}$ denote the maximum possible regret for the MDP. Using the

PIL: \mathcal{I}_N^π , the true regret is bounded as: $\text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \leq 2\mathcal{R}_{\max} \cdot \sup_{\pi} \sqrt{1 - \exp\left(-\frac{\mathcal{I}_N^\pi}{1 - \gamma}\right)}$

We observe from Theorem 1 that the rate at which regret decreases with N is governed by the rate the PIL decreases, which is known as the *information rate*. This measures how much information the model has gained from an incremental amount of data. Fast information rates imply highly informative posteriors can be learned using minimal data as regret will decrease at least as fast. How fast the information rate is depends on the exact model specification, prior and underlying MDP. Formulating our bound in terms of the PIL ties the regret to the KL divergence over the reward-state model: $\text{KL}(P_{R,S}^*(s, a) \| P_{R,S}(s, a, \mathcal{D}_N))$. Not only is this mathematically more convenient, yielding a simpler bound, but the PIL is *easy to estimate in practice* meaning the information rate can be monitored offline as a proxy for downstream regret and to carry out hyperparameter tuning associated with the model and approximate inference method, partially resolving key **issue I**. Using Theorem 1, we can study the PIL \mathcal{I}_N^π for different classes of models which allows us to understand how regret will evolve given the model choice. This also provides a frequentist justification for many Bayesian approaches. We now characterise the information rate for parametric models, which allows us to recover the optimal $\mathcal{O}(1/\sqrt{N})$ regret convergence rate for Bayesian offline RL:

Theorem 2. Let the data be drawn from the underlying true distribution $\mathcal{D}_N \sim P_{Data}^*$. Under standard local asymptotic normality assumptions (see Assumption 1 in Section D.3), there exists some constant $0 < C < \infty$ such that for sufficiently large N : $\mathbb{E}_{\mathcal{D}_N \sim P_{Data}^*} [\text{Regret}(\mathcal{M}^*, \mathcal{D}_N)] \leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-\frac{Cd}{(1-\gamma)N}\right)} = \mathcal{O}\left(\frac{1}{\sqrt{(1-\gamma)N}}\right)$.

4.4 GAUSSIAN WORLD MODELS

Many methods specify Gaussian reward and state transition models of the form:

$$P_R(s, a, \theta) = \mathcal{N}(r_\theta(s, a), \sigma_r^2(s, a)), \quad P_S(s, a, \theta) = \mathcal{N}(s'_\theta(s, a), I\sigma_s^2(s, a)), \quad (2)$$

with isotropic variance characterised by σ_r^2 and σ_s^2 , mean reward function $r_\theta(s, a)$ and mean state transition function $s'_\theta(s, a)$. Using a Gaussian world model, we find the PIL takes a convenient and intuitive form. Let $r(s, a, \mathcal{D}_N) := \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [r_\theta(s, a)]$ and $s'(s, a, \mathcal{D}_N) := \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [s_\theta(s, a)]$ denote the Bayesian mean reward and state transition functions and $r^*(s, a)$ and $s^*(s, a)$ denote the true mean functions. We define the mean squared error between the true and Bayesian mean functions as:

$$\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) := \mathbb{E}_{(s,a) \sim \rho_\pi^*} \left[\frac{\|r(s, a, \mathcal{D}_N) - r^*(s, a)\|_2^2}{2\sigma_r^2(s, a)} + \frac{\|s'(s, a, \mathcal{D}_N) - s^*(s, a)\|_2^2}{2\sigma_s^2(s, a)} \right], \quad (3)$$

and the predictive variance as:

$$\mathcal{V}(\mathcal{D}_N) := \mathbb{E}_{(s,a) \sim \rho_\pi^*} \left[\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\frac{\|r(s, a, \mathcal{D}_N) - r_\theta(s, a)\|_2^2}{2\sigma_r^2(s, a)} + \frac{\|s'(s, a, \mathcal{D}_N) - s'_\theta(s, a)\|_2^2}{2\sigma_s^2(s, a)} \right] \right]. \quad (4)$$

We now re-write the PIL for the Gaussian world model using these two terms:

Proposition 1. Using the Gaussian world model in Eq. (2), it follows: $\mathcal{I}_N^\pi = \mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) + \mathcal{V}(\mathcal{D}_N)$.

Proposition 1 shows that the PIL is governed by i) the MSE of the point estimate $\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*)$, which characterises how quickly the Bayesian mean function converges to the true function; and ii) the predictive variance $\mathcal{V}(\mathcal{D}_N)$, which characterises the epistemic uncertainty in the model. For frequentist methods using point estimates like the MLE, there is no characterisation of epistemic uncertainty, meaning $\mathcal{V}(\mathcal{D}_N) = 0$. The PIL can easily be estimated by estimating $\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*)$ using the empirical MSE with offline data and estimating $\mathcal{V}(\mathcal{D}_N)$ using posterior sampling.

5 FROM THEORY TO PRACTICE

Our frequentist analysis in Section 4.3 provides valuable intuition about how we might expect regret to change depending on the choice of model, however it cannot address key **issues I and II** from Section 1. This is because asymptotic results are only theoretical guarantees that apply with high probability in the limit of large data across a space of MDPs; their main use is a theoretical object to provide a way to compare asymptotic behaviour of algorithms rather than for obtaining reliable regret estimation or a being a metric to tune algorithms as they cannot actually be calculated. As a concrete example, the regret bound of MOREL (Kidambi et al., 2020) is: $\text{Regret}_N \leq C/\sqrt{N} + m$ which is impossible to calculate in the offline RL setting because 1) C is not specified, it is just proved to exist; 2) the constant m requires knowing: the hitting time of all unknown states in the original MDP of the optimal policy, the smallest nonzero elements of the initial state distribution and the smallest elements of state-transition distribution; and 3) the bound only holds for large enough N , which is unspecified. Even if these quantities were known, we would only have a bound and not an estimate of the true regret. For these reasons, tackling **issues I and II** is a challenging problem to solve.

5.1 SOREL

We now introduce SOREL in Algorithm 1, our algorithm for reliable regret estimation and offline hyperparameter tuning. In our SOREL framework, there are three sets of hyperparameters: ϕ_I the **model** (such as the architecture for a neural-network function approximator); ϕ_{II} the **approximate inference method** (such as the number of ensemble members for RP); and ϕ_{III} the **BAMDP solver** (the hyper-parameters of a Bayesian meta-learning algorithm like RNN-PPO). Sets ϕ_I and ϕ_{II} are tuned jointly to both minimise the PIL and ensure a roughly even split between the predictive variance and MSE loss terms. Set ϕ_{III} is then tuned to minimise approximate regret based on the now-fixed model and approximate posterior: for each combination of hyper-parameters, we learn a policy using the BAMDP solver, and choose the combination whose policy leads to the lowest approximate regret. $\mathcal{R}_{\text{Deploy}}$ denotes the desired level of regret of the deployed policy.

Fixing Issue I - PIL Monitoring: To tune sets ϕ_I and ϕ_{II} , we monitor the change in PIL: \mathcal{I}_N^π . Our goal is to select hyperparameters that minimises the PIL whilst ensuring the the MSE term (c.f. Eq. (3)) closely matches the predictive variance term (c.f. Eq. (4)): $\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) \approx \mathcal{V}(\mathcal{D}_N)$. Misalignment of predictive variance and MSE indicates either an overfitting/underfitting issue with model hyperparameters in set ϕ_I and/or an issue with uncertainty estimation due to approximate inference hyperparameters in set ϕ_{II} . Moreover, under/overestimating uncertainty will lead to poor regret estimation, which is why we tune sets ϕ_I and ϕ_{II} first in Algorithm 1. Our empirical evaluations in Section 6.2 confirm that when $\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) \approx \mathcal{V}(\mathcal{D}_N)$, the approximate regret aligns strongly with true regret.

Fixing Issues I and II - Regret Approximation: A simple approach for bounding regret would be to estimate the PIL from the offline data and then apply our theoretical upper bound in Theorem 1. This is likely to be too conservative for most applications as it protects against the worst case MDP that the agent could encounter. In particular, it is very sensitive to errors in the model, especially as $\gamma \rightarrow 1$, which is an artifact of model errors accumulating over all future timesteps in the regret analysis. Instead, we approximate the regret using the posterior predictive median:

$$\text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \approx \hat{R}_{\max} - \hat{\mathbb{M}}_{\theta \sim P_{\Theta}(\mathcal{D}_N), h_{\infty} \sim P_{\infty}^{\pi}(\theta)} [R(h_{\infty})], \quad (5)$$

where $\hat{\mathbb{M}}_{\theta \sim P_{\Theta}(\mathcal{D}_N), h_{\infty} \sim P_{\infty}^{\pi}(\theta)} [R(h_{\infty})]$ denotes the median predictive return based on sampling from the (approximate) posterior and rolling out the Bayes-optimal policy and \hat{R}_{\max} is estimated from the maximum return in the offline dataset - full details and an overview of alternative metrics that can be derived from the posterior to approximate the regret with varying degrees of conservatism are found in Section C.2. We hypothesise that the sample median offers a good compromise: neither overly conservative nor overly susceptible to being skewed by a policy that performs well on only a subset of posterior samples. Our empirical evaluations support this hypothesis in Section 6.2. The posterior predictive median also allows us to tune hyperparameter set ϕ_{III} , selecting hyperparameters to learn a policy that achieves the lowest approximate regret as shown in Algorithm 1.

5.2 TOReL

SOReL’s offline hyperparameter tuning methods are directly applicable to general offline RL approaches, allowing us to address **issue I** for existing offline methods. We now adapt these methods to derive a general tuning for offline reinforcement learning approach called TOReL, shown in Algorithm 2. A policy is learned offline using a planning algorithm, denoted by ORL. There thus exists a corresponding set of hyperparameters associated with ϕ_{III} the **offline planner**. For model-based methods with uncertainty estimation like MOREL Kidambi et al. (2020) and MOPO Yu et al. (2020), we can exactly adapt SOReL’s PIL tuning method to the parameters associated with: ϕ_I the **dynamics model** and ϕ_{II} **uncertainty estimation**. For all other methods, we introduce and learn a dynamics model and an approximate inference method like in SOReL and jointly tune the corresponding hyperparameters ϕ_I and ϕ_{II} to minimise the PIL without requiring the even split between the predictive variance and MSE loss terms. Since the policy learned with ORL is typically neither Bayes-Optimal nor robust to model uncertainty, we expect that applying SOReL’s regret approximation method to more general methods in TOReL will not yield an accurate estimate of the regret in terms of its absolute value. Instead, we treat the approximate regret in Eq. (5) as a *regret metric* that is positively correlated with true regret, and use this to tune ORL parameters ϕ_{III} . Our empirical evaluations in Section 6.1 support this hypothesis. We note that in model-free methods, the dynamics model and an approximate inference method are not used in policy learning, only to aid regret metric calculation.

Algorithm 1 SOReL($P_{\Theta}, \mathcal{D}_N, \mathcal{R}_{\text{Deploy}}$)

```

 $R_N \leftarrow \hat{R}_{\max}$ 
while  $R_N > \mathcal{R}_{\text{Deploy}}$  do
  Hyperparameter tuning:
   $\phi_I, \phi_{II} \leftarrow \arg \min_{\phi_I, \phi_{II}} \text{PIL}(\phi_I, \phi_{II}, \mathcal{D}_N)$ 
  s.t.  $\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) \approx \mathcal{V}(\mathcal{D}_N)$ 
   $\phi_{III} \leftarrow \arg \min_{\phi_{III}} \text{ApproxRegret}(\phi_I, \phi_{II}, \phi_{III}, \mathcal{D}_N)$ 
  Policy Learning and Regret Approximation:
   $\pi_{\text{Bayes}}^* \leftarrow \text{SolveBAMDP}(\phi_I, \phi_{II}, \phi_{III}, \mathcal{D}_N)$ 
   $R_N \leftarrow \text{ApproxRegret}(\phi_I, \phi_{II}, \phi_{III}, \mathcal{D}_N)$ 
end while
return  $\pi_{\text{Bayes}}^*$ 

```

Algorithm 2 TOReL($P_{\Theta}, \mathcal{D}_N$)

```

 $\phi_I, \phi_{II} \leftarrow \arg \min_{\phi_I, \phi_{II}} \text{PIL}(\phi_I, \phi_{II}, \mathcal{D}_N)$ 
  [ s.t.  $\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) \approx \mathcal{V}(\mathcal{D}_N)$ , (model-based) ]
 $\phi_{III} \leftarrow \arg \min_{\phi_{III}} \text{RegretMetric}(\phi_I, \phi_{II}, \phi_{III}, \mathcal{D}_N)$ 
 $\pi_{\text{TOReL}}^* \leftarrow \begin{cases} \text{ORL}(\phi_{III}, \mathcal{D}_N), & \text{(model-free)} \\ \text{ORL}(\phi_I, \phi_{II}, \phi_{III}, \mathcal{D}_N), & \text{(model-based)} \end{cases}$ 
return  $\pi_{\text{TOReL}}^*$ 

```

6 EXPERIMENTS

The goal of experimental section is to confirm that our methods successfully solves key **issues I and II**. In Section 6.1, we confirm that TOReL resolves **issue I** for existing ORL algorithms, consistently identifying hyperparameters with a lower regret than the average regret of randomly chosen hyperparameters using *offline data alone*. In Section 6.2, we confirm that SOReL is the first approach that can resolve **issue II**, accurately approximating regret over a range of regret curves.

6.1 TOReL IS AN EFFECTIVE OFFLINE HYPERPARAMETER TUNER FOR ORL

For this experiment, our goal is to use TOReL to identify hyperparameters for ReBRAC Tarasov et al. (2023) and IQL Kostrikov et al. (2021) (two model-free algorithms), and MOPO Yu et al. (2020) and MOREL Kidambi et al. (2020) (two model-based algorithms). We use Unifloral’s online hyperparameter tuning framework as it matches or achieves better performance than papers originally report through its sophisticated online UCB hyperparameter search methods Jackson et al. (2025, Section C). Details are given in Section F.4. In Fig. 2 and Fig. 6 in Section E, we compare the regret of the TOReL-selected hyperparameter combination to the true regret, which we define as the expected regret over all possible hyperparameter combinations. We also compare against the oracle regret: the minimum regret achieved by any hyperparameter combination. We evaluate each algorithm on 8 offline datasets: 200K randomly sampled transitions from each of our three brax datasets, and in the D4RL Fu et al. (2020) Adroit (pen-expert and hammer-expert) and locomotion datasets (halfcheetah-medium-expert, hopper-medium and walker2d-medium-replay) suggested by Jackson et al. (2025). Full results are in Table 3 of Section E. Table 1 shows ReBRAC+TOReL as a consistently high-achieving combination, reaching near-oracle performance on every dataset. For two-thirds of the tasks and algorithms there is statistically significant ($p < 0.05$), strong ($r > |0.5|$) positive ($r > 0$) Pearson correlation between the ensemble median regret metric and the true regret (Table 5 and Fig. 9 in Section E). Where no strong positive correlation is observed (possibly due to limited hyperparameter coverage) the average TOReL regret (0.433) is still lower than the corresponding true regret (0.458).

Our final experiment analyses the number of samples saved using TOReL rather than the UCB bandit-based online hyperparameter selection algorithm proposed by Jackson et al. (2025). We tune hyperparameters for ReBRAC, as ReBRAC achieves the lowest regret across all tasks and algorithms. Results for the D4RL and brax datasets are depicted in Fig. 3. TOReL offers significant savings in terms of sample complexity compared to existing online hyperparameter tuning methods: for the D4RL datasets, 20K to >200K additional online samples are spared, while for the brax datasets >200K are spared, essentially preventing a doubling of the size of the offline dataset.

6.2 SOReL IS A SAFE ALGORITHM FOR ORL

We demonstrate how SOReL can be implemented as a *safe* ORL algorithm. Our goal is test whether SOReL accurately approximates regret over a range of datasets and collection policies including transitions from poor, medium and expert regions of performance to produce high, medium and low regret curves. We evaluate in 5 environments: two gymnax environments and three brax environments. Referring back to Fig. 1b, we progressively include more offline data to learn a policy until a safe level of approximate regret is achieved. For our implementation, we use a variation of the standard Gaussian world model presented in Section 4.4, randomised priors (Osband et al., 2018) for

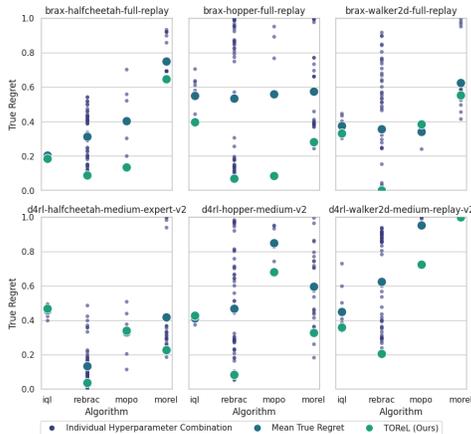


Figure 2: TOReL-selected hyperparameter regret versus mean hyperparameter regret.

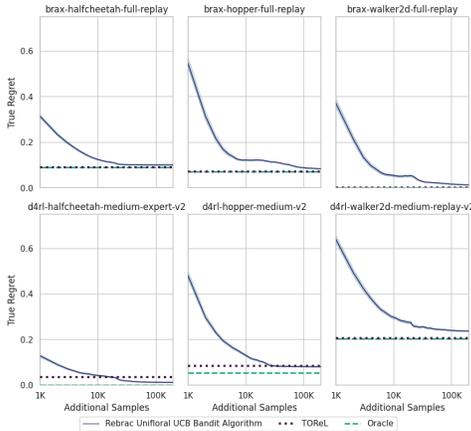


Figure 3: TOReL compared to UCB bandit-based online hyperparameter selection. The x-axis shows the additional samples required during online tuning. Note that the size of the D4RL datasets are 1998K, 1000K and 302K (left to right). UCB 95th percentile confidence interval shaded.

Task	Algo.	Oracle	TOReL	Oracle Mean	TOReL Mean	True
brax-halfcheetah-full-replay	ReBRAC	0.089	0.089	0.262	0.264	0.417
brax-hopper-full-replay	ReBRAC	0.070	0.070	0.193	0.209	0.554
brax-walker-full-replay	ReBRAC	0.000	0.000	0.241	0.317	0.425
d4rl-halfcheetah-medium-expert-v2	ReBRAC	0.000	0.036	0.176	0.268	0.336
d4rl-hopper-medium-v2	ReBRAC	0.053	0.083	0.380	0.323	0.580
d4rl-walker2d-medium-replay-v2	ReBRAC	0.204	0.206	0.567	0.572	0.757
d4rl-pen-expert-v1	ReBRAC	0.033	0.188	0.510	0.564	0.570
d4rl-hammer-expert-v1	ReBRAC	0.086	0.159	0.585	0.604	0.684

Table 1: TOReL Regret Summary Statistics (lower is better): bold indicates where TOReL is within 5% of the corresponding Oracle. Left: algorithm chosen if the Oracle can choose over both hyperparameters *and* algorithms; corresponding oracle regret and regret of the TOReL-chosen hyperparameters for that algorithm. Middle: oracle and TOReL regrets averaged over all algorithms. Right: true regret averaged over all algorithms. For these tasks, ReBRAC+TOReL is consistently the best.

approximate inference and RNN-PPO (Schulman et al., 2017) to solve the BAMDP. Implementation and dataset details are found in Section F.

In practice, each time additional offline data is incorporated, the model, approximate inference and BAMDP hyperparameters should be newly tuned. To avoid too high a computational burden in our experiments, we use fixed model and approximate inference hyperparameters, highlighting in red the region where the approximate regret may be unreliable, and only tune the BAMDP hyperparameters using the approximate regret for one seed and offline dataset size (Fig. 26 in Section E.2). While we deploy each policy in the true environment to validate our approximate regret, in practice, the policy would only be deployed once the approximate regret is sufficiently low. Fig. 4, showing results for halfcheetah-full-replay, along with all other results found in Section E.2, confirm that *SOReL’s approximate regret is a good proxy for the true regret*, allowing for the safe deployment of the Bayes-optimal policy. Using the regret and PIL, all hyperparameters can be tuned entirely offline and the practitioner can identify any issues (whether with the offline-dataset, the approximate inference method, or the model) prior to deployment. We also highlight the *generalisability* of our algorithm: while the policy used to collect the halfcheetah dataset achieves an expected episodic return of around 1800 (Fig. 30 in Section E), *SOReL’s* policy (learned on a subset of the offline dataset) achieves a normalised regret of around 0.28 in the true environment (bottom of Fig. 4), corresponding to an undiscounted episode return of just under 2500. As expected (Section 5.1), our experiments show that the utility of the upper bound depends critically on the model being accurate enough relative to the discount factor. More details on a non-trivial upper-bound, along with results for gymtax and the remaining brax environments and ablations of different ensemble metrics that can be used to approximate regret with varying degrees of conservatism are found in Section E.

7 CONCLUSION

High online sample complexity and lack of performance guarantees of existing methods present a major barrier to the widespread adoption of offline RL. In this paper, we introduce *SOReL* and *TOReL*, two theoretically grounded approaches to tackle these core issues. For *SOReL*, we introduce a model-based Bayesian approach for offline RL and exploit predictive uncertainty to approximate regret. To tune hyperparameters and ensure accurate regret quantification, we minimise the PIL. In *TOReL*, we extend our fully offline hyperparameter tuning algorithm to general offline RL methods. Our empirical evaluations confirm *SOReL* is a reliable method for safe offline RL with accurate regret quantification and *TOReL* achieves near-oracle performance with *offline data alone*, resulting in significant savings in online samples for hyperparameter tuning without sacrificing performance.

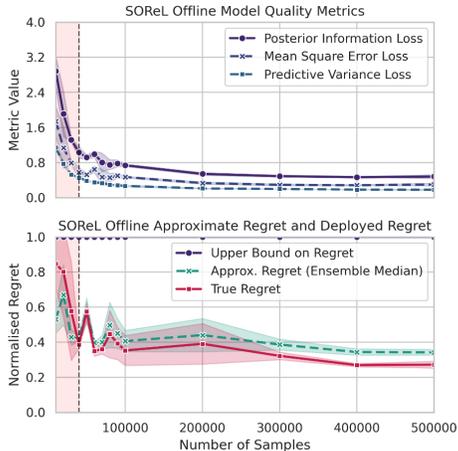


Figure 4: *SOReL* applied to brax-halfcheetah-full-replay to identify when the policy can be deployed. Shaded red indicates where $\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) \not\approx \mathcal{V}(\mathcal{D}_N)$ (for a threshold of 0.25), and hence the approximate regret may be unreliable. Mean and standard deviation given over 3 seeds.

REFERENCES

- 486
487
488 Alekh Agarwal, Nan Jiang, and Sham M. Kakade. Reinforcement learning: Theory and algorithms.
489 2022. URL https://rltheorybook.github.io/rltheorybook_AJKS.pdf. 4.3
- 490 J. Aitchison. Goodness of prediction fit. *Biometrika*, 62(3):547–554, 1975. ISSN 00063444,
491 14643510. URL <http://www.jstor.org/stable/2335509>. D.2
- 492
493 Ahmed M. Alaa and Mihaela van der Schaar. Bayesian nonparametric causal inference: Information
494 rates and learning algorithms. *IEEE Journal of Selected Topics in Signal Processing*, 12(5):
495 1031–1046, 2018. doi: 10.1109/JSTSP.2018.2848230. D.2
- 496 Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based of-
497 fline reinforcement learning with diversified q-ensemble. In M. Ranzato, A. Beygelz-
498 imer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural*
499 *Information Processing Systems*, volume 34, pages 7436–7447. Curran Associates, Inc.,
500 2021. URL [https://proceedings.neurips.cc/paper_files/paper/2021/](https://proceedings.neurips.cc/paper_files/paper/2021/file/3d3d286a8d153a4a58156d0e02d8570c-Paper.pdf)
501 [file/3d3d286a8d153a4a58156d0e02d8570c-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/3d3d286a8d153a4a58156d0e02d8570c-Paper.pdf). 3
- 502 Mihaela Aslan. Asymptotically minimax bayes predictive densities. *The Annals of Statistics*, 34(6):
503 2921–2938, 2006. ISSN 00905364. URL <http://www.jstor.org/stable/25463538>.
504 D.2, D.3
- 505
506 Philip J Ball, Cong Lu, Jack Parker-Holder, and Stephen Roberts. Augmented world models facilitate
507 zero-shot dynamics generalization from a single offline environment. In Marina Meila and Tong
508 Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume
509 139 of *Proceedings of Machine Learning Research*, pages 619–629. PMLR, 18–24 Jul 2021. URL
510 <https://proceedings.mlr.press/v139/ball121a.html>. 1, 3
- 511 Andrew R Barron. Information-theoretic characterization of bayes performance and the choice of
512 priors in parametric and nonparametric problems. In *Bayesian Statistics 6: Proceedings of the*
513 *Sixth Valencia International Meeting June 6-10, 1998*. Oxford University Press, 08 1999. ISBN
514 9780198504856. doi: 10.1093/oso/9780198504856.003.0002. URL [https://doi.org/10.](https://doi.org/10.1093/oso/9780198504856.003.0002)
515 [1093/oso/9780198504856.003.0002](https://doi.org/10.1093/oso/9780198504856.003.0002). D.2, D.3
- 516 A.R. Barron. *The Exponential Convergence of Posterior Probabilities with Implications for Bayes*
517 *Estimators of Density Functions*. Department of Statistics, University of Illinois, 1988. URL
518 <https://books.google.co.uk/books?id=8raEnQAACAAJ>. D.2, D.3
- 519
520 R.F. Bass. *Real Analysis for Graduate Students*, chapter 21. Createspace Ind Pub, 2013. ISBN
521 9781481869140. URL <https://books.google.co.uk/books?id=s6mVlgEACAAJ>.
522 D.3, 2, 3
- 523 Jacob Beck, Risto Vuorio, Evan Zheran Liu, Zheng Xiong, Luisa Zintgraf, Chelsea Finn, and Shimon
524 Whiteson. A survey of meta-reinforcement learning, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2301.08028)
525 [2301.08028](https://arxiv.org/abs/2301.08028). 4.2
- 526
527 Richard Bellman. A problem in the sequential design of experiments. *Sankhyā: The Indian Journal*
528 *of Statistics (1933-1960)*, 16(3/4):221–229, 1956. ISSN 00364452. URL [http://www.jstor.](http://www.jstor.org/stable/25048278)
529 [org/stable/25048278](http://www.jstor.org/stable/25048278). 2.1
- 530 Richard Bellman. Dynamic programming and stochastic control processes. *Informa-*
531 *tion and Control*, 1(3):228–239, 1958. ISSN 0019-9958. doi: [https://doi.org/10.](https://doi.org/10.1016/S0019-9958(58)80003-0)
532 [1016/S0019-9958\(58\)80003-0](https://doi.org/10.1016/S0019-9958(58)80003-0). URL [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S0019995858800030)
533 [article/pii/S0019995858800030](https://www.sciencedirect.com/science/article/pii/S0019995858800030). 2.1
- 534 Blair Bilodeau, Dylan J. Foster, and Daniel M. Roy. Minimax rates for conditional density
535 estimation via empirical entropy. *The Annals of Statistics*, 2021. URL [https://api.](https://api.semanticscholar.org/CorpusID:237592759)
536 [semanticscholar.org/CorpusID:237592759](https://api.semanticscholar.org/CorpusID:237592759). D.2, D.3
- 537
538 J. L. Bretagnolle and Catherine Huber. Estimation des densités: risque minimax. *Zeitschrift für*
539 *Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 47:119–137, 1978. URL [https://api.](https://api.semanticscholar.org/CorpusID:122597694)
[semanticscholar.org/CorpusID:122597694](https://api.semanticscholar.org/CorpusID:122597694). D.1, 1

- 540 Jiayu Chen, Wentse Chen, and Jeff Schneider. Bayes adaptive monte carlo tree search for offline
541 model-based reinforcement learning, 2024. URL <https://arxiv.org/abs/2410.11234>.
542 3
- 543
- 544 Kamil Ciosek, Vincent Fortuin, Ryota Tomioka, Katja Hofmann, and Richard Turner. Conservative
545 uncertainty estimation by fitting prior networks. In *Eighth International Conference on Learning*
546 *Representations*, 04 2020. 4.2, C.1
- 547 B.S. Clarke and A.R. Barron. Information-theoretic asymptotics of bayes methods. *IEEE transactions*
548 *on information theory*, 36(3):453–471, 1990. ISSN 0018-9448. D.2, D.3
- 549
- 550 J. L. Doob. Application of the theory of martingales. In *Le calcul des probabilités et ses applications*
551 *[The calculus of probabilities and its applications]*, number 13 in CNRS International Colloquia,
552 pages 23–27. Centre National de la Recherche Scientifique, Paris, 1949. (Lyon, France, 28 June–3
553 July 1948). MR:33460. Zbl:0041.45101. D.2
- 554
- 555 Alvin Drake. Observation of a markov process through a noisy channel. *PhD Thesis*, 1962. B
- 556 Yan Duan, John Schulman, Xi Chen, Peter Bartlett, Ilya Sutskever, and Peter Abbeel. Fast reinforce-
557 ment learning via slow reinforcement learning. 1987. 4.2
- 558
- 559 Michael O’Gordon Duff. *Optimal Learning: Computational Procedures for Bayes-Adaptive Markov*
560 *Decision Processes*. PhD thesis, 2002. AAI3039353. 1, 4.1, 4.2, B
- 561
- 562 Mattie Fellows, Brandon Kaplowitz, Christian Schroeder de Witt, and Shimon Whiteson. Bayesian
563 exploration networks. In *ICML*, 2024. 4.2
- 564
- 565 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep
566 data-driven reinforcement learning, 2020. 6.1
- 567
- 568 Scott Fujimoto and Shixiang (Shane) Gu. A minimalist approach to offline reinforcement learn-
569 ing. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors,
570 *Advances in Neural Information Processing Systems*, volume 34, pages 20132–20145. Curran Asso-
571 ciates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/a8166da05c5a094f7dc03724b41886e5-Paper.pdf. 3
- 572
- 573 Arthur Guez, David Silver, and Peter Dayan. Efficient bayes-adaptive reinforcement learning
574 using sample-based search. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Wein-
575 berger, editors, *Advances in Neural Information Processing Systems*, volume 25. Curran As-
576 sociates, Inc., 2012. URL <https://proceedings.neurips.cc/paper/2012/file/35051070e572e47d2c26c241ab88307f-Paper.pdf>. 4.2
- 577
- 578 Arthur Guez, David Silver, and Peter Dayan. Scalable and efficient bayes-adaptive reinforcement
579 learning based on monte-carlo tree search. *Journal of Artificial Intelligence Research*, 48:841–883,
580 10 2013. doi: 10.1613/jair.4117. 4.2
- 581
- 582 Arthur Guez, Nicolas Heess, David Silver, and Peter Dayan. Bayes-adaptive simulation-based search
583 with value function approximation. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and
584 K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Cur-
585 ran Associates, Inc., 2014. URL https://proceedings.neurips.cc/paper_files/paper/2014/file/74d863ca4a12ccca50a754a3b277dbf7-Paper.pdf. 3
- 586
- 587 J. A. Hartigan. The maximum likelihood prior. *The Annals of statistics*, 26(6):2083–2103, 1998.
588 ISSN 0090-5364. D.2, D.3
- 589
- 590 Matthew Thomas Jackson, Uljad Berdica, Jarek Liesen, Shimon Whiteson, and Jakob Nicolaus
591 Foerster. A clean slate for offline reinforcement learning. *arXiv preprint arXiv:2504.11453*, 2025.
592 1, 1, 2.1, 3, 4.2, 6.1, F.2, F.4
- 593
- A Jesson, C Lu, N Beltran-Velez, A Filos, J Foerster, and Y Gal. Relu to the rescue: Improve your
on-policy actor-critic with positive advantages. 2024. 10

- 594 Leslie Pack Kaelbling, Michael L. Littman, and Anthony R. Cassandra. Planning and acting in
595 partially observable stochastic domains. *Artif. Intell.*, 101(1–2):99–134, may 1998. ISSN 0004-
596 3702. B
- 597 Robert E Kass, Luke Tierney, and Joseph B Kadane. The validity of posterior expansions based on
598 laplace’s method. *Bayesian and Likelihood Methods in Statistics and Economics*, pages 473–488,
599 1990. URL <https://www.stat.cmu.edu/~kass/papers/validity.pdf>. D.3
- 601 Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel:
602 Model-based offline reinforcement learning. In H. Larochelle, M. Ranzato, R. Hadsell,
603 M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33,
604 pages 21810–21823. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/file/
605 f7efa4f864ae9b88d43527f4b14f750f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/f7efa4f864ae9b88d43527f4b14f750f-Paper.pdf). 1, 1, 3, 4.3, 5, 5.2, 6.1,
606 D.4
- 608 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
609 *arXiv:1412.6980*, 2014. F.2
- 610 Fumiyasu Komaki. On asymptotic properties of predictive distributions. *Biometrika*, 83(2):299–313,
611 06 1996. ISSN 0006-3444. doi: 10.1093/biomet/83.2.299. URL [https://doi.org/10.
612 1093/biomet/83.2.299](https://doi.org/10.1093/biomet/83.2.299). D.2, D.3
- 614 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-
615 learning. *CoRR*, abs/2110.06169, 2021. URL <https://arxiv.org/abs/2110.06169>. 1,
616 1, 3, 6.1
- 617 Danie G. Krige. A statistical approach to some basic mine valuation problems on the witwatersand.
618 *Journal of the Chemical, Metallurgical and Mining Society of South Africa*, 52:119–139, 1951.
619 4.1, B
- 620 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline
621 reinforcement learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors,
622 *Advances in Neural Information Processing Systems*, volume 33, pages 1179–1191. Curran Associ-
623 ates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/
624 2020/file/0d2b2061826a5df3221116a5085a6052-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/0d2b2061826a5df3221116a5085a6052-Paper.pdf). 3
- 626 Sascha Lange, Thomas Gabel, and Martin Riedmiller. *Batch Reinforcement Learning*, pages 45–73.
627 Springer Berlin Heidelberg, Berlin, Heidelberg, 2012. ISBN 978-3-642-27645-3. doi: 10.1007/
628 978-3-642-27645-3_2. URL https://doi.org/10.1007/978-3-642-27645-3_2. 1,
629 2.1
- 630 Lucien Le Cam. On some asymptotic properties of maximum likelihood estimates and related bayes’
631 estimates. volume 1, pages 277–300, 1953. D.2, D.3
- 632 Yann LeCun, Leon Bottou, Genevieve B Orr, and Klaus-Robert Mueller. Efficient backprop. In
633 *Neural Networks: Tricks of the Trade*, pages 9–50. Springer, 1998. F.2
- 634 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,
635 review, and perspectives on open problems, 2020. URL [https://arxiv.org/abs/2005.
636 01643](https://arxiv.org/abs/2005.01643). 1, 2.1
- 637 D. V. Lindley. Approximate bayesian methods. *Trabajos de Estadística Y de Investigación Operativa*,
638 31(1):223–245, 1980. D.3
- 641 Chris Lu, Jakub Kuba, Alistair Letcher, Luke Metz, Christian Schroeder de Witt, and Jakob Foerster.
642 Discovered policy optimisation. *Advances in Neural Information Processing Systems*, 35:16455–
643 16468, 2022a. F.3
- 644 Cong Lu, Philip Ball, Jack Parker-Holder, Michael Osborne, and Stephen J. Roberts. Revisiting design
645 choices in offline model based reinforcement learning. In *International Conference on Learning*
646 *Representations*, 2022b. URL <https://openreview.net/forum?id=zz9hXVhf40>. 1,
647 3

- 648 J. J. Martin. *Bayesian decision problems and Markov chains [by] J. J. Martin*. Wiley New York,
649 1967. 1, 4.1, 4.2, B
- 650
- 651 Kevin Murphy. Reinforcement learning: An overview, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2412.05265)
652 2412.05265. 1, 2.1
- 653
- 654 Ian Osband and Benjamin Van Roy. Why is posterior sampling better than optimism for reinforcement
655 learning? In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International*
656 *Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages
657 2701–2710. PMLR, 06–11 Aug 2017. URL [https://proceedings.mlr.press/v70/](https://proceedings.mlr.press/v70/osband17a.html)
658 osband17a.html. 4.2
- 659 Ian Osband, John Aslanides, and Albin Cassirer. Randomized prior functions for deep rein-
660 forcement learning. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi,
661 and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages
662 8617–8629. Curran Associates, Inc., 2018. URL [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/8080-randomized-prior-functions-for-deep-reinforcement-learning.pdf)
663 8080-randomized-prior-functions-for-deep-reinforcement-learning.
664 pdf. 4.2, 6.2, C.1
- 665 Tom Le Paine, Cosmin Paduraru, Andrea Michi, Çağlar Gülçehre, Konrad Zolna, Alexander Novikov,
666 Ziyu Wang, and Nando de Freitas. Hyperparameter selection for offline reinforcement learning.
667 *CoRR*, abs/2007.09055, 2020. URL <https://arxiv.org/abs/2007.09055>. 3
- 668
- 669 K. B. Petersen and M. S. Pedersen. The matrix cookbook, nov 2012. URL [http://localhost/](http://localhost/pubdb/p.php?3274)
670 pubdb/p.php?3274. Version 20121115. 3
- 671 Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John
672 Wiley & Sons, Inc., USA, 1st edition, 1994. ISBN 0471619779. 2.1
- 673
- 674 Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*.
675 The MIT Press, 2006. URL <https://gaussianprocess.org/gpml/>. 4.1, B
- 676
- 677 Gareth O. Roberts and Jeffrey S. Rosenthal. General state space Markov chains and MCMC
678 algorithms. *Probability Surveys*, 1(none):20–71, 2004. doi: 10.1214/154957804100000024. URL
679 <https://doi.org/10.1214/154957804100000024>. D.3
- 680
- 681 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
682 optimization algorithms. *CoRR*, abs/1707.06347, 2017. 1, 4.2, 6.2
- 683
- 684 Anya Sims, Cong Lu, Jakob Nicolaus Foerster, and Yee Whye Teh. The edge-of-reach problem in
685 offline model-based reinforcement learning. In *The Thirty-eighth Annual Conference on Neural*
686 *Information Processing Systems*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=3dn1hINA6o)
687 3dn1hINA6o. 1, 3
- 688
- 689 Richard D. Smallwood and Edward J. Sondik. The optimal control of partially observable markov
690 processes over a finite horizon. *Operations Research*, 21(5):1071–1088, 1973. ISSN 0030364X,
691 15265463. URL <http://www.jstor.org/stable/168926>. B
- 692
- 693 Matthew Smith, Lucas Maystre, Zhenwen Dai, and Kamil Ciosek. A strong baseline for batch
694 imitation learning, 2023. URL <https://arxiv.org/abs/2302.02788>. 3
- 695
- 696 Bharath K. Sriperumbudur, Kenji Fukumizu, Arthur Gretton, Bernhard Scholkopf, and Gert R. G.
697 Lanckriet. On integral probability metrics, ϕ -divergences and binary classification. *arXiv:*
698 *Information Theory*, 2009. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:14114329)
699 14114329. D.1
- 700
- 701 Yihao Sun, Jiayi Zhang, Chengxing Jia, Haoxin Lin, Junyin Ye, and Yang Yu. Model-Bellman
702 inconsistency for model-based offline reinforcement learning. In Andreas Krause, Emma
703 Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors,
704 *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Pro-*
705 *ceedings of Machine Learning Research*, pages 33177–33194. PMLR, 23–29 Jul 2023. URL
706 <https://proceedings.mlr.press/v202/sun23q.html>. 1, 3

- 702 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press,
703 second edition, 2018. URL <http://incompleteideas.net/book/the-book-2nd.html>. 2.1
- 704
705 Csaba Szepesvári. Algorithms for Reinforcement Learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 4(1):1–103, 2010. ISSN 1939-4608. doi: 10.2200/S00268ED1V01Y201005AIM009. URL <http://www.morganclaypool.com/doi/abs/10.2200/S00268ED1V01Y201005AIM009>. 2.1
- 706
707
708
709 Denis Tarasov, Vladislav Kurenkov, Alexander Nikulin, and Sergey Kolesnikov. Revisiting the minimalist approach to offline reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=vqGwslLeEw>. 1, 1, 3, 6.1
- 710
711
712
713 Luke Tierney and Joseph B. Kadane. Accurate approximations for posterior moments and marginal densities. *Journal of the American Statistical Association*, 81(393):82–86, 1986. ISSN 0162-1459. D.3
- 714
715
716
717 Luke Tierney, Robert E. Kass, and Joseph B. Kadane. Fully exponential laplace approximations to expectations and variances of nonpositive functions. *Journal of the American Statistical Association*, 84(407):710–716, 1989. ISSN 0162-1459. D.3
- 718
719
720
721 A. W. van der Vaart. *Bayes Procedures*, page 138–152. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1998. doi: 10.1017/CBO9780511802256.011. D.2, D.3
- 722
723
724
725 Aad van der Vaart and Harry van Zanten. Information rates of nonparametric gaussian process methods. *J. Mach. Learn. Res.*, 12(null):2095–2119, July 2011. ISSN 1532-4435. D.2, D.3
- 726
727
728
729
730 Han Wang, Archit Sakhadeo, Adam M White, James M Bell, Vincent Liu, Xutong Zhao, Puer Liu, Tadashi Kozuno, Alona Fyshe, and Martha White. No more pesky hyperparameters: Offline hyperparameter tuning for RL. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=AiOUi3440V>. 3
- 731
732
733
734 Ruosong Wang, Dean P. Foster, and Sham M. Kakade. What are the statistical limits of offline rl with linear function approximation?, 2020. URL <https://arxiv.org/abs/2010.11895>. 3
- 735
736
737 Norbert Wiener. Differential-space. *Journal of Mathematics and Physics*, 2(1-4):131–174, 1923. doi: <https://doi.org/10.1002/sapm192321131>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/sapm192321131>. 4.1, B
- 738
739
740
741 Yuhong Yang and Andrew Barron. Information-theoretic determination of minimax rates of convergence. *The Annals of Statistics*, 27(5):1564 – 1599, 1999. doi: 10.1214/aos/1017939142. URL <https://doi.org/10.1214/aos/1017939142>. D.2, D.3
- 742
743
744
745
746 Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 14129–14142. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/a322852ce0df73e204b7e67cbbef0d0a-Paper.pdf. 1, 3, 4.3, 5.2, 6.1
- 747
748
749
750
751
752 Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn. Combo: Conservative offline model-based policy optimization. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 28954–28967. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/f29a179746902e331572c483c45e5086-Paper.pdf. 1, 3
- 753
754
755 Baohe Zhang, Raghu Rajan, Luis Pineda, Nathan Lambert, André Biedenkapp, Kurtland Chua, Frank Hutter, and Roberto Calandra. On the importance of hyperparameter optimization for model-based reinforcement learning. In Arindam Banerjee and Kenji Fukumizu, editors, *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of

756 *Proceedings of Machine Learning Research*, pages 4015–4023. PMLR, 13–15 Apr 2021. URL
757 <https://proceedings.mlr.press/v130/zhang21n.html>. 1
758
759 Luisa Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarin Gal, Katja Hofmann, and
760 Shimon Whiteson. Varibad: A very good method for bayes-adaptive deep rl via meta-learning. In
761 *International Conference on Learning Representations*, 2020. URL [https://openreview.](https://openreview.net/pdf?id=Hkl9JlBYvr)
762 [net/pdf?id=Hkl9JlBYvr](https://openreview.net/pdf?id=Hkl9JlBYvr). 4.2
763 Åström, Karl Johan. Optimal Control of Markov Processes with Incomplete State Information I.
764 10:174–205, 1965. ISSN 0022-247X. doi: {10.1016/0022-247X(65)90154-X}. URL [https:](https://lup.lub.lu.se/search/files/5323668/8867085.pdf)
765 [//lup.lub.lu.se/search/files/5323668/8867085.pdf](https://lup.lub.lu.se/search/files/5323668/8867085.pdf). B
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A BROADER IMPACT

This paper presents work whose goal is to improve the safety and efficacy of offline RL. Our work therefore takes significant steps towards the development of safe offline RL methods with accurate regret guarantees that act in ways that are predictable. We also hope our paper helps to re-frame the discourse on offline RL to focus on safety, practical efficacy and sample efficiency.

More generally, any advancement in RL should be seen in the context of general advancements to machine learning. Whilst machine learning has the potential to develop useful tools to benefit humanity, it must be carefully integrated into underlying political and social systems to avoid negative consequences for people living within them. A discussion of this complex topic lies beyond the scope of this work.

B PRIMER ON BAYESIAN RL

A Bayesian epistemology characterises the agent’s uncertainty in the MDP through distributions over any unknown variable (Martin, 1967; Duff, 2002). In our learning problem, a Bayesian first specifies a model $P_{R,S}(s_t, a_t)$ over the unknown state transition and reward distribution, representing a hypothesis space of possible environment dynamics. We focus on a parametric model: $p(r_t, s_{t+1} | s_t, a_t, \theta)$ with each $\theta \in \Theta \subseteq \mathbb{R}^d$ representing a hypothesis about the MDP \mathcal{M}^* , however our results can easily be generalised to non-parametric methods like Gaussian process regression (Rasmussen and Williams, 2006; Wiener, 1923; Krige, 1951). A prior distribution over the parameter space P_Θ is specified, which represents the initial *a priori* belief in the true value of $P_{R,S}^*(s, a)$ before the agent has observed any transitions. Priors are a powerful aspect of Bayesian RL, allowing practitioners to provide the agent with any information about the MDP and transfer knowledge between agents and domains. Given a history h_t , the prior is updated to a posterior $P_\Theta(h_t)$, representing the agent’s beliefs in the MDP’s dynamics once h_t has been observed. For each history, the posterior is used to *marginalise* across all hypotheses according to the agent’s uncertainty, yielding the predictive state transition-reward distribution $P_{R,S}(h_t, a_t) = \mathbb{E}_{\theta \sim P_\Theta(h_t)} [P_{R,S}(s_t, a_t, \theta)]$ which characterise the epistemic and aleatoric uncertainty in $P_{R,S}^*(s_t, a_t)$. Given $P_{R,S}(h_t, a_t)$, we reason over counterfactual future trajectories using the predictive distribution over trajectories P_t^π and define the BRL objective as:

$$J_{\text{Bayes}}^\pi(P_\Theta) := \mathbb{E}_{h_\infty \sim P_\infty^\pi} \left[\sum_{i=0}^{\infty} \gamma^i r_i \right].$$

Let $\Pi_{\mathcal{H}}$ denote the space of all history-conditioned policies. A corresponding optimal policy is known as a Bayes-optimal policy, which we denote as $\pi_{\text{Bayes}}^*(\cdot) \in \Pi_{\text{Bayes}}^*(P_\Theta) := \arg \max_{\pi \in \Pi_{\mathcal{H}}} J_{\text{Bayes}}^\pi(P_\Theta)$. Unlike in frequentist RL, Bayesian variables depend on histories obtained through posterior marginalisation; hence the posterior is often known as the *belief state*, which augments each ground state s_t like in a partially observable Markov decision process (Drake, 1962; Åström, Karl Johan, 1965; Smallwood and Sondik, 1973; Kaelbling et al., 1998). Analogously to the state-transition distribution in frequentist RL, we define a *belief transition* distribution $P_{\mathcal{H}}(h_t, a_t)$ using the predictive state transition-reward distribution, which yields *Bayes-adaptive MDP* (BAMDP) (Duff, 2002): $\mathcal{M}_{\text{Bayes}}(P_\Theta) := \langle \mathcal{H}, \mathcal{A}, P_0, P_{\mathcal{H}}(h, a), \gamma \rangle$. The BAMDP is solved using planning methods to obtain a Bayes-optimal policy, which naturally balances exploration with exploitation: after every timestep, the agent’s uncertainty is characterised via the posterior conditioned on the history h_t , which includes all future trajectories to marginalise over. Via the belief transition, the BRL objective accounts for how the posterior evolves on every timestep, and hence any Bayes-optimal policy π_{Bayes}^* is optimal not only according to the epistemic uncertainty of a fixed belief but accounts for how the epistemic uncertainty evolves at every future timestep, decaying according to the discount factor.

C DERIVATIONS

C.1 NEGATIVE LOG LIKELIHOOD LOSS FUNCTION WITH APPROXIMATE INFERENCE

Assume a dataset of N input-output pairs:

$$\mathcal{D}_N := \{(x_0, y_0), (x_1, y_1), \dots, (x_{N-1}, y_{N-1})\},$$

and a multivariate Gaussian regression model:

$$p(y|x, \theta) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_\theta(x)|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (y - \mu_\theta(x)) \Sigma_\theta^{-1} (y - \mu_\theta(x))^T \right),$$

where D is the number of dimensions. Here our model $\text{NN}_\theta : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{Y})$ is a neural network parametrised by $\theta \in \Theta$ that outputs a Gaussian distribution $\mathcal{N}(\mu_\theta, \Sigma_\theta)$ over \mathcal{Y} . Assuming independent dimensions, such that the covariance matrix is diagonal:

$$p(y|x, \theta) = \prod_{d=0}^{D-1} \frac{1}{\sqrt{2\pi\sigma_{\theta_d}^2(x)}} \exp\left(-\frac{1}{2\sigma_{\theta_d}^2(x)}(y_d - \mu_{\theta_d}(x))^2\right).$$

We can then fit our model by minimising the negative log likelihood loss:

$$\begin{aligned} \mathcal{L}(\text{NLL}(\theta)) &:= -\log p(\mathcal{D}_N|\theta), \\ &= \sum_{i=0}^{N-1} \left(\frac{D}{2} \log(2\pi) + \frac{1}{2} \sum_{d=0}^{D-1} \left(\log \sigma_{\theta_d}^2(x_i) + \frac{(y_{i_d} - \mu_{\theta_d}(x_i))^2}{\sigma_{\theta_d}^2(x_i)} \right) \right), \\ &= \left[\sum_{i=0}^{N-1} \frac{1}{N} \left(\frac{D}{2} \log(2\pi) + \frac{1}{2} \sum_{d=0}^{D-1} \left(\log \sigma_{\theta_d}^2(x_i) + \frac{(y_{i_d} - \mu_{\theta_d}(x_i))^2}{\sigma_{\theta_d}^2(x_i)} \right) \right) \right], \\ &= \mathbb{E}_{i \sim \mathcal{U}_N} \left[\frac{D}{2} \log(2\pi) + \frac{1}{2} \sum_{d=0}^{D-1} \left(\log \sigma_{\theta_d}^2(x_i) + \frac{(y_{i_d} - \mu_{\theta_d}(x_i))^2}{\sigma_{\theta_d}^2(x_i)} \right) \right], \\ &\doteq \mathbb{E}_{i \sim \mathcal{U}_N} \left[\sum_{d=0}^{D-1} \left(\log \sigma_{\theta_d}^2(x_i) + \frac{(y_{i_d} - \mu_{\theta_d}(x_i))^2}{\sigma_{\theta_d}^2(x_i)} \right) \right], \end{aligned}$$

where recall \mathcal{U}_N is the uniform distribution over $\{0, 1, \dots, N-1\}$. Our final line means equality up to a constant, as we can ignore the $\frac{D}{2} \log(2\pi)$ term for optimisation because it is independent of θ .

We use RP ensembles for our approximate posterior (Osband et al., 2018; Ciosek et al., 2020); here an ensemble of M separate model weights $\{\theta_0, \theta_1, \dots, \theta_{M-1}\}$ are randomly initialised and are optimised in parallel, summing over the corresponding negative log likelihoods. When training, we optimise the log-variance rather than the variance for numerical stability and to ensure that the variance remains positive. This allows us to simultaneously optimise maximum and minimum log-variance parameters for each dimension across the ensemble, which we use to soft-clamp the log-variances output by individual models, preventing any individual model becoming overly confident or too uncertain in one dimension. Our final loss function is then given by:

$$\begin{aligned} \mathcal{L}(\theta, \mathcal{D}_N) &= \sum_{j=0}^{M-1} \left(\mathbb{E}_{i \sim \mathcal{U}_N} \left[\sum_{d=0}^{D-1} \left(\xi_{\theta_{d_j}}(x_i) + \frac{(y_{i_d} - \mu_{\theta_{d_j}}(x_i))^2}{\exp(\xi_{\theta_{d_j}}(x_i))} \right) \right] \right) \\ &\quad + c \cdot \sum_{d=0}^{D-1} (\xi_{\theta_{d_{\max}}} - \xi_{\theta_{d_{\min}}}), \end{aligned}$$

where $\xi_{\theta_{d_j}} = \xi_{\theta_{d_{\min}}} + \left[1 + \exp(\xi_{\theta_{d_{\max}}} - [1 + \exp(\log \sigma_{\theta_{d_j}}^2 - \xi_{\theta_{d_{\max}}})] - \xi_{\theta_{d_{\min}}}) \right]$ and $\xi_{\theta_{d_{\min}}}$ and $\xi_{\theta_{d_{\max}}}$ are respectively the minimum and maximum log-variance parameters optimised across the ensemble, c is the log-variance difference coefficient used to control the clamping term, and M is the number of models in the ensemble. $\mathcal{L}(\theta, \mathcal{D}_N)$ can be minimised by using Monte Carlo minibatch gradient descent with a minibatch \mathcal{M}_n of $n < N$ samples drawn uniformly from \mathcal{D}_N .

C.2 REGRET APPROXIMATORS

Predictive Variance: We now show how the true regret can be approximated using the Bayesian predictive variance of returns. We start with the bound on from Ineq. 9. Defining the discounted return $R(h_\infty) := \sum_{i=0}^{\infty} \gamma^i r_i$:

$$\begin{aligned} |J^\pi(\mathcal{M}^*) - J_{\text{Bayes}}^\pi(P_\Theta(\mathcal{D}_N))| &= |J^\pi(\mathcal{M}^*) - J_{\text{Bayes}}^\pi(P_\Theta(\mathcal{D}_N))|, \\ &= |J^\pi(\mathcal{M}^*) - \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)]|, \\ &= |\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [J^\pi(\mathcal{M}^*) - R(h_\infty)]|, \\ &= \sqrt{\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [J^\pi(\mathcal{M}^*) - R(h_\infty)]^2}. \end{aligned}$$

Applying Jensen’s inequality:

$$|J^\pi(\mathcal{M}^*) - J_{\text{Bayes}}^\pi(P_\Theta(\mathcal{D}_N))| \leq \sqrt{\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [(J^\pi(\mathcal{M}^*) - R(h_\infty))^2]}. \quad (6)$$

We now recognise that the mean squared error term in Eq. (6) relies on knowing the true MDP dynamics $J^\pi(\mathcal{M}^*)$. We can approximate this term using the predictive variance over returns:

$$\begin{aligned} & \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [(J^\pi(\mathcal{M}^*) - R(h_\infty))^2] \\ & \approx \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [(J_{\text{Bayes}}^\pi(\mathcal{D}_N) - R(h_\infty))^2], \\ & = \mathbb{V}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)], \end{aligned}$$

which can be estimated using the dataset \mathcal{D}_N , yielding:

$$|J^\pi(\mathcal{M}^*) - J_{\text{Bayes}}^\pi(P_\Theta(\mathcal{D}_N))| \leq \sqrt{\mathbb{V}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)]}.$$

Finally, using Ineq. 9 this justifies our approximation for estimating the regret:

$$\text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \approx 2\sqrt{\mathbb{V}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)]}.$$

To improve the approximation, we conservatively upper-bound the regret based on alternate ensemble statistics with varying degrees of conservatism to prevent associating a low regret with a policy that performs equally poorly in all members of the ensemble:

$$\begin{aligned} \text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \approx \max & \left[2\sqrt{\mathbb{V}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)]}, \right. \\ & \left. \hat{R}_{\max} - \hat{\mathbb{M}}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)] \right], \end{aligned}$$

for example, here $\hat{\mathbb{M}}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)]$ denotes the median predictive returns based on sampling from the (approximate) posterior and rolling out the Bayes-optimal policy. \hat{R}_{\max} is estimated from the maximum return in the offline dataset. In Fig. 27 and Fig. 29, we plot different ensemble statistics (the ensemble mean, median, maximum and minimum regrets) that can be used to inform the approximate regret: we shade the regret based on the range of these statistics in purple. As long as the true environment falls in the space spanned by the posterior (model ensemble), the true regret is guaranteed to lie within this range. By ensuring that the (normalised) predictive variance is at least as large as the (normalised) MSE in the PIL, the space spanned by the approximate posterior via model ensemble is approximately large enough, relative to the model error. Below we order different ensemble statistics from least to most conservative:

$$\begin{aligned} \text{Regret}(\mathcal{M}^*, \mathcal{D}_N) & \approx \hat{R}_{\max} - \hat{R}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)], \\ & \leq \hat{R}_{\max} - \hat{\mathbb{E}}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)], \\ & \approx \hat{R}_{\max} - \hat{\mathbb{M}}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)], \\ & \leq \hat{R}_{\max} - \hat{r}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)]. \end{aligned}$$

Here $\hat{R}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)]$, $\hat{\mathbb{E}}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)]$, and $\hat{r}_{\theta \sim P_\Theta(\mathcal{D}_N), h_\infty \sim P_\infty^\pi(\theta)} [R(h_\infty)]$, respectively denote the maximum, mean, and minimum predictive returns based on sampling from the (approximate) posterior and rolling out the Bayes-optimal policy. We note that the minimum predictive return leads to the maximum predictive regret: "Ensemble Max" in Fig. 27 and Fig. 29 refer to the maximum predictive *regrets* rather than maximum predictive *returns*. Empirically, we find that the ensemble median alone is a good proxy for the true regret, being neither overly conservative, nor overly susceptible to being skewed by a policy that performs well on only a subset of posterior samples (as the mean might be). Using the variance, a policy with high variance but also high mean will be associated with a high approximate regret, which is the case for brax-hopper-full-replay (Fig. 29c), where the variance actually overestimates the regret.

972 D PROOFS

973 D.1 PRIMER ON TOTAL VARIATIONAL DISTANCE

974 We measure distance between two probability distributions P_X and Q_X using the total variational
975 (TV) distance, defined as:

$$976 \text{TV}(P_X \| Q_X) := \sup_E |P_X(E) - Q_X(E)|.$$

977 The TV distance takes the supremum over all events E to find the event that gives rise to the
978 maximum difference in probability between two distributions. A key property of the TV distance
979 is: $0 \leq \text{TV}(P_X \| Q_X) \leq 1$. If $\text{TV}(P_X \| Q_X) = 0$, then $P_X = Q_X$ as there is no event that
980 both distributions don't assign the same probability mass to. If $\text{TV}(P_X \| Q_X) = 1$, then the dis-
981 tributions assign completely different mass to at least one event. The TV distance can be related
982 to the Kullback-Leibler (KL) divergence using the Bretagnolle-Huber (Bretagnolle and Huber,
983 1978) inequality: $\text{TV}(P_X \| Q_X) \leq \sqrt{1 - \exp(-\text{KL}(P_X \| Q_X))} \leq 1$, which preserves the property
984 that $0 \leq \text{TV}(P_X \| Q_X) \leq 1$. The TV distance can be shown (Sriperumbudur et al., 2009) to be
985 equivalent to the integral probability metric under the ∞ -norm, which we will make use of in our
986 theorems:

$$987 \text{TV}(P_X \| Q_X) = \frac{1}{2} \sup_{f \in \mathcal{F}: \mathcal{X} \rightarrow [-1,1]} |\mathbb{E}_{x \sim P_X} [f(x)] - \mathbb{E}_{x \sim Q_X} [f(x)]|, \quad (7)$$

988 In this form, the supremum is taken over the space of all functions that are bounded by unity, that is
989 $\|f\|_\infty = 1$.

990 D.2 PROOF OF THEOREM 1

991 Let the predictive distribution over history h_t using the posterior $P_\Theta(\mathcal{D}_N)$ be $P_{t,\pi}(\mathcal{D}_N)$, which has
992 density:

$$993 p_\pi(h_t, \mathcal{D}_N) := p_0(s_0) \prod_{i=0}^{t-1} \pi(a_i | h_i) p(r_i | h_i, a_i, \mathcal{D}_N) p(s_{i+1} | h_i, a_i, \mathcal{D}_N).$$

994 According to the Bernstein-von Mises theorem (Doob, 1949; Le Cam, 1953; Vaart, 1998), as the
995 posterior becomes more informative it concentrates around a smaller (and more tractable) subset of
996 the hypothesis space. Not only does this ease the computational burden of solving the BRL objective,
997 but in the limit $N \rightarrow \infty$, the Bayesian RL objective using the true posterior will approach the true
998 expected discounted return for the MDP: $J^\pi(P_\Theta(\mathcal{D}_N)) \xrightarrow{N \rightarrow \infty} J^\pi(\mathcal{M}^*)$. In this limit, any Bayes-
999 optimal policy will be an optimal policy for the true MDP, achieving the highest expected returns
1000 once deployed. To make progress towards quantifying how much offline data we need to achieve an
1001 acceptable level of regret, we first relate the true regret $\text{Regret}(\mathcal{M}^*, \mathcal{D}_N)$ to the TV distance between
1002 the true $P_{i,\pi}^*$ and predictive $P_{i,\pi}(\mathcal{D}_N)$ history distributions:

1003 **Lemma 1.** Let $\mathcal{R}_{\max} := \frac{(r_{\max} - r_{\min})}{1 - \gamma}$ denote the maximum possible regret for the MDP. For a prior
1004 $P_\Theta(\mathcal{D}_N)$, the true regret can be bounded as:

$$1005 \text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \leq 2\mathcal{R}_{\max} \cdot \sup_\pi \mathbb{E}_{i \sim \mathcal{G}(\gamma)} [\text{TV}(P_{i+1,\pi}^* \| P_{i+1,\pi}^\pi(\mathcal{D}_N))]. \quad (8)$$

1006 *Proof.* We start from the definition of the true regret:

$$1007 \text{Regret}(\mathcal{M}^*, \mathcal{D}_N) := J^{\pi^*}(\mathcal{M}^*) - J^{\pi_{\text{Bayes}}^*}(\mathcal{M}^*, \mathcal{D}_N).$$

1008 We now bound the difference between $J^{\pi^*}(\mathcal{M}^*)$ and $J^{\pi_{\text{Bayes}}^*}(\mathcal{M}^*, \mathcal{D}_N)$ in terms of the difference
1009 between $J^\pi(\mathcal{M}^*)$ and $J_{\text{Bayes}}^\pi(\mathcal{D}_N)$:

$$1010 \begin{aligned} & \text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \\ &= J^{\pi^*}(\mathcal{M}^*) - J_{\text{Bayes}}^{\pi^*}(P_\Phi(\mathcal{D}_N)) + J_{\text{Bayes}}^{\pi^*}(P_\Phi(\mathcal{D}_N)) - J^{\pi_{\text{Bayes}}^*}(\mathcal{M}^*, \mathcal{D}_N), \\ &\leq J^{\pi^*}(\mathcal{M}^*) - J_{\text{Bayes}}^{\pi^*}(P_\Phi(\mathcal{D}_N)) + J_{\text{Bayes}}^{\pi_{\text{Bayes}}^*}(P_\Phi(\mathcal{D}_N)) - J^{\pi_{\text{Bayes}}^*}(\mathcal{M}^*, \mathcal{D}_N), \\ &\leq \sup_\pi |J^\pi(\mathcal{M}^*) - J_{\text{Bayes}}^\pi(P_\Theta(\mathcal{D}_N))| + \sup_\pi |J_{\text{Bayes}}^\pi(P_\Theta(\mathcal{D}_N)) - J^\pi(\mathcal{M}^*)|, \\ &= 2 \sup_\pi |J^\pi(\mathcal{M}^*) - J_{\text{Bayes}}^\pi(P_\Theta(\mathcal{D}_N))|, \end{aligned} \quad (9)$$

where the second line follows from $J_{\text{Bayes}}^{\pi^*}(P_{\Phi}(\mathcal{D}_N)) \leq J_{\text{Bayes}}^{\pi^*}(P_{\Phi}(\mathcal{D}_N))$ by definition. Now our goal is to bound $|J^{\pi}(\mathcal{M}^*) - J_{\text{Bayes}}^{\pi}(P_{\Theta}(\mathcal{D}_N))|$:

$$\begin{aligned} |J^{\pi}(\mathcal{M}^*) - J_{\text{Bayes}}^{\pi}(P_{\Theta}(\mathcal{D}_N))| &= \left| \mathbb{E}_{h_{\infty} \sim P_{\infty}^*, \pi} \left[\sum_{i=0}^{\infty} \gamma^i r_i \right] - \mathbb{E}_{h_{\infty} \sim P_{\infty}^{\pi}(\mathcal{D}_N)} \left[\sum_{i=0}^{\infty} \gamma^i r_i \right] \right|, \\ &= \left| \sum_{i=0}^{\infty} \gamma^i \mathbb{E}_{h_{i+1} \sim P_{i+1}^*, \pi} [r_i] - \sum_{i=0}^{\infty} \gamma^i \mathbb{E}_{h_{i+1} \sim P_{i+1}^{\pi}(\mathcal{D}_N)} [r_i] \right|, \\ &= \left| \sum_{i=0}^{\infty} \gamma^i \left(\mathbb{E}_{h_{i+1} \sim P_{i+1}^*, \pi} [r_i] - \mathbb{E}_{h_{i+1} \sim P_{i+1}^{\pi}(\mathcal{D}_N)} [r_i] \right) \right|, \\ &\leq \sum_{i=0}^{\infty} \gamma^i \left| \mathbb{E}_{h_{i+1} \sim P_{i+1}^*, \pi} [r_i] - \mathbb{E}_{h_{i+1} \sim P_{i+1}^{\pi}(\mathcal{D}_N)} [r_i] \right|. \end{aligned}$$

Using Ineq. 10 from Lemma 2, we now bound each difference $\left| \mathbb{E}_{h_{i+1} \sim P_{i+1}^*, \pi} [r_i] - \mathbb{E}_{h_{i+1} \sim P_{i+1}^{\pi}(\mathcal{D}_N)} [r_i] \right|$ in terms of total variational distance between P_{i+1}^*, π and $P_{i+1}^{\pi}(\mathcal{D}_N)$:

$$\begin{aligned} |J^{\pi}(\mathcal{M}^*) - J_{\text{Bayes}}^{\pi}(P_{\Theta}(\mathcal{D}_N))| &\leq (r_{\max} - r_{\min}) \cdot \sum_{i=0}^{\infty} \gamma^i \text{TV} (P_{i+1}^*, \pi \| P_{i+1}^{\pi}(\mathcal{D}_N)), \\ &= \frac{r_{\max} - r_{\min}}{1 - \gamma} \cdot \sum_{i=0}^{\infty} (1 - \gamma) \gamma^i \text{TV} (P_{i+1}^*, \pi \| P_{i+1}^{\pi}(\mathcal{D}_N)), \\ &= \frac{r_{\max} - r_{\min}}{1 - \gamma} \cdot \mathbb{E}_{i \sim \mathcal{G}(\gamma)} [\text{TV} (P_{i+1}^*, \pi \| P_{i+1}^{\pi}(\mathcal{D}_N))], \\ &= \mathcal{R}_{\max} \cdot \mathbb{E}_{i \sim \mathcal{G}(\gamma)} [\text{TV} (P_{i+1}^*, \pi \| P_{i+1}^{\pi}(\mathcal{D}_N))], \end{aligned}$$

where $\mathcal{G}(\gamma)$ is the geometric distribution. Finally, substituting into Ineq. 9 yields our desired result:

$$\begin{aligned} \text{Regret}(\mathcal{M}^*, \mathcal{D}_N) &\leq 2 \sup_{\pi} \left[(\mathcal{R}_{\max} \cdot \mathbb{E}_{i \sim \mathcal{G}(\gamma)} [\text{TV} (P_{i+1}^*, \pi \| P_{i+1}^{\pi}(\mathcal{D}_N))]) \right], \\ &= 2 \mathcal{R}_{\max} \cdot \sup_{\pi} \mathbb{E}_{i \sim \mathcal{G}(\gamma)} [\text{TV} (P_{i+1}^*, \pi \| P_{i+1}^{\pi}(\mathcal{D}_N))]. \end{aligned}$$

□

We remark that Lemma 1 holds for any general reward-transition model given bounded rewards. The bound in Ineq. 8 proves the true regret is governed by the geometric average of TV distances: $\mathbb{E}_{i \sim \mathcal{G}(\gamma)} [\text{TV} (P_{i+1}^*, \pi \| P_{i+1}^{\pi}(\mathcal{D}_N))]$. As each term $\text{TV} (P_{i+1}^*, \pi \| P_{i+1}^{\pi}(\mathcal{D}_N))$ measures the distance between the true and predictive distributions over history h_i of length i , the discounting factor γ determines how much long term histories contribute to regret.

Intuitively, the more mass the posterior places close to the true value $\theta^* \in \Theta^*$, the smaller each TV distance becomes, with regret tending to zero for $P_{i+1, \pi}(\mathcal{D}_N) \approx P_{i+1, \pi}^* \implies \text{TV} (P_{i+1, \pi}^* \| P_{i+1, \pi}^{\pi}(\mathcal{D}_N)) \approx 0$. Conversely, a strong but highly incorrect prior will concentrate mass around MDPs whose dynamics oppose the true dynamics, yielding $\text{TV} (P_{i+1, \pi}^* \| P_{i+1, \pi}^{\pi}(\mathcal{D}_N)) \approx 1$ for all i , achieving the highest possible regret: $\mathcal{R}_{\max} := (r_{\max} - r_{\min}) / (1 - \gamma)$. The resulting Bayes-optimal policy would choose actions that encourage negative reward-seeking behaviour, being farthest from optimal in terms of expected returns.

Our next lemma

Lemma 2. *For bounded reward functions:*

$$\left| \mathbb{E}_{h_{i+1} \sim P_{i+1}^*, \pi} [r_i] - \mathbb{E}_{h_{i+1} \sim P_{i+1}^{\pi}(\mathcal{D}_N)} [r_i] \right| \leq (r_{\max} - r_{\min}) \cdot \text{TV} (P_{i+1}^*, \pi \| P_{i+1}^{\pi}(\mathcal{D}_N)). \quad (10)$$

Proof. We start by subtracting and adding $\frac{r_{\max} + r_{\min}}{2}$ to the left hand side of Ineq. 10:

$$\begin{aligned}
& \left| \mathbb{E}_{h_{i+1} \sim P_{i+1, \pi}^*} [r_i] - \mathbb{E}_{h_{i+1} \sim P_{i+1}^\pi(\mathcal{D}_N)} [r_i] \right| \\
&= \left| \mathbb{E}_{h_{i+1} \sim P_{i+1, \pi}^*} \left[r_i - \frac{r_{\max} + r_{\min}}{2} \right] - \mathbb{E}_{h_{i+1} \sim P_{i+1}^\pi(\mathcal{D}_N)} \left[r_i - \frac{r_{\max} + r_{\min}}{2} \right] \right|, \\
&= \left| \mathbb{E}_{h_{i+1} \sim P_{i+1, \pi}^*} \left[\frac{2r_i - (r_{\max} + r_{\min})}{2} \right] - \mathbb{E}_{h_{i+1} \sim P_{i+1}^\pi(\mathcal{D}_N)} \left[\frac{2r_i - (r_{\max} + r_{\min})}{2} \right] \right|, \\
&= \frac{(r_{\max} - r_{\min})}{2} \cdot \left| \mathbb{E}_{h_{i+1} \sim P_{i+1, \pi}^*} \left[\frac{2r_i - (r_{\max} + r_{\min})}{r_{\max} - r_{\min}} \right] - \mathbb{E}_{h_{i+1} \sim P_{i+1}^\pi(\mathcal{D}_N)} \left[\frac{2r_i - (r_{\max} + r_{\min})}{r_{\max} - r_{\min}} \right] \right|, \\
&= \frac{(r_{\max} - r_{\min})}{2} \cdot \left| \mathbb{E}_{h_{i+1} \sim P_{i+1, \pi}^*} [r_{\text{norm}}(h_{i+1})] - \mathbb{E}_{h_{i+1} \sim P_{i+1}^\pi(\mathcal{D}_N)} [r_{\text{norm}}(h_{i+1})] \right|, \tag{11}
\end{aligned}$$

where:

$$r_{\text{norm}}(h_{i+1}) := \frac{2r_i - (r_{\max} + r_{\min})}{r_{\max} - r_{\min}}.$$

Now, as $r_{\text{norm}} : \mathcal{H}_{i+1} \rightarrow [-1, 1]$, we can bound Eq. (11) using the integral probability metric form of the TV distance (see Eq. (7)), yielding our desired result:

$$\begin{aligned}
& \left| \mathbb{E}_{h_{i+1} \sim P_{i+1, \pi}^*} [r_i] - \mathbb{E}_{h_{i+1} \sim P_{i+1}^\pi(\mathcal{D}_N)} [r_i] \right| \\
&\leq \frac{r_{\max} - r_{\min}}{2} \cdot \sup_{f \in \mathcal{F}: \mathcal{H}_{i+1} \rightarrow [-1, 1]} \left| \mathbb{E}_{h_{i+1} \sim P_{i+1, \pi}^*} [f(h_{i+1})] - \mathbb{E}_{h_{i+1} \sim P_{i+1}^\pi(\mathcal{D}_N)} [f(h_{i+1})] \right|, \\
&= (r_{\max} - r_{\min}) \cdot \text{TV} (P_{i+1, \pi}^* \| P_{i+1}^\pi(\mathcal{D}_N)).
\end{aligned}$$

□

We proved in Lemma 1 that the rate of convergence of the sum of discounted TV distances between the true and predictive history distributions governs the rate of decrease in regret decreases with increasing data. Using the Bretagnolle-Huber inequality (see Section D.1), we now relate the sum of discounted TV distances to a sum of KL divergences, allowing us to control the expected regret using the PIL.

Theorem 1. *Using the PIL \mathcal{I}_N^π , the true regret is bounded as:*

$$\text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \leq 2\mathcal{R}_{\max} \cdot \sup_{\pi} \sqrt{1 - \exp\left(-\frac{\mathcal{I}_N^\pi}{1 - \gamma}\right)}$$

Proof. Starting with the bounded derived in Ineq. 8 of Lemma 1, we apply the Bretagnolle-Huber inequality (Bretagnolle and Huber, 1978) (see Section D.1) to bound the TV distance terms using the KL divergence:

$$\begin{aligned}
\text{Regret}(\mathcal{M}^*, \mathcal{D}_N) &\leq 2\mathcal{R}_{\max} \cdot \sup_{\pi} \mathbb{E}_{i \sim \mathcal{G}(\gamma)} [\text{TV} (P_{i+1, \pi}^* \| P_{i+1}^\pi(\mathcal{D}_N))], \\
&\leq 2\mathcal{R}_{\max} \cdot \sup_{\pi} \mathbb{E}_{i \sim \mathcal{G}(\gamma)} \left[\sqrt{1 - \exp(-\text{KL} (P_{i+1, \pi}^* \| P_{i+1}^\pi(\mathcal{D}_N)))} \right]. \tag{12}
\end{aligned}$$

We make two observations. Firstly, as the KL divergence is convex in its second argument and $P_{i+1, \pi}(\mathcal{D}_N) = \mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [P_{i+1}^\pi(\theta)]$, we can bound each KL divergence term using Jensen's inequality as:

$$\text{KL} (P_{i+1, \pi}^* \| P_{i+1}^\pi(\mathcal{D}_N)) \leq \mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [\text{KL} (P_{i+1, \pi}^* \| P_{i+1}^\pi(\theta))].$$

Secondly, as the function $f(x) = \sqrt{1 - \exp(-x)}$ is monotonically increasing in x , it follows that $f(x) \leq f(x')$ for any $x \leq x'$, hence:

$$\sqrt{1 - \exp(-\text{KL} (P_{i+1, \pi}^* \| P_{i+1}^\pi(\mathcal{D}_N)))} \leq \sqrt{1 - \exp(-\mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [\text{KL} (P_{i+1, \pi}^* \| P_{i+1}^\pi(\theta))])}$$

1134 Applying this bound to Ineq. 12 yields:
 1135
 1136

$$1137 \mathbb{E}_{i \sim \mathcal{G}(\gamma)} \left[\sqrt{1 - \exp(-\text{KL}(P_{i+1,\pi}^* \| P_{i+1}^\pi(\mathcal{D}_N)))} \right]$$

$$1138 \leq \mathbb{E}_{i \sim \mathcal{G}(\gamma)} \left[\sqrt{1 - \exp(-\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [\text{KL}(P_{i+1,\pi}^* \| P_{i+1}^\pi(\theta))])} \right].$$

1142
 1143 As the function $f(x) = \sqrt{1 - \exp(-x)}$ is concave in x , we can apply Jensen's inequality, yielding:
 1144
 1145

$$1146 \mathbb{E}_{i \sim \mathcal{G}(\gamma)} \left[\sqrt{1 - \exp(-\text{KL}(P_{i+1,\pi}^* \| P_{i+1}^\pi(\mathcal{D}_N)))} \right]$$

$$1147 \leq \sqrt{1 - \exp(-\mathbb{E}_{i \sim \mathcal{G}(\gamma)} [\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [\text{KL}(P_{i+1,\pi}^* \| P_{i+1}^\pi(\theta))])]}.$$
 (13)

1151
 1152 Examining the KL divergence term:
 1153
 1154

$$1155 \text{KL}(P_{i+1,\pi}^* \| P_{i+1}^\pi(\theta))$$

$$1156 = \mathbb{E}_{h_{i+1} \sim P_{i+1,\pi}^*} \left[\log \left(\frac{d_0(s_0) \prod_{j=0}^i \pi(a_j | h_j) p^*(r_j | s_j, a_j) p^*(s_{j+1} | s_j, a_j)}{d_0(s_0) \prod_{j=0}^i \pi(a_j | h_j) p(r_j | s_j, a_j, \theta) p(s_{j+1} | s_j, a_j, \theta)} \right) \right],$$

$$1157 = \mathbb{E}_{h_{i+1} \sim P_{i+1,\pi}^*} \left[\log \left(\frac{\prod_{j=0}^i p^*(r_j | s_j, a_j) p^*(s_{j+1} | s_j, a_j)}{\prod_{j=0}^i p(r_j | s_j, a_j, \theta) p(s_{j+1} | s_j, a_j, \theta)} \right) \right],$$

$$1158 = \mathbb{E}_{h_{i+1} \sim P_{i+1,\pi}^*} \left[\sum_{j=0}^i \left(\log p^*(r_j | s_j, a_j) - \log p(r_j | s_j, a_j, \theta) \right) \right.$$

$$1159 \left. + \log p^*(s_{j+1} | s_j, a_j) - \log p(s_{j+1} | s_j, a_j, \theta) \right],$$

$$1160 = \sum_{j=0}^i \mathbb{E}_{h_j \sim P_{j,\pi}^*} \left[\left(\log p^*(r_j | s_j, a_j) - \log p(r_j | s_j, a_j, \theta) \right) \right.$$

$$1161 \left. + \log p^*(s_{j+1} | s_j, a_j) - \log p(s_{j+1} | s_j, a_j, \theta) \right],$$

$$1162 = \sum_{j=0}^i \mathbb{E}_{s_j, a_j \sim P_{j,\pi}^*} \left[\mathbb{E}_{r_j, s_{j+1} \sim P_{R,S}^*(s_j, a_j)} \left[\left(\log p^*(r_j | s_j, a_j) - \log p(r_j | s_j, a_j, \theta) \right) \right. \right.$$

$$1163 \left. + \log p^*(s_{j+1} | s_j, a_j) - \log p(s_{j+1} | s_j, a_j, \theta) \right] \right],$$

$$1164 = \sum_{j=0}^i \mathbb{E}_{s_j, a_j \sim P_{j,\pi}^*} \left[\mathbb{E}_{r_j, s_{j+1} \sim P_{R,S}^*(s_j, a_j)} \left[\left(\log p^*(r_j, s_{j+1} | s_j, a_j) \right) \right. \right.$$

$$1165 \left. - \log p(r_j, s_{j+1} | s_j, a_j, \theta) \right] \right],$$

$$1166 = \sum_{j=0}^i \mathbb{E}_{s, a \sim P_{j,\pi}^*} [\text{KL}(P_{R,S}^*(s, a) \| P_{R,S}(s, a, \theta))],$$

1188 hence:

$$\begin{aligned}
1189 & \mathbb{E}_{i \sim \mathcal{G}(\gamma)} \left[\mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} \left[\text{KL} \left(P_{i+1, \pi}^* \| P_{i+1}^{\pi}(\theta) \right) \right] \right] \\
1190 & = \mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} \left[\mathbb{E}_{i \sim \mathcal{G}(\gamma)} \left[\sum_{j=0}^i \mathbb{E}_{s, a \sim P_{j, \pi}^*} \left[\text{KL} \left(P_{R, S}^*(s, a) \| P_{R, S}(s, a, \theta) \right) \right] \right] \right], \\
1191 & = \mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} \left[\sum_{i=0}^{\infty} (1 - \gamma) \gamma^i \sum_{j=0}^i \mathbb{E}_{s, a \sim P_{j, \pi}^*} \left[\text{KL} \left(P_{R, S}^*(s, a) \| P_{R, S}(s, a, \theta) \right) \right] \right], \\
1192 & = \mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} \left[\sum_{i=0}^{\infty} (1 - \gamma) \gamma^i (i + 1) \sum_{j=0}^i \frac{1}{i + 1} \mathbb{E}_{s, a \sim P_{j, \pi}^*} \left[\text{KL} \left(P_{R, S}^*(s, a) \| P_{R, S}(s, a, \theta) \right) \right] \right], \\
1193 & = \frac{1}{1 - \gamma} \mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} \left[\sum_{i=0}^{\infty} (1 - \gamma)^2 \gamma^i (i + 1) \mathbb{E}_{j \sim \mathcal{U}_i} \left[\mathbb{E}_{s, a \sim P_{j, \pi}^*} \left[\text{KL} \left(P_{R, S}^*(s, a) \| P_{R, S}(s, a, \theta) \right) \right] \right] \right], \\
1194 & = \frac{1}{1 - \gamma} \mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} \left[\mathbb{E}_{i \sim \text{AG}(\gamma)} \left[\mathbb{E}_{j \sim \mathcal{U}_i} \left[\mathbb{E}_{s, a \sim P_{j, \pi}^*} \left[\text{KL} \left(P_{R, S}^*(s, a) \| P_{R, S}(s, a, \theta) \right) \right] \right] \right] \right]. \quad (14)
\end{aligned}$$

1206 Now, as $\rho_{\pi}^* = \mathbb{E}_{i \sim \text{AG}(\gamma)} \left[\mathbb{E}_{j \sim \mathcal{U}_i} \left[P_{j, \pi}^* \right] \right]$ is the arithemico-geometric ergodic state-action distribu-
1207 tion, we can simplify Eq. (14) to yield:

$$\begin{aligned}
1208 & \mathbb{E}_{i \sim \mathcal{G}(\gamma)} \left[\mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} \left[\text{KL} \left(P_{i+1, \pi}^* \| P_{i+1}^{\pi}(\theta) \right) \right] \right] \\
1209 & = \frac{1}{1 - \gamma} \mathbb{E}_{\mathcal{D}_N \sim P_{\text{Data}}} \left[\mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} \left[\mathbb{E}_{s, a \sim \rho_{\pi}^*} \left[\text{KL} \left(P_{R, S}^*(s, a) \| P_{R, S}(s, a, \theta) \right) \right] \right] \right], \\
1210 & = \frac{1}{1 - \gamma} \mathbb{E}_{\mathcal{D}_N \sim P_{\text{Data}}} \left[\mathbb{E}_{s, a \sim \rho_{\pi}^*} \left[\mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} \left[\text{KL} \left(P_{R, S}^*(s, a) \| P_{R, S}(s, a, \theta) \right) \right] \right] \right], \\
1211 & = \frac{1}{1 - \gamma} \mathcal{I}_N^{\pi},
\end{aligned}$$

1212 hence, substituting into Ineq. 13, we obtain:

$$\mathbb{E}_{i \sim \mathcal{G}(\gamma)} \left[\sqrt{1 - \exp \left(-\text{KL} \left(P_{i+1, \pi}^* \| P_{i+1}^{\pi}(\mathcal{D}_N) \right) \right)} \right] \leq \sqrt{1 - \exp \left(-\frac{\mathcal{I}_N^{\pi}}{1 - \gamma} \right)}.$$

1213 Finally, substituting into Eq. (12) yields our desired result:

$$\begin{aligned}
1214 & \text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \leq 2\mathcal{R}_{\max} \cdot \sup_{\pi} \left| \mathbb{E}_{i \sim \mathcal{G}(\gamma)} \left[\sqrt{1 - \exp \left(-\text{KL} \left(P_{i+1, \pi}^* \| P_{i+1}^{\pi}(\mathcal{D}_N) \right) \right)} \right] \right|, \\
1215 & \leq 2\mathcal{R}_{\max} \cdot \sup_{\pi} \sqrt{1 - \exp \left(-\frac{\mathcal{I}_N^{\pi}}{1 - \gamma} \right)}.
\end{aligned}$$

1216 \square

1217 The PIL has an intuitive information-geometric interpretation: the inner expectation
1218 $\mathbb{E}_{s, a \sim \rho_{\pi}^*} \left[\text{KL} \left(P_{R, S}^*(s, a) \| P_{R, S}(s, a, \theta) \right) \right]$ measures the distance between the model and the true
1219 distribution in terms of the information lost when approximating $P_{R, S}^*(s, a)$ with $P_{R, S}(s, a, \theta)$,
1220 averaged across all states. The PIL thus measures how close the posterior's belief is to the truth
1221 according to the average information lost under the posterior expectation. We observe that via
1222 Jensen's inequality, the PIL is an upper bound on the classic KL risk (sometimes known as expected
1223 relative entropy) from Bayesian asymptotics and regret analysis (Aitchison, 1975; Clarke and Barron,
1224 1990; Komaki, 1996; Hartigan, 1998; Barron, 1988; 1999; Yang and Barron, 1999; van der Vaart and
1225 van Zanten, 2011; Aslan, 2006; Alaa and van der Schaar, 2018; Bilodeau et al., 2021).

1226 By substituting in our definition of the Gaussian world model, we now find a convenient form for the
1227 PIL:

1228 **Proposition 2.** *Using the Gaussian world model in Eq. (2), it follows:*

$$\mathcal{I}_N^{\pi} = \mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) + \mathcal{V}(\mathcal{D}_N).$$

1242 *Proof.* We substitute the Gaussian world model into the KL divergence to yield:

$$\begin{aligned}
& \text{KL} (P_{R,S}^*(s, a) \| P_{R,S}(s, a, \theta)) \\
&= \mathbb{E}_{r, s' \sim P_{R,S}^*(s, a)} \left[\log \left(\exp \left(-\frac{\|r^*(s, a) - r\|_2^2}{2\sigma_r^2} \right) \exp \left(-\frac{\|s^{*'}(s, a) - s'\|_2^2}{2\sigma_s^2} \right) \right) \right] \\
&\quad - \mathbb{E}_{r, s' \sim P_{R,S}(s, a)} \left[\log \left(\exp \left(-\frac{\|r_\theta(s, a) - r\|_2^2}{2\sigma_r^2} \right) \exp \left(-\frac{\|s'_\theta(s, a) - s'\|_2^2}{2\sigma_s^2} \right) \right) \right], \\
&= \mathbb{E}_{r, s' \sim P_{R,S}^*(s, a)} \left[\frac{\|r_\theta(s, a) - r\|_2^2 - \|r^*(s, a) - r\|_2^2}{2\sigma_r^2} \right. \\
&\quad \left. + \frac{\|s'_\theta(s, a) - s'\|_2^2 - \|s^{*'}(s, a) - s'\|_2^2}{2\sigma_s^2} \right], \\
&= \mathbb{E}_{r, s' \sim P_{R,S}^*(s, a)} \left[\frac{r_\theta(s, a)^2 - 2rr_\theta(s, a) - r^*(s, a)^2 + 2rr^*(s, a)}{2\sigma_r^2} \right. \\
&\quad \left. + \frac{\|s'_\theta(s, a)\|_2^2 - 2s'^\top s'_\theta(s, a) - \|s^{*'}(s, a)\|_2^2 + 2s'^\top s^{*'}(s, a)}{2\sigma_s^2} \right], \\
&= \frac{r_\theta(s, a)^2 - 2r^*(s, a)r_\theta(s, a) - r^*(s, a)^2 + 2r^*(s, a)^2}{2\sigma_r^2} \\
&\quad + \frac{\|s'_\theta(s, a)\|_2^2 - 2s^{*'\top} s'_\theta(s, a) - \|s^{*'}(s, a)\|_2^2 + 2\|s^{*'}(s, a)\|_2^2}{2\sigma_s^2}, \\
&= \frac{r_\theta(s, a)^2 - 2r^*(s, a)r_\theta(s, a) + r^*(s, a)^2}{2\sigma_r^2} \\
&\quad + \frac{\|s'_\theta(s, a)\|_2^2 - 2s^{*'\top} s'_\theta(s, a) + \|s^{*'}(s, a)\|_2^2}{2\sigma_s^2}.
\end{aligned}$$

1274
1275
1276 Now, taking expectations with respect to the posterior:

$$\begin{aligned}
& \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\text{KL} (P_{R,S}^*(s, a) \| P_{R,S}(s, a, \theta)) \right] \\
&= \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\frac{r_\theta(s, a)^2 - 2r^*(s, a)r_\theta(s, a) + r^*(s, a)^2}{2\sigma_r^2} \right] \\
&\quad + \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\frac{\|s'_\theta(s, a)\|_2^2 - 2s^{*'\top} s'_\theta(s, a) + \|s^{*'}(s, a)\|_2^2}{2\sigma_s^2} \right], \\
&= \frac{\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [r_\theta(s, a)^2] - 2r^*(s, a)r(s, a, \mathcal{D}_N) + r^*(s, a)^2}{2\sigma_r^2} \\
&\quad + \frac{\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [\|s'_\theta(s, a)\|_2^2] - 2s^{*'\top} s'(s, a, \mathcal{D}_N) + \|s^{*'}(s, a)\|_2^2}{2\sigma_s^2}.
\end{aligned}$$

1292
1293
1294 Now, we use the variance identity for both the reward and state functions: $\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [r_\theta(s, a)^2] =$
1295 $\mathbb{V}_{\theta \sim P_\Theta(\mathcal{D}_N)} [r_\theta(s, a)] + r_\theta(s, a, \mathcal{D}_N)^2$ and $\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [\|s'_\theta(s, a)\|_2^2] = \mathbb{V}_{\theta \sim P_\Theta(\mathcal{D}_N)} [\|s'_\theta(s, a)\|_2] +$

1296 $\|s'(s, a, \mathcal{D}_N)\|_2^2$ yielding:

$$\begin{aligned}
1297 & \mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [\text{KL}(P_{R,S}^*(s, a) \| P_{R,S}(s, a, \theta))] \\
1298 &= \frac{\mathbb{V}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [r_{\theta}(s, a)] + r_{\theta}(s, a, \mathcal{D}_N)^2 - 2r^*(s, a)r(s, a, \mathcal{D}_N) + r^*(s, a)^2}{2\sigma_r^2} \\
1299 &+ \frac{\mathbb{V}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [\|s'_{\theta}(s, a)\|_2] + \|s'(s, a, \mathcal{D}_N)\|_2^2 - 2s^{*\prime}(s, a)^{\top} s'(s, a, \mathcal{D}_N) + \|s^{*\prime}(s, a)\|_2^2}{2\sigma_s^2}, \\
1300 &= \frac{\mathbb{V}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [r_{\theta}(s, a)] + (r_{\theta}(s, a, \mathcal{D}_N)^2 - r^*(s, a)^2)}{2\sigma_r^2} \\
1301 &+ \frac{\mathbb{V}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [\|s'_{\theta}(s, a)\|_2] + \|s'(s, a, \mathcal{D}_N) - s^{*\prime}(s, a)\|_2^2}{2\sigma_s^2}, \\
1302 &= \frac{\mathbb{V}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [r_{\theta}(s, a)]}{2\sigma_r^2} + \frac{\mathbb{V}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [\|s'_{\theta}(s, a)\|_2]}{2\sigma_s^2} \\
1303 &+ \frac{(r_{\theta}(s, a, \mathcal{D}_N) - r^*(s, a))^2}{2\sigma_r^2} + \frac{\|s'(s, a, \mathcal{D}_N) - s^{*\prime}(s, a)\|_2^2}{2\sigma_s^2}, \\
1304 &= \mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) + \mathcal{V}(\mathcal{D}_N),
\end{aligned}$$

1305 and hence:

$$\begin{aligned}
1306 & \mathcal{I}_N^{\pi} = \mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [\text{KL}(P_{R,S}^*(s, a) \| P_{R,S}(s, a, \theta))], \\
1307 &= \mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) + \mathcal{V}(\mathcal{D}_N).
\end{aligned}$$

1308 \square

1309 D.3 PROOF OF THEOREM 2

1310 We first introduce some simplifying notation for the expected cross entropy, log likelihood and corresponding gradients and Hessian:

$$\begin{aligned}
1311 & \ell(\theta) := \mathbb{E}_{s, a \sim \rho_{\pi}^*, r, s' \sim P_{R,S}^*(s, a)} [\log p(r, s' | s, a, \theta)], \\
1312 & \ell^* := \max_{\theta \in \Theta} \ell(\theta) = \mathbb{E}_{s, a \sim \rho_{\pi}^*, r, s' \sim P_{R,S}^*(s, a)} [\log p^*(r, s' | s, a)], \\
1313 & \ell_N(\theta) := \frac{1}{N} \sum_{i=0}^{N-1} \log p(r_i, s'_i | s_i, a_i, \theta), \\
1314 & g_{i,N}^* := \sqrt{N} \nabla_{\theta} \ell_N(\theta) \Big|_{\theta=\theta_i^*}, \\
1315 & H_i^* := \nabla_{\theta}^2 \ell(\theta) \Big|_{\theta=\theta_i^*}.
\end{aligned}$$

1316 We now introduce key regularity assumptions for our parametric model that are required to derive the convergence rate for PIL. They are relatively mild and commonplace in the asymptotic statistics literature (Le Cam, 1953; Barron, 1988; Clarke and Barron, 1990; Komaki, 1996; Hartigan, 1998; Barron, 1999; Aslan, 2006).

1317 **Assumption 1.** We assume that:

1318 *i* There exists at least one parametrisation that corresponds to the true environment dynamics with:

$$1319 \left| \mathbb{E}_{s, a \sim \rho_{\pi}^*, r, s' \sim P_{R,S}^*(s, a)} [\log p^*(r, s' | s, a)] \right| < \infty$$

1320 and $|\ell^* - \ell(\theta)|$ is bounded P_{Θ} -almost surely.

1321 *ii* $\ell_N(\theta)$ and $\ell(\theta)$ are C^2 -continuous in θ .

1322 *iii* There are $K < \infty$ maximising points θ_i^* :

$$1323 \{\theta_1^*, \theta_2^*, \dots, \theta_K^*\} = \arg \max_{\theta \in \Theta} \ell(\theta).$$

1324 For each maximiser θ_i^* , there exists a small region $\Theta_i^* := \{\theta \in \Theta \mid \|\theta_i^* - \theta\| \leq \epsilon\}$ for some $\epsilon > 0$ such that θ_i^* is the unique maximiser in Θ_i^* , θ_i^* is in the interior of Θ_i^* , $\nabla_{\theta}^2 \ell(\theta_i^*)$ is negative definite, invertible and the regions are disjoint: $\bigcap_{i=1}^K \Theta_i^* = \emptyset$.

1350 iv The prior $p(\theta)$ is Lipschitz continuous in θ with support over Θ .

1351 v The sampling regime ensures that the strong law of large numbers holds for all maximisers θ_i^* for
1352 the Hessian, and uniformly for $\theta \in \Theta$ for the likelihood, that is:

$$1353 \ell_N(\theta) \xrightarrow{\text{Unif. a.s.}} \ell(\theta), \quad \nabla_{\theta}^2 \ell_N(\theta_i^*) \xrightarrow{\text{a.s.}} \nabla_{\theta}^2 \ell(\theta_i^*).$$

1354 The central limit theorem applies to the gradient at each θ_i^* , that is:

$$1355 \sqrt{N} \nabla_{\theta} \ell_N(\theta_i^*) \xrightarrow{d} \mathcal{N}(0, \Sigma_i^g),$$

1356 where $\Sigma_i^g = \mathbb{E}_{s, a \sim \rho_{\pi}^*, r, s' \sim P_{R, S}(s, a)} [\nabla_{\theta} \log p(r, s' | s, a, \theta_i^*) \nabla_{\theta} \log p(r, s' | s, a, \theta_i^*)^{\top}]$ with
1357 $\|\Sigma_i^g\| < \infty$.

1358 Our assumptions are mild. Assumption 1i is our strictest assumption and is included for ease of
1359 exposition. We generalise our theory in Section D.4 to relax this assumption and also account
1360 for incomplete Bayes-optimal policy learning. Assumption 1ii ensures that a second order Taylor
1361 series expansion can be applied to obtain an asymptotic expansion around the maximising points.
1362 Assumption 1iii is much more general than most settings, which only consider problems with a single
1363 maximiser. The invertibility of the matrix can easily be guaranteed in Bayesian methods by the use
1364 of a prior that can re-condition a low rank matrix that may results from linearly dependent data.
1365 Assumption 1iv ensures that the prior places sufficient mass on the true parametrisation. The sampling
1366 and model would need to be very irregular for Assumption 1v not to hold; stochastic optimisation
1367 methods used to find statistics like the MAP will fail if this assumption did not hold. Assumption 1v
1368 holds automatically if sampling is either i.i.d. from $s, a \sim \rho_{\pi}^*$ (see e.g. Bass (2013)) or from an
1369 aperiodic and irreducible Markov chain with stationary distribution ρ_{π}^* (see e.g. Roberts and Rosenthal
1370 (2004)). In both sample regimes, noting that $\mathbb{E}_{s, a \sim \rho_{\pi}^*, r, s' \sim P_{R, S}(s, a)} [\nabla_{\theta} \log p(r, s' | s, a, \theta_i^*)] = 0$, it's
1371 clear the (long run) covariance of $\nabla_{\theta} \log p(r, s' | s, a, \theta_i^*)$ is Σ_i^g .

1372 Our first lemma borrows techniques from Vaart (1998, Chapter 10). This approach is similar to
1373 asymptotic integral expansion approaches that apply Laplace's method (Lindley, 1980; Tierney and
1374 Kadane, 1986; Tierney et al., 1989; Kass et al., 1990) except we expand around the global maximising
1375 values of $\ell(\theta)$ rather than the maximising values of the likelihood $\ell_N(\theta)$ to obtain an asymptotic
1376 expression for the posterior:

1377 **Lemma 3.** Under Assumption 1 and using the notation introduced at the start of Section D.3:

$$1378 \frac{\int_{\Theta_i^*} (\ell^* - \ell(\theta)) \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta_i^*} \exp(N\ell_N(\theta)) p(\theta) d\theta} = \mathcal{O}\left(\frac{d - g_{i, N}^{* \top} H_i^{* - 1} g_{i, N}^*}{N}\right),$$

1379 almost surely.

1380 *Proof.* We start by applying the transformation of variables $\theta' = f(\theta) := \sqrt{N}(\theta - \theta_i^*)$ to integrals
1381 in the numerator and denominator with:

$$1382 \theta = f^{-1}(\theta') = \theta_i^* + \frac{1}{\sqrt{N}}\theta', \quad |\det \nabla_{\theta} f^{-1}(\theta')| = N^{-\frac{d}{2}}, \quad \Theta' := f(\Theta_i^*),$$

1383 yielding:

$$1384 \frac{\int_{\Theta_i^*} (\ell^* - \ell(\theta)) \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta_i^*} \exp(N\ell_N(\theta)) p(\theta) d\theta} = \frac{\int_{\Theta'} \left(\ell^* - \ell\left(\theta = \theta_i^* + \frac{1}{\sqrt{N}}\theta'\right)\right) \exp\left(N\ell_N\left(\theta = \theta_i^* + \frac{1}{\sqrt{N}}\theta'\right)\right) p'(\theta') d\theta'}{\int_{\Theta'} \exp\left(N\ell_N\left(\theta = \theta_i^* + \frac{1}{\sqrt{N}}\theta'\right)\right) p'(\theta') d\theta'}, \quad (15)$$

1385 where $p'(\theta') := p\left(\theta = \theta_i^* + \frac{1}{\sqrt{N}}\theta'\right)$. Now and making a Taylor series expansion of $\ell(\theta)$ about θ_i^* :

$$1386 \ell(\theta) = \ell^* + \underbrace{\nabla_{\theta} \ell(\theta_i^*)}_{=0}^{\top} (\theta - \theta_i^*) + (\theta - \theta_i^*)^{\top} H_i^* (\theta - \theta_i^*) + \mathcal{O}(\|\theta - \theta_i^*\|^3),$$

$$1387 = \ell^* + (\theta - \theta_i^*)^{\top} H_i^* (\theta - \theta_i^*) + \mathcal{O}(\|\theta - \theta_i^*\|^3),$$

1404 hence:

$$1405 \ell \left(\theta = \theta_i^* + \frac{1}{\sqrt{N}} \theta' \right) = \ell^* + \frac{1}{N} \theta'^{\top} H_i^* \theta' + \mathcal{O} \left(N^{-\frac{3}{2}} \right).$$

1408 Using the notation $H_N^* := \nabla_{\theta}^2 \ell_N(\theta)|_{\theta=\theta_i^*}$ and making a Taylor series expansion of $\ell_N(\theta)$ about θ_i^* :

$$1409 \ell_N(\theta) = \ell_N(\theta_i^*) + \nabla_{\theta} \ell_N(\theta_i^*)^{\top} (\theta - \theta_i^*) + (\theta - \theta_i^*)^{\top} \nabla_{\theta}^2 \ell_N(\theta_i^*) (\theta - \theta_i^*) + \mathcal{O}(\|\theta - \theta_i^*\|^3),$$

1412 hence:

$$1413 N \ell_N \left(\theta = \theta_i^* + \frac{1}{\sqrt{N}} \theta' \right) = N \ell_N(\theta_i^*) + \sqrt{N} \nabla_{\theta} \ell_N(\theta_i^*)^{\top} \theta' + \theta'^{\top} \nabla_{\theta}^2 \ell_N(\theta_i^*) \theta' + \mathcal{O} \left(\frac{1}{\sqrt{N}} \right).$$

1415 Substituting into Eq. (15) yields:

$$\begin{aligned} 1417 & \frac{\int_{\Theta^*} (\ell^* - \ell(\theta)) \exp(N\ell(\theta)) p(\theta) d\theta}{\int_{\Theta^*} \exp(N\ell(\theta)) p(\theta) d\theta} \\ 1418 &= - \frac{\int_{\Theta'} \theta'^{\top} H_i^* \theta' \exp \left(N \ell_N(\theta_i^*) + g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' + \mathcal{O} \left(\frac{1}{\sqrt{N}} \right) \right) p'(\theta') d\theta'}{\int_{\Theta'} \exp \left(N \ell_N(\theta_i^*) + g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' + \mathcal{O} \left(\frac{1}{\sqrt{N}} \right) \right) p'(\theta') d\theta'} \mathcal{O} \left(\frac{1}{N} \right), \\ 1420 &= - \frac{\int_{\Theta'} \theta'^{\top} H_i^* \theta' \exp \left(g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' + \mathcal{O} \left(\frac{1}{\sqrt{N}} \right) \right) p'(\theta') d\theta'}{\int_{\Theta'} \exp \left(N g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' + \mathcal{O} \left(\frac{1}{\sqrt{N}} \right) \right) p'(\theta') d\theta'} \mathcal{O} \left(\frac{1}{N} \right), \\ 1422 &= - \frac{\int_{\Theta'} \theta'^{\top} H_i^* \theta' \exp \left(g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) \exp \left(\mathcal{O} \left(\frac{1}{\sqrt{N}} \right) \right) p'(\theta') d\theta'}{\int_{\Theta'} \exp \left(N g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) \exp \left(\mathcal{O} \left(\frac{1}{\sqrt{N}} \right) \right) p'(\theta') d\theta'} \mathcal{O} \left(\frac{1}{N} \right), \\ 1424 &= \mathcal{O} \left(- \frac{1}{N} \frac{\int_{\Theta'} \theta'^{\top} H_i^* \theta' \exp \left(g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) p'(\theta') d\theta'}{\int_{\Theta'} \exp \left(N g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) p'(\theta') d\theta'} \right) \end{aligned} \quad (16)$$

1434 where we have multiplied top and bottom by $\exp(-N\ell_N(\theta_i^*))$ to derive the second equality and
1435 used the fact that $0 < \exp \left(\mathcal{O} \left(\frac{1}{\sqrt{N}} \right) \right) < \infty$ to derive the final line. Now, as the prior is Lipschitz,
1436 we make a Taylor series expansion about θ_i^* :

$$1437 p(\theta) = p(\theta_i^*) + \mathcal{O}(\|\theta - \theta_i^*\|),$$

1439 hence:

$$1440 p'(\theta') = p \left(\theta = \theta_i^* + \frac{1}{\sqrt{N}} \theta' \right) = p(\theta_i^*) + \mathcal{O} \left(\frac{1}{\sqrt{N}} \right).$$

1443 This allows us to find an asymptotic expression for Eq. (16):

$$\begin{aligned} 1444 & \frac{\int_{\Theta'} \theta'^{\top} H_i^* \theta' \exp \left(g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) p'(\theta') d\theta'}{\int_{\Theta'} \exp \left(g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) p'(\theta') d\theta'} \\ 1445 &= \frac{\int_{\Theta'} \theta'^{\top} H_i^* \theta' \exp \left(g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) p'(\theta_i^*) d\theta' \left(1 + \mathcal{O} \left(\frac{1}{\sqrt{N}} \right) \right)}{\int_{\Theta'} \exp \left(g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) p'(\theta_i^*) d\theta' \left(1 + \mathcal{O} \left(\frac{1}{\sqrt{N}} \right) \right)} \\ 1448 &= \frac{\int_{\Theta'} \theta'^{\top} H_i^* \theta' \exp \left(g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) p'(\theta_i^*) d\theta'}{\int_{\Theta'} \exp \left(g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) p'(\theta_i^*) d\theta'} \left(1 + \mathcal{O} \left(\frac{1}{\sqrt{N}} \right) \right) \\ 1449 &= \frac{\int_{\Theta'} \theta'^{\top} H_i^* \theta' \exp \left(g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) d\theta'}{\int_{\Theta'} \exp \left(g_{i,N}^{*\top} \theta' + \theta'^{\top} H_N^* \theta' \right) d\theta'} \left(1 + \mathcal{O} \left(\frac{1}{\sqrt{N}} \right) \right). \end{aligned}$$

We re-write the exponential term to recover a quadratic form:

$$\begin{aligned} & \exp\left(g_{i,N}^{\star\top} \theta' + \theta'^\top H_N^* \theta'\right) \\ &= \exp\left(\left(\frac{1}{2} H_N^{\star-1} g_{i,N}^* + \theta'\right)^\top H_N^* \left(\frac{1}{2} H_N^{\star-1} g_{i,N}^* + \theta'\right) - \frac{1}{4} g_{i,N}^{\star\top} H_N^{\star-1} g_{i,N}^*\right). \end{aligned}$$

Substituting yields:

$$\begin{aligned} & \frac{\int_{\Theta'} \theta'^\top H_i^* \theta' \exp\left(g_{i,N}^{\star\top} \theta' + \theta'^\top H_N^* \theta'\right) d\theta'}{\int_{\Theta'} \exp\left(g_{i,N}^{\star\top} \theta' + \theta'^\top H_N^* \theta'\right) d\theta'} \\ &= \frac{\int_{\Theta'} \theta'^\top H_i^* \theta' \exp\left(\left(\frac{1}{2} H_N^{\star-1} g_{i,N}^* + \theta'\right)^\top H_N^* \left(\frac{1}{2} H_N^{\star-1} g_{i,N}^* + \theta'\right)\right) d\theta'}{\int_{\Theta'} \exp\left(\left(\frac{1}{2} H_N^{\star-1} g_{i,N}^* + \theta'\right)^\top H_N^* \left(\frac{1}{2} H_N^{\star-1} g_{i,N}^* + \theta'\right)\right) d\theta'}. \end{aligned}$$

In this form, we notice the expectation is that of a Gaussian $\mathcal{N}(\mu = -\frac{1}{2} H_N^{\star-1} g_{i,N}^*, \Sigma = -H_i^{\star-1})$ restricted to Θ' . Noting that in the limit $\Theta' \xrightarrow{N \rightarrow \infty} \mathbb{R}^d$, hence:

$$\begin{aligned} & \frac{\int_{\Theta'} \theta'^\top H_i^* \theta' \exp\left(\left(\theta' + \frac{1}{2} H_N^{\star-1} g_{i,N}^*\right)^\top H_N^* \left(\theta' + \frac{1}{2} H_N^{\star-1} g_{i,N}^*\right)\right) d\theta'}{\int_{\Theta'} \exp\left(\left(\theta' + \frac{1}{2} H_N^{\star-1} g_{i,N}^*\right)^\top H_N^* \left(\theta' + \frac{1}{2} H_N^{\star-1} g_{i,N}^*\right)\right) d\theta'} \\ &= \mathcal{O}\left(\frac{\int_{\mathbb{R}^d} \theta'^\top H_i^* \theta' \exp\left(\left(\theta' + \frac{1}{2} H_N^{\star-1} g_{i,N}^*\right)^\top H_N^* \left(\theta' + \frac{1}{2} H_N^{\star-1} g_{i,N}^*\right)\right) d\theta'}{\int_{\mathbb{R}^d} \exp\left(\left(\theta' + \frac{1}{2} H_N^{\star-1} g_{i,N}^*\right)^\top H_N^* \left(\theta' + \frac{1}{2} H_N^{\star-1} g_{i,N}^*\right)\right) d\theta'}\right), \\ &= \mathcal{O}\left(\mathbb{E}_{\theta' \sim \mathcal{N}\left(-\frac{1}{2} H_N^{\star-1} g_{i,N}^*, -H_N^{\star-1}\right)} \left[\theta'^\top H_i^* \theta'\right]\right). \end{aligned} \quad (17)$$

Putting everything together, we have:

$$\begin{aligned} & \frac{\int_{\Theta_i^*} (\ell^* - \ell(\theta)) \exp(N \ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta_i^*} \exp(N \ell_N(\theta)) p(\theta) d\theta} \\ &= \mathcal{O}\left(-\frac{1}{N} \frac{\int_{\Theta'} \theta'^\top H_i^* \theta' \exp\left(g_{i,N}^{\star\top} \theta' + \theta'^\top H_N^* \theta'\right) p'(\theta') d\theta'}{\int_{\Theta'} \exp\left(g_{i,N}^{\star\top} \theta' + \theta'^\top H_N^* \theta'\right) p'(\theta') d\theta'}\right), \quad \text{Eq. (16)} \\ &= \mathcal{O}\left(-\frac{1}{N} \mathbb{E}_{\theta' \sim \mathcal{N}\left(-\frac{1}{2} H_N^{\star-1} g_{i,N}^*, -H_i^{\star-1}\right)} \left[\theta'^\top H_i^* \theta'\right]\right), \quad \text{Eq. (17)} \end{aligned}$$

Using standard results for the multivariate Gaussian (Petersen and Pedersen, 2012) yields our desired result:

$$\begin{aligned} -\frac{1}{N} \mathbb{E}_{\theta' \sim \mathcal{N}\left(-\frac{1}{2} H_N^{\star-1} g_{i,N}^*, -H_i^{\star-1}\right)} \left[\theta'^\top H_i^* \theta'\right] &= \frac{\text{Tr}\left(H_i^* H_N^{\star-1}\right) - \frac{1}{4} g_{i,N}^{\star\top} H_N^{\star-1\top} H_i^* H_N^{\star-1} g_{i,N}^*}{N}, \\ &= \mathcal{O}\left(\frac{\text{Tr}(I) - g_{i,N}^{\star\top} H_i^{\star-1} g_{i,N}^*}{N}\right), \\ &= \mathcal{O}\left(\frac{d - g_{i,N}^{\star\top} H_i^{\star-1} g_{i,N}^*}{N}\right), \end{aligned}$$

almost surely, where we have used the strong law of large numbers on the empirical Hessian from Assumption 1 to derive the second line. \square

In our final Lemma, we show that regions that are not close to the maximising points diminish exponentially in posterior probability as N grows large.

Lemma 4. Under Assumption 1, $\mathbb{E}_{\mathcal{D}_N \sim P_{\text{Data}}^*} [P(\bar{\Theta}|\mathcal{D}_N)] = \mathcal{O}(\exp(-N))$.

Proof. We start by splitting the posterior expectation into integrals over $\bar{\Theta}$ and $\Theta \setminus \bar{\Theta}$:

$$\begin{aligned} P(\bar{\Theta}|\mathcal{D}_N) &= \frac{\int_{\bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta} \exp(N\ell_N(\theta)) p(\theta) d\theta} \\ &= \frac{\int_{\bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta + \int_{\Theta \setminus \bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta}. \end{aligned}$$

Dividing top and bottom by $\int_{\bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta$:

$$P(\bar{\Theta}|\mathcal{D}_N) = \frac{1}{1 + \frac{\int_{\Theta \setminus \bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta}}.$$

Hence if we can show there exists some $N' < \infty$ and a function $C \exp(cN)$ with positive constants c and C that lower bounds the ratio:

$$C \exp(cN) \leq \frac{\int_{\Theta \setminus \bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta}$$

almost surely for all $N \geq N'$, then it follows:

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_N \sim P_{\text{Data}}^*} [P(\bar{\Theta}|\mathcal{D}_N)] &= \mathcal{O}\left(\frac{1}{1 + \exp(N)}\right), \\ &= \mathcal{O}(\exp(-N)). \end{aligned}$$

From Assumption 1, each θ_i^* maximises $\ell(\theta)$ with $\sup_{\theta \in \bar{\Theta}} \ell(\theta') < \ell(\theta_i^*)$. As $\ell(\theta)$ is continuous, there thus exists a small, closed ball $B(\theta_j^*, r) := \{\theta \mid \|\theta_j^* - \theta\| \leq r\}$ of radius $r > 0$ centred on some θ_j^* such that $\sup_{\theta' \in \bar{\Theta}} \ell(\theta') < \min_{\theta'' \in B(\theta_j^*, r)} \ell(\theta'')$. From Assumption 1, the uniform strong

law of large numbers holds with $\ell_N(\theta) \xrightarrow{\text{Unif. a.s.}} \ell(\theta)$. By the definition of the limit and continuity of $\ell_N(\theta)$, there thus exists some finite N' such that $\sup_{\theta' \in \bar{\Theta}} \ell_N(\theta') < \min_{\theta'' \in B(\theta_j^*, \frac{r}{2})} \ell_N(\theta'')$ for all $N \geq N'$ almost surely, where $B(\theta_j^*, \frac{r}{2})$ is a ball of half radius $\frac{r}{2}$. Noting that $B(\theta_j^*, \frac{r}{2}) \subset \Theta \setminus \bar{\Theta}$ and $0 \leq \exp(N\ell_N(\theta))$, this allows us to lower bound the integral:

$$\begin{aligned} \int_{\Theta \setminus \bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta &\geq \int_{B(\theta_j^*, \frac{r}{2})} \exp(N\ell_N(\theta)) p(\theta) d\theta, \\ &\geq \exp\left(N \min_{\theta'' \in B(\theta_j^*, \frac{r}{2})} \ell_N(\theta'')\right) \int_{B(\theta_j^*, \frac{r}{2})} p(\theta) d\theta, \\ &= \exp\left(N \min_{\theta'' \in B(\theta_j^*, \frac{r}{2})} \ell_N(\theta'')\right) P\left(B\left(\theta_j^*, \frac{r}{2}\right)\right). \end{aligned}$$

We can also upper bound the integral:

$$\begin{aligned} \int_{\bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta &\leq \exp\left(N \sup_{\theta' \in \bar{\Theta}} \ell_N(\theta')\right) \int_{\bar{\Theta}} p(\theta) d\theta, \\ &= \exp\left(N \sup_{\theta' \in \bar{\Theta}} \ell_N(\theta')\right) P(\bar{\Theta}). \end{aligned}$$

Using these results, we lower bound the ratio as:

$$\begin{aligned} \frac{\int_{\Theta \setminus \bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta} &\geq \frac{\exp\left(N \min_{\theta'' \in B(\theta_j^*, \frac{r}{2})} \ell_N(\theta'')\right) P\left(B\left(\theta_j^*, \frac{r}{2}\right)\right)}{\exp\left(N \sup_{\theta' \in \bar{\Theta}} \ell_N(\theta')\right) P(\bar{\Theta})}, \\ &= \exp\left(N \left(\min_{\theta'' \in B(\theta_j^*, \frac{r}{2})} \ell_N(\theta'') - \sup_{\theta' \in \bar{\Theta}} \ell_N(\theta') \right)\right) \frac{P\left(B\left(\theta_j^*, \frac{r}{2}\right)\right)}{P(\bar{\Theta})}. \end{aligned}$$

1566 Let $\frac{P(B(\theta_j^*, \frac{\tau}{2}))}{P(\bar{\Theta})} = C > 0$ from Assumption 1. As there exists some N' such that
 1567
 1568 $\min_{\theta'' \in B(\theta_j^*, \frac{\tau}{2})} \ell_N(\theta'') > \sup_{\theta' \in \bar{\Theta}} \ell_N(\theta')$ for all $N > N'$, we have shown exists some positive
 1569 constants $c > 0$ and $C > 0$ such that

$$1570 \quad C \exp(cN) \leq \frac{\int_{\Theta \setminus \bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\bar{\Theta}} \exp(N\ell_N(\theta)) p(\theta) d\theta},$$

1571 for all $N > N'$ almost surely, as required. \square

1572
 1573
 1574 We now present our proof of Theorem 2. Here we split the posterior expectation up into small
 1575 regions close to maximising points and regions away from maximising. We then apply our two
 1576 lemmas to each region. Our result then follows by an application the central limit theorem under
 1577 Assumption 1.

1578 **Theorem 2.** *Let the data be drawn from the underlying true distribution $\mathcal{D}_N \sim P_{Data}^*$. Under*
 1579 *Assumption 1, there exists some constant $0 < C < \infty$ such that for sufficiently large N :*

$$1580 \quad \mathbb{E}_{\mathcal{D}_N \sim P_{Data}^*} [\text{Regret}(\mathcal{M}^*, \mathcal{D}_N)] \leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-\frac{Cd}{(1-\gamma)N}\right)}. \quad (18)$$

1581
 1582
 1583 *Proof.* Using the notation introduced at the start of Section D.3, we write the PIL as:

$$1584 \quad \mathcal{I}_N^\pi := \mathbb{E}_{s,a \sim \rho_\pi^*} \left[\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\text{KL} \left(P_{R,S}^*(s,a) \| P_{R,S}(s,a,\theta) \right) \right] \right],$$

$$1585 \quad = \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\log p(r, s' | s, a, \theta^*) - \log p(r, s' | s, a, \theta) \right] \right],$$

$$1586 \quad = \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [\ell^* - \ell(\theta)],$$

1587 Under this same notation, we write the posterior density as:

$$1588 \quad p(\theta | \mathcal{D}_N) = \frac{\exp(N\ell_N(\theta)) p(\theta)}{\int_{\Theta} \exp(N\ell_N(\theta)) p(\theta) d\theta}. \quad (19)$$

1589 Now, under Assumption 1, we split the inner expectation into small regions Θ_i^* around each maximising
 1590 point θ_i^* and the remainder of the parameter space $\bar{\Theta} := \Theta \setminus \bigcup_{i=1}^K \Theta_i^*$:

$$1591 \quad \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [\ell^* - \ell(\theta)] = \sum_{i=1}^K \int_{\Theta_i^*} (\ell^* - \ell(\theta)) p(\theta | \mathcal{D}_N) d\theta + \int_{\bar{\Theta}} (\ell^* - \ell(\theta)) p(\theta | \mathcal{D}_N) d\theta. \quad (20)$$

1592 Using Eq. (19), we now re-write each integral in the summation term of Eq. (20) as:

$$1593 \quad \int_{\Theta_i^*} (\ell^* - \ell(\theta)) p(\theta | \mathcal{D}_N) d\theta = \frac{\int_{\Theta_i^*} (\ell^* - \ell(\theta)) \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta} \exp(N\ell_N(\theta)) p(\theta) d\theta},$$

$$1594 \quad = \frac{\int_{\Theta_i^*} (\ell^* - \ell(\theta)) \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta} \exp(N\ell_N(\theta)) p(\theta) d\theta} \cdot \frac{\int_{\Theta_i^*} \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta_i^*} \exp(N\ell_N(\theta)) p(\theta) d\theta},$$

$$1595 \quad = \frac{\int_{\Theta_i^*} (\ell^* - \ell(\theta)) \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta_i^*} \exp(N\ell_N(\theta)) p(\theta) d\theta} \cdot \frac{\int_{\Theta_i^*} \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta} \exp(N\ell_N(\theta)) p(\theta) d\theta},$$

$$1596 \quad = \frac{\int_{\Theta_i^*} (\ell^* - \ell(\theta)) \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta_i^*} \exp(N\ell_N(\theta)) p(\theta) d\theta} \cdot P(\Theta_i^* | \mathcal{D}_N),$$

$$1597 \quad \leq \frac{\int_{\Theta_i^*} (\ell^* - \ell(\theta)) \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta_i^*} \exp(N\ell_N(\theta)) p(\theta) d\theta}. \quad (21)$$

1598 where we have used $0 \leq P(\Theta_i^* | \mathcal{D}_N) \leq 1$ from Kolmogorov's axioms to bound the final line.

1599 For the last term in Eq. (20), we note that $\ell^* - \ell(\theta)$ is bounded P_Θ -almost surely from Assumption 1,
 1600 hence there exists some $\ell^\dagger < \infty$ such that:

$$1601 \quad \int_{\bar{\Theta}} (\ell^* - \ell(\theta)) p(\theta | \mathcal{D}_N) d\theta \leq \int_{\bar{\Theta}} \ell^\dagger p(\theta | \mathcal{D}_N) d\theta,$$

$$1602 \quad = \ell^\dagger P(\bar{\Theta} | \mathcal{D}_N). \quad (22)$$

Using Ineqs. 21 and 22, we bound Eq. (20) as:

$$\mathbb{E}_{\theta \sim P_{\Theta}(\mathcal{D}_N)} [\ell^* - \ell(\theta)] \leq \sum_{i=1}^K \frac{\int_{\Theta_i^*} (\ell^* - \ell(\theta)) \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta_i^*} \exp(N\ell_N(\theta)) p(\theta) d\theta} + \ell^\dagger P(\bar{\Theta}|\mathcal{D}_N),$$

and hence the PIL can be bounded as:

$$\mathcal{I}_N^\pi \leq \sum_{i=1}^K \frac{\int_{\Theta_i^*} (\ell^* - \ell(\theta)) \exp(N\ell_N(\theta)) p(\theta) d\theta}{\int_{\Theta_i^*} \exp(N\ell_N(\theta)) p(\theta) d\theta} + \ell^\dagger P(\bar{\Theta}|\mathcal{D}_N).$$

Applying Lemma 3 and Lemma 4 under Assumption 1 yields:

$$\begin{aligned} \mathcal{I}_N^\pi &= \sum_{i=1}^K \mathcal{O} \left(\frac{d - g_{i,N}^{*\top} H_i^{*-1} g_{i,N}^*}{N} \right) + \ell^\dagger \mathcal{O}(\exp(-N)), \\ &= \mathcal{O} \left(\frac{d - \sum_{i=1}^K g_{i,N}^{*\top} H_i^{*-1} g_{i,N}^*}{N} \right). \end{aligned} \quad (23)$$

almost surely. As $f(x) := 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-\frac{x}{(1-\gamma)}\right)}$ is monotonic in x and $\frac{d - \sum_{i=1}^K g_{i,N}^{*\top} H_i^{*-1} g_{i,N}^*}{N} \geq 0$, Eq. (23) implies there exists some positive $0 < C < \infty$ such that:

$$\text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-C \frac{d - \sum_{i=1}^K g_{i,N}^{*\top} H_i^{*-1} g_{i,N}^*}{(1-\gamma)N}\right)},$$

almost surely for large enough N . Under Assumption 1, $g_{i,N}^* \xrightarrow{d} \mathcal{N}(0, \Sigma_i^g)$. As $f(x)$ is also a bounded, continuous function and concave, we can apply the Portmanteau Theorem (see for example Bass (2013, Chapter 21.7)) followed by Jensen's inequality to yield:

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_N \sim P_{\text{Data}}^*} [\text{Regret}(\mathcal{M}^*, \mathcal{D}_N)] &\leq 2\mathcal{R}_{\max} \cdot \mathbb{E}_{g_i \sim \mathcal{N}(0, \Sigma_i^g)} \left[\sqrt{1 - \exp\left(-C \frac{d - \sum_{i=1}^K g_i^\top H_i^{*-1} g_i}{(1-\gamma)N}\right)} \right], \\ &\leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-C \frac{d - \sum_{i=1}^K \mathbb{E}_{g_i \sim \mathcal{N}(0, \Sigma_i^g)} [g_i^\top H_i^{*-1} g_i]}{(1-\gamma)N}\right)}, \\ &= 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-C \frac{d - \sum_{i=1}^K \text{Tr}(\Sigma_i^g H_i^{*-1})}{(1-\gamma)N}\right)}. \end{aligned} \quad (24)$$

Now, examining the Hessian:

$$\begin{aligned} H(\theta) &= \nabla_\theta^2 \mathbb{E}_{s, a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s, a)} [\log p(r, s' | s, a, \theta)] \\ &= \nabla_\theta \mathbb{E}_{s, a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s, a)} [\nabla_\theta \log p(r, s' | s, a, \theta)], \\ &= \nabla_\theta \mathbb{E}_{s, a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s, a)} \left[\frac{\nabla_\theta p(r, s' | s, a, \theta)}{p(r, s' | s, a, \theta)} \right], \\ &= \mathbb{E}_{s, a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s, a)} \left[\nabla_\theta \frac{\nabla_\theta p(r, s' | s, a, \theta)}{p(r, s' | s, a, \theta)} \right], \\ &= \mathbb{E}_{s, a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s, a)} \left[\frac{\nabla_\theta^2 p(r, s' | s, a, \theta)}{p(r, s' | s, a, \theta)} \right] \\ &\quad - \mathbb{E}_{s, a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s, a)} \left[\frac{\nabla_\theta p(r, s' | s, a, \theta)}{p(r, s' | s, a, \theta)} \frac{\nabla_\theta p(r, s' | s, a, \theta)^\top}{p(r, s' | s, a, \theta)} \right], \end{aligned} \quad (25)$$

Hence at $\theta = \theta_i^*$, the first term of Eq. (25) is:

$$\begin{aligned}
\mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\frac{\nabla_\theta^2 p(r, s' | s, a, \theta)}{p(r, s' | s, a, \theta_i^*)} \right] &= \mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\frac{\nabla_\theta^2 p(r, s' | s, a, \theta)|_{\theta=\theta_i^*}}{p^*(r, s' | s, a)} \right], \\
&= \mathbb{E}_{s,a \sim \rho_\pi^*} \left[\int_{\mathbb{R} \times \mathcal{S}} \nabla_\theta^2 p(r, s' | s, a, \theta)|_{\theta=\theta_i^*} d(r, s') \right], \\
&= \mathbb{E}_{s,a \sim \rho_\pi^*} \left[\nabla_\theta^2 \int_{\mathbb{R} \times \mathcal{S}} p(r, s' | s, a, \theta) d(r, s') |_{\theta=\theta_i^*} \right], \\
&= \nabla_\theta^2 1|_{\theta=\theta_i^*}, \\
&= 0,
\end{aligned}$$

hence:

$$\begin{aligned}
H(\theta_i^*) &= 0 - \mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\frac{\nabla_\theta p(r, s' | s, a, \theta_i^*)}{p(r, s' | s, a, \theta_i^*)} \frac{\nabla_\theta p(r, s' | s, a, \theta_i^*)^\top}{p(r, s' | s, a, \theta_i^*)} \right], \\
&= -\mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\nabla_\theta \log p(r, s' | s, a, \theta_i^*) \nabla_\theta \log p(r, s' | s, a, \theta_i^*)^\top \right], \\
&= -\Sigma_i^g.
\end{aligned}$$

Using this result, each $\text{Tr}(\Sigma_i^g H_i^{*-1}) = \text{Tr}(-I) = -d$. Substituting yields:

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}_N \sim P_{\text{Data}}^*} [\text{Regret}(\mathcal{M}^*, \mathcal{D}_N)] &\leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-C \frac{(k+1)d}{(1-\gamma)N}\right)}, \\
&\leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-C' \frac{d}{(1-\gamma)N}\right)},
\end{aligned}$$

for some $0 < C' < \infty$ and sufficiently large N , as required. \square

We note that Theorem 2 applies to the Gaussian world model introduced in Section 4.4 with neural network mean functions with C^2 -continuous activations (tanh, identity, sigmoid, softplus, SiLU, SELU, GELU...) using a Gaussian or uniform prior truncated to a compact parameter space and similarly well-behaved parametric models. The resulting differences in performance only arise from the choice of prior, model representability and coverage of dataset, which affect Bayesian and frequentist methods equally. The information rate coincides with the optimal ‘minimax’ convergence rate of frequentist parametric density estimators (Yang and Barron, 1999; Bilodeau et al., 2021). Similar results for the information rate have been found for nonparametric models such a Gaussian processes (van der Vaart and van Zanten, 2011).

Using our result in Theorem 2, we plot the normalised regret bound (i.e. taking $\mathcal{R}_{\max} = 0.5$) in Ineq. 18 for increasing dimensionality (blue) and decreasing γ (copper) in Fig. 5. Our bound reveals an S-shaped curve with three distinct phases as number of data points N increases: an initial plateau, a sudden decrease in regret follow by a slow exponential decay towards a regret of zero. The plateau indicates that a minimum amount of data is needed before any benefit can be realised in terms of regret. This is to be expected because initially the only information about the parameter values is given by the prior, which has no guarantee of accuracy under our analysis. Once a threshold of data points has been reached, the data can start to overwhelm the prior, resulting in a sudden decrease in regret. The higher the dimensionality of the model, the greater this data limit is - represented in Fig. 5 by the plateau length increasing with greater d (blue curves). Due to overspecification in models, this limit is likely to be set by the effective dimension of

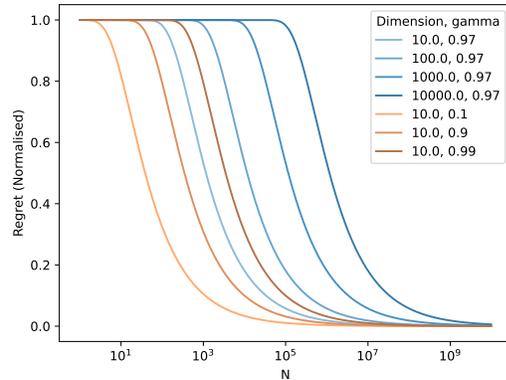


Figure 5: Normalised Regret Curves for $C = 1$

the problem (which may be much lower than d) as many parameters will be redundant, however the effective dimension is typically not possible to ascertain a priori. Finally, we observe that increasing the discount factor γ leads to a longer regret plateau (copper curves) due to any error in the model dynamics being compounded over a longer horizon at test time.

D.4 EXTENSIONS FOR MODEL MISSPECIFICATION AND SUB-OPTIMAL POLICY LEARNING

We now generalise our theorems to include the effects of model misspecification, that is models that cannot fully represent the true environment dynamics, and sub-optimal Bayesian policy learning, that is the effect of using a policy that does not fully optimise the Bayesian RL objective. We use the dagger notation to denote the maximum cross entropy parametrisation:

$$\theta_i^\dagger \in \arg \max_{\theta \in \Theta} \ell(\theta) = \arg \min_{\theta \in \Theta} \mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} [\text{KL}(P_{R,S}^*(s,a) | P_{R,S}(s,a, \theta))].$$

To characterise the degree of model misspecification, we use the KL divergence:

$$\epsilon_{\text{miss}} := \min_{\theta \in \Theta} \mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} [\text{KL}(P_{R,S}^*(s,a) | P_{R,S}(s,a, \theta))].$$

We also introduce the following simplifying dagger notation for the expected cross entropy and corresponding gradients and Hessian under the optimal parameter:

$$\begin{aligned} \ell^\dagger &:= \max_{\theta \in \Theta} \ell(\theta) = \mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} [\log p(r, s' | s, a, \theta_i^\dagger)], \\ g_{i,N}^\dagger &:= \sqrt{N} \nabla_\theta \ell_N(\theta) \big|_{\theta=\theta_i^\dagger}, \\ H_i^\dagger &:= \nabla_\theta^2 \ell(\theta) \big|_{\theta=\theta_i^\dagger}. \end{aligned}$$

We now relax Assumption 2 to allow for model misspecification:

Assumption 2. *We assume that:*

- i *The maximum likelihood is finite $|\ell^\dagger| < \infty$ and $|\ell^\dagger - \ell(\theta)|$ is bounded P_Θ -almost surely.*
- ii *$\ell_N(\theta)$ and $\ell(\theta)$ are C^2 -continuous in θ .*
- iii *There are $K < \infty$ maximising points θ_i^\dagger :*

$$\{\theta_1^\dagger, \theta_2^\dagger, \dots, \theta_K^\dagger\} = \arg \max_{\theta \in \Theta} \ell(\theta).$$

For each maximiser θ_i^\dagger , there exists a small region $\Theta_i^\dagger := \{\theta \in \Theta \mid \|\theta_i^\dagger - \theta\| \leq \epsilon\}$ for some $\epsilon > 0$ such that θ_i^\dagger is the unique maximiser in Θ_i^\dagger , θ_i^\dagger is in the interior of Θ_i^\dagger , $\nabla_\theta^2 \ell(\theta_i^\dagger)$ is negative definite, invertible and the regions are disjoint: $\bigcap_{i=1}^K \Theta_i^\dagger = \emptyset$.

iv *The prior $p(\theta)$ is Lipschitz continuous in θ with support over Θ .*

- v *The sampling regime ensures that the strong law of large numbers holds for all maximisers θ_i^\dagger for the Hessian, and uniformly for $\theta \in \Theta$ for the likelihood, that is:*

$$\ell_N(\theta) \xrightarrow{\text{Unif. a.s.}} \ell(\theta), \quad \nabla_\theta^2 \ell_N(\theta_i^\dagger) \xrightarrow{\text{a.s.}} \nabla_\theta^2 \ell(\theta_i^\dagger).$$

The central limit theorem applies to the gradient at each θ_i^\dagger , that is:

$$\sqrt{N} \nabla_\theta \ell_N(\theta_i^\dagger) \xrightarrow{d} \mathcal{N}(0, \Sigma_i^g),$$

where $\Sigma_i^g = \mathbb{E}_{s,a \sim \rho_\pi^, r, s' \sim P_{R,S}^*(s,a)} [\nabla_\theta \log p(r, s' | s, a, \theta_i^\dagger) \nabla_\theta \log p(r, s' | s, a, \theta_i^\dagger)^\top]$ with $\|\Sigma_i^g\| < \infty$.*

Finally, we account for let the Bayes sub-optimality be defined as

$$\epsilon_{\text{Bayes}} := \left| J_{\text{Bayes}}^{\pi^*}(P_\Phi(\mathcal{D}_N)) - J_{\text{Bayes}}^{\hat{\pi}}(P_\Phi(\mathcal{D}_N)) \right|. \quad (26)$$

Lemma 5. Let $\mathcal{R}_{\max} := \frac{(r_{\max} - r_{\min})}{1 - \gamma}$ denote the maximum possible regret for the MDP and the Bayes sub-optimality ϵ_{Bayes} be defined as in Eq. (26). For a prior $P_{\Theta}(\mathcal{D}_N)$, the true regret can be bounded as:

$$\text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \leq 2\mathcal{R}_{\max} \cdot \sup_{\pi} \mathbb{E}_{i \sim \mathcal{G}(\gamma)} [TV(P_{i+1, \pi}^* \| P_{i+1, \pi}^{\pi}(\mathcal{D}_N))] + \epsilon_{\text{Bayes}}.$$

Proof. We start from the definition of the true regret under Bayes sub-optimality:

$$\begin{aligned} \text{Regret}(\mathcal{M}^*, \mathcal{D}_N) &:= J^{\pi^*}(\mathcal{M}^*) - J^{\hat{\pi}}(\mathcal{M}^*, \mathcal{D}_N), \\ &= J^{\pi^*}(\mathcal{M}^*) - J_{\text{Bayes}}^{\pi^*}(P_{\Phi}(\mathcal{D}_N)) + J_{\text{Bayes}}^{\pi^*}(P_{\Phi}(\mathcal{D}_N)) - J^{\hat{\pi}}(\mathcal{M}^*, \mathcal{D}_N), \\ &\leq J^{\pi^*}(\mathcal{M}^*) - J_{\text{Bayes}}^{\pi^*}(P_{\Phi}(\mathcal{D}_N)) + J_{\text{Bayes}}^{\hat{\pi}}(P_{\Phi}(\mathcal{D}_N)) - J^{\hat{\pi}}(\mathcal{M}^*, \mathcal{D}_N), \\ &= J^{\pi^*}(\mathcal{M}^*) - J_{\text{Bayes}}^{\pi^*}(P_{\Phi}(\mathcal{D}_N)) + J_{\text{Bayes}}^{\hat{\pi}}(P_{\Phi}(\mathcal{D}_N)) - J^{\hat{\pi}}(\mathcal{M}^*, \mathcal{D}_N) \\ &\quad + J_{\text{Bayes}}^{\pi^*}(P_{\Phi}(\mathcal{D}_N)) - J_{\text{Bayes}}^{\hat{\pi}}(P_{\Phi}(\mathcal{D}_N)), \\ &\leq \sup_{\pi} |J^{\pi}(\mathcal{M}^*) - J_{\text{Bayes}}^{\pi}(P_{\Theta}(\mathcal{D}_N))| + \sup_{\pi} |J_{\text{Bayes}}^{\pi}(P_{\Theta}(\mathcal{D}_N)) - J^{\pi}(\mathcal{M}^*)| \\ &\quad + |J_{\text{Bayes}}^{\pi^*}(P_{\Phi}(\mathcal{D}_N)) - J_{\text{Bayes}}^{\hat{\pi}}(P_{\Phi}(\mathcal{D}_N))|, \\ &= 2 \sup_{\pi} |J^{\pi}(\mathcal{M}^*) - J_{\text{Bayes}}^{\pi}(P_{\Theta}(\mathcal{D}_N))| + \epsilon_{\text{Bayes}}, \end{aligned}$$

where the second line follows from $J_{\text{Bayes}}^{\pi^*}(P_{\Phi}(\mathcal{D}_N)) \leq J_{\text{Bayes}}^{\hat{\pi}}(P_{\Phi}(\mathcal{D}_N))$ by definition. We then bound $|J^{\pi}(\mathcal{M}^*) - J_{\text{Bayes}}^{\pi}(P_{\Theta}(\mathcal{D}_N))|$ using Lemma 1 to obtain our desired result. \square

Theorem 3. Let the data be drawn from the underlying true distribution $\mathcal{D}_N \sim P_{\text{Data}}^*$. Under Assumption 2, there exists some constant $0 < C < \infty$ such that for sufficiently large N :

$$\mathbb{E}_{\mathcal{D}_N \sim P_{\text{Data}}^*} [\text{Regret}(\mathcal{M}^*, \mathcal{D}_N)] \leq 2\mathcal{R}_{\max} \cdot \exp\left(1 - \sqrt{C \left(\frac{d}{(1-\gamma)N} + \frac{\epsilon_{\text{miss}}}{1-\gamma}\right)}\right) + \epsilon_{\text{Bayes}}.$$

Proof. Starting with Lemma 5, we obtain:

$$\text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-C \left(\frac{d^2}{(1-\gamma)N} + \frac{\epsilon_{\text{miss}}}{(1-\gamma)}\right)\right)} + \epsilon_{\text{Bayes}}.$$

Next we apply Theorem 1 to bound the first term, obtaining:

$$\text{Regret}(\mathcal{M}^*, \mathcal{D}_N) \leq 2\mathcal{R}_{\max} \cdot \sup_{\pi} \sqrt{1 - \exp\left(-\frac{\mathcal{I}_{\pi}^{\pi}}{1-\gamma}\right)} + \epsilon_{\text{Bayes}}.$$

We write the PIL to include misspecification as:

$$\begin{aligned}
\mathcal{I}_N^\pi &:= \mathbb{E}_{s,a \sim \rho_\pi^*} \left[\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\text{KL} \left(P_{R,S}^*(s,a) \| P_{R,S}(s,a,\theta) \right) \right] \right], \\
&= \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\log p(r, s' | s, a, \theta^*) - \log p(r, s' | s, a, \theta) \right] \right], \\
&= \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\log p(r, s' | s, a, \theta^*) - \log p(r, s' | s, a, \theta_i^\dagger) \right. \right. \\
&\quad \left. \left. + \log p(r, s' | s, a, \theta_i^\dagger) - \log p(r, s' | s, a, \theta) \right] \right], \\
&= \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\log p(r, s' | s, a, \theta^*) - \log p(r, s' | s, a, \theta_i^\dagger) \right] \right. \\
&\quad \left. + \mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\log p(r, s' | s, a, \theta_i^\dagger) - \log p(r, s' | s, a, \theta) \right] \right], \\
&= \mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\log p(r, s' | s, a, \theta^*) - \log p(r, s' | s, a, \theta_i^\dagger) \right] \\
&\quad + \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\log p(r, s' | s, a, \theta_i^\dagger) - \log p(r, s' | s, a, \theta) \right] \right], \\
&= \mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\text{KL} \left(P_{R,S}^*(s,a) | P_{R,S}(s,a,\theta_i^\dagger) \right) \right] \\
&\quad + \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\log p(r, s' | s, a, \theta_i^\dagger) - \log p(r, s' | s, a, \theta) \right] \right], \\
&= \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\ell^\dagger - \ell(\theta) \right] + \mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\text{KL} \left(P_{R,S}^*(s,a) | P_{R,S}(s,a,\theta_i^\dagger) \right) \right], \\
&= \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\ell^\dagger - \ell(\theta) \right] + \min_{\theta} \mathbb{E}_{s,a \sim \rho_\pi^*, r, s' \sim P_{R,S}^*(s,a)} \left[\text{KL} \left(P_{R,S}^*(s,a) | P_{R,S}(s,a,\theta) \right) \right], \\
&= \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\ell^\dagger - \ell(\theta) \right] + \epsilon_{\text{miss}}.
\end{aligned}$$

To bound the first term $\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\ell^\dagger - \ell(\theta) \right]$, we follow the remainder of the proof of Theorem 2 to Ineq. 24, replacing \star notion with \dagger to yield:

$$\mathcal{I}_N^\pi = \epsilon_{\text{miss}} + \mathcal{O} \left(\frac{d - \sum_{i=1}^K g_{i,N}^\dagger \top H_i^{\dagger^{-1}} g_{i,N}^\dagger}{N} \right). \quad (27)$$

almost surely. As $f(x) := 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-\frac{x}{(1-\gamma)}\right)}$ is monotonic in x and $\frac{d - \sum_{i=1}^K g_{i,N}^\dagger \top H_i^{\dagger^{-1}} g_{i,N}^\dagger}{N} + \epsilon_{\text{miss}} \geq 0$, Eq. (27) implies there exists some positive $0 < C < \infty$ such that:

$$\text{Regret}(\mathcal{M}^\dagger, \mathcal{D}_N) \leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-C \left(\frac{d - \sum_{i=1}^K g_{i,N}^\dagger \top H_i^{\dagger^{-1}} g_{i,N}^\dagger}{(1-\gamma)N} + \frac{\epsilon_{\text{miss}}}{(1-\gamma)} \right)\right)},$$

almost surely for large enough N . Under Assumption 2, $g_{i,N}^\dagger \xrightarrow{d} \mathcal{N}(0, \Sigma_i^g)$. As $f(x)$ is also a bounded, continuous function and concave, we can apply the Portmanteau Theorem (see for example Bass (2013, Chapter 21.7)) followed by Jensen's inequality to yield:

$$\begin{aligned}
&\mathbb{E}_{\mathcal{D}_N \sim P_{\text{Data}}^\dagger} [\text{Regret}(\mathcal{M}^\dagger, \mathcal{D}_N)] \\
&\leq 2\mathcal{R}_{\max} \cdot \mathbb{E}_{g_i \sim \mathcal{N}(0, \Sigma_i^g)} \left[\sqrt{1 - \exp\left(-C \left(\frac{d - \sum_{i=1}^K g_i \top H_i^{\dagger^{-1}} g_i}{(1-\gamma)N} + \frac{\epsilon_{\text{miss}}}{(1-\gamma)} \right)\right)} \right], \\
&\leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-C \left(\frac{d - \sum_{i=1}^K \mathbb{E}_{g_i \sim \mathcal{N}(0, \Sigma_i^g)} \left[g_i \top H_i^{\dagger^{-1}} g_i \right]}{(1-\gamma)N} + \frac{\epsilon_{\text{miss}}}{(1-\gamma)} \right)\right)}, \\
&= 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-C \left(\frac{d - \sum_{i=1}^K \text{Tr} \left(\Sigma_i^g H_i^{\dagger^{-1}} \right)}{(1-\gamma)N} + \frac{\epsilon_{\text{miss}}}{(1-\gamma)} \right)\right)}.
\end{aligned}$$

Now $\text{Tr}(\sum_i^g H_i^*{}^{-1}) = \mathcal{O}(d^2)$, hence

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_N \sim P_{\text{Data}}^\dagger} [\text{Regret}(\mathcal{M}^\dagger, \mathcal{D}_N)] &\leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-C \left(\frac{(k+1)d^2}{(1-\gamma)N} + \frac{\epsilon_{\text{miss}}}{(1-\gamma)}\right)\right)}, \\ &\leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-C' \left(\frac{d^2}{(1-\gamma)N} + \frac{\epsilon_{\text{miss}}}{(1-\gamma)}\right)\right)}, \end{aligned}$$

for some $0 < C' < \infty$ and sufficiently large N , as required. \square

Taking the limit $N \rightarrow \infty$, we see the residual term due to misspecification and sub-optimal policy learning is:

$$\mathbb{E}_{\mathcal{D}_N \sim P_{\text{Data}}^\dagger} [\text{Regret}(\mathcal{M}^\dagger, \mathcal{D}_N)] \leq 2\mathcal{R}_{\max} \cdot \sqrt{1 - \exp\left(-\frac{\epsilon_{\text{miss}} C}{1-\gamma}\right)} + \epsilon_{\text{Bayes}}.$$

We compare against prior work such as MOREL (Kidambi et al., 2020, Corollary 2), where the residual misspecification term is:

$$\frac{4\gamma\mathcal{R}_{\max}}{1-\gamma} \epsilon_{TV}$$

where ϵ_{TV} characterises the misspecification in terms of total variational distance instead of KL divergence of our method. Crucially, we see that our bound is much less sensitive to γ ; our bound is $\mathcal{O}\left(\frac{1}{\sqrt{1-\gamma}}\right)$ in comparison to $\mathcal{O}\left(\frac{\gamma}{1-\gamma}\right)$ of MOREL meaning our bound is tighter as $\gamma \rightarrow 1$.

E FURTHER RESULTS

E.1 TOReL

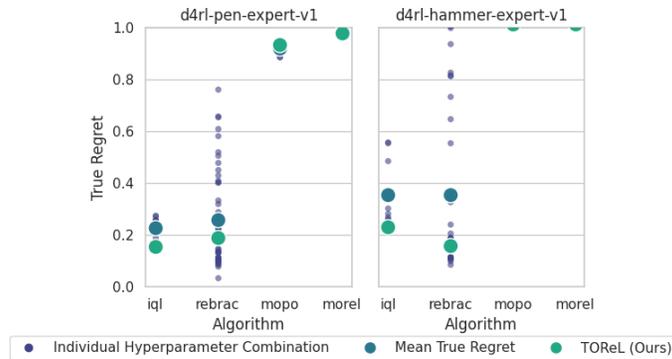


Figure 6: TOReL-hyperparameter regret versus mean hyperparameter regret for Adroit (lower is better).

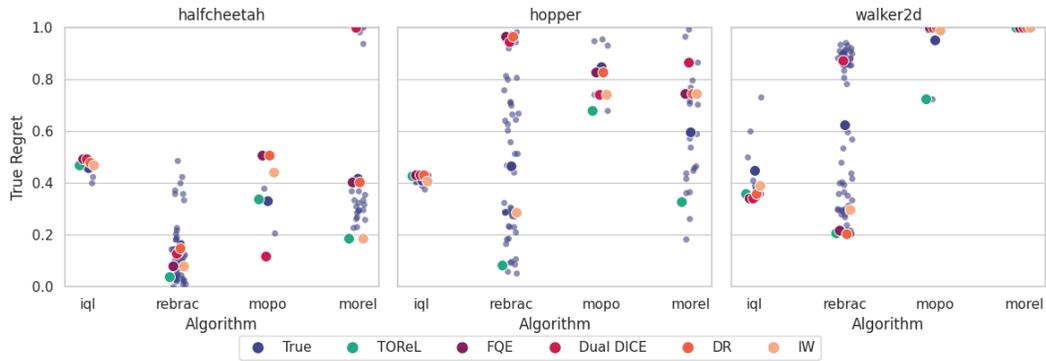


Figure 7: Hyperparameters selected by different OPE metrics for d4rl-halfcheetah-medium-expert-v2, d4rl-hopper-medium-v2 and d4rl-walker-medium-replay-v2 (lower is better).

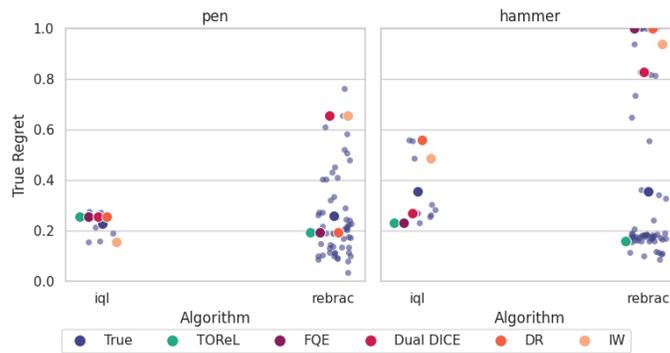


Figure 8: Hyperparameters selected by different OPE metrics for d4rl-pen-expert-v1 and d4rl-hammer-expert-v1 (lower is better).

Task		IQL	ReBRAC	MOPO	MOReL	
1998						
1999						
2000	brax-	True	0.203	0.312	0.403	0.751
2001	halfcheetah-	Oracle	0.186	0.089	0.133	0.641
2002	full-replay	TOReL	<i>0.186</i>	<i>0.089</i>	<i>0.133</i>	<i>0.648</i>
2003	brax-	True	0.550	0.534	0.558	0.575
2004	hopper-	Oracle	0.377	0.070	0.082	0.243
2005	full-replay	TOReL	<i>0.397</i>	<i>0.070</i>	<i>0.086</i>	<i>0.282</i>
2006	brax-	True	0.374	0.357	0.342	0.625
2007	walker-	Oracle	0.304	0.000	0.243	0.415
2008	full-replay	TOReL	<i>0.331</i>	<i>0.000</i>	0.384	<i>0.554</i>
2009						
2010		True	0.469	0.134	0.331	0.418
2011		Oracle	0.400	0.000	0.116	0.187
2012	d4rl-	TOReL	<i>0.469</i>	<i>0.036</i>	0.339	<i>0.187</i>
2013	halfcheetah-	FQE	0.492	0.079	0.507	0.402
2014	medium-expert-v2	DICE	0.492	0.127	<i>0.116</i>	1.000
2015		DR	0.480	0.147	0.507	0.402
2016		IW	0.469	0.079	0.441	0.187
2017		True	0.411	0.467	0.848	0.595
2018		Oracle	0.375	0.053	0.681	0.183
2019	d4rl-	TOReL	0.428	<i>0.083</i>	<i>0.681</i>	<i>0.327</i>
2020	hopper-	FQE	0.432	0.967	0.829	0.744
2021	medium-v2	DICE	0.432	0.945	0.743	0.865
2022		DR	0.432	0.967	0.829	0.744
2023		IW	<i>0.406</i>	0.285	0.743	0.744
2024		True	0.450	0.625	0.952	1.000
2025		Oracle	0.339	0.204	0.724	1.000
2026	d4rl-	TOReL	<i>0.358</i>	<i>0.204</i>	<i>0.724</i>	<i>1.000</i>
2027	walker2d-	FQE	0.339	0.217	1.000	1.000
2028	medium-replay-v2	DICE	<i>0.341</i>	0.872	1.000	1.000
2029		DR	0.358	<i>0.204</i>	1.000	1.000
2030		IW	0.388	0.295	0.990	1.000
2031						
2032						

Table 2: True, oracle, TOReL and OPE regrets across tasks. Bold indicates where TOReL identifies the oracle hyperparameters, while italic indicates where TOReL identifies hyperparameters with a regret lower than the true regret. ReBRAC+TOReL outperforms all algorithms on every dataset. Green indicates where TOReL is the best OPE metric. Orange indicates where another OPE metric beats or ties with TOReL. OPE metrics for the brax datasets are undetermined, as the brax datasets do not contain full trajectories (required for DR and IW) or initial state flags (required for DICE). Given TOReL’s strong performance on the brax datasets (with ReBRAC + TOReL achieving the Oracle hyper-parameters for each dataset), FQE could only have performed on par with it.

Task		IQL	ReBRAC	MOPO	MOReL
d4rl- pen- expert-v1	True	0.226	0.258	0.919	0.985
	Oracle	0.154	0.033	0.885	0.968
	TOReL	0.254	<i>0.192</i>		
	FQE	0.254	0.192		
	DICE	0.254	0.654		
	DR	0.254	0.192		
	IW	0.154	0.654		
d4rl- hammer- expert-v1	True	0.355	0.354	1.000	1.000
	Oracle	0.229	0.086	1.000	1.000
	TOReL	0.229	<i>0.159</i>		
	FQE	0.229	1.000		
	DICE	0.268	0.827		
	DR	0.558	1.000		
	IW	0.485	0.938		

Table 3: True, oracle, TOReL and OPE regrets across tasks. Bold indicates where TOReL identifies the oracle hyperparameters, while italic indicates where TOReL identifies hyperparameters with a regret lower than the true regret. Green indicates where TOReL is the best OPE metric. Orange indicates where another OPE metric beats or ties with TOReL. All metrics are unreliable for the pen and hammer datasets due to the quality of the data. We refer the reader to the corresponding scatter-plots.

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

Task			IQL	ReBRAC	MOPO	MOReL	
brax-halfcheetah-full-replay	TOReL	r	0.29	0.92	0.98	0.93	
		p	0.448	0.000	0.000	0.000	
brax-hopper-full-replay	TOReL	r	0.18	0.98	1.00	0.98	
		p	0.635	0.000	0.00	0.000	
brax-walker-full-replay	TOReL	r	0.32	nan	-0.68	0.99	
		p	0.406	nan	0.133	0.000	
d4rl-halfcheetah-medium-expert-v2	TOReL	r	0.29	0.90	0.02	0.98	
		p	0.442	0.000	0.975	0.000	
	FQE	r	-0.21	0.66	-0.32	0.97	
		p	0.591	0.000	0.530	0.000	
	DICE	r	-0.47	0.19	0.74	0.65	
		p	0.204	0.166	0.094	0.001	
	DR	r	-0.24	0.32	-0.27	0.96	
		p	0.537	0.017	0.602	0.000	
	IW	r	-0.33	-0.44	0.07	0.97	
		p	0.393	0.001	0.901	0.000	
	d4rl-hopper-medium-v2	TOReL	r	0.29	0.98	0.98	0.94
			p	0.443	0.000	0.001	0.000
FQE		r	0.16	0.27	0.87	0.37	
		p	0.674	0.044	0.025	0.0078	
DICE		r	0.19	0.34	0.31	0.09	
		p	0.623	0.009	0.554	0.673	
DR		r	0.26	0.20	0.78	0.34	
		p	0.499	0.135	0.065	0.099	
IW		r	-0.39	0.27	0.24	0.24	
		p	0.301	0.042	0.647	0.264	
d4rl-walker2d-medium-replay-v2		TOReL	r	0.65	0.78	1.00	-0.53
			p	0.057	0.000	0.000	0.008
	FQE	r	0.41	0.89	-0.48	-0.05	
		p	0.276	0.000	0.334	0.810	
	DICE	r	0.52	-0.04	-0.14	-0.07	
		p	0.148	0.700	0.785	0.744	
	DR	r	0.19	0.82	-0.29	-0.06	
		p	0.631	0.000	0.571	0.780	
	IW	r	-0.06	0.86	0.40	0.96	
		p	0.875	0.000	0.427	0.000	

Table 4: Pearson correlation (r) and statistical significance (p) between TOReL regret metrics and true regrets for different hyperparameter combinations. Even where no strong positive correlation is observed (possibly due to limited hyperparameter coverage), the TOReL regret is lower than the true regret averaged over those tasks. Green indicates where TOReL has strong ($r > |0.5|$), statically significant ($p < 0.05$) positive correlation, while orange indicates where other OPE metrics have strong, statistically significant correlation.

Task			IQL	ReBRAC
d4rl-pen-expert-v1	TOReL	r	nan	nan
		p	nan	nan
	FQE	r	nan	nan
		p	nan	nan
	DICE	r	0.49	0.00
		p	0.178	0.980
	DR	r	nan	0.26
		p	nan	0.056
	IW	r	0.26	-0.66
		p	0.497	0.000
d4rl-hammer-expert-v1	TOReL	r	0.70	0.30
		p	0.035	0.023
	FQE	r	0.74	0.11
		p	0.023	0.422
	DICE	r	-0.12	-0.14
		p	0.755	0.313
	DR	r	-0.35	-0.04
		p	0.349	0.745
	IW	r	-0.40	-0.66
		p	0.287	0.000

Table 5: Pearson correlation (r) and statistical significance (p) between TOReL regret metrics and true regrets for different hyperparameter combinations. Green indicates where TOReL has strong ($r > |0.5|$), statically significant ($p < 0.05$) positive correlation, while orange indicates where other OPE metrics have strong, statistically significant correlation. All OPE metrics perform poorly due to the nature of the datasets.

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

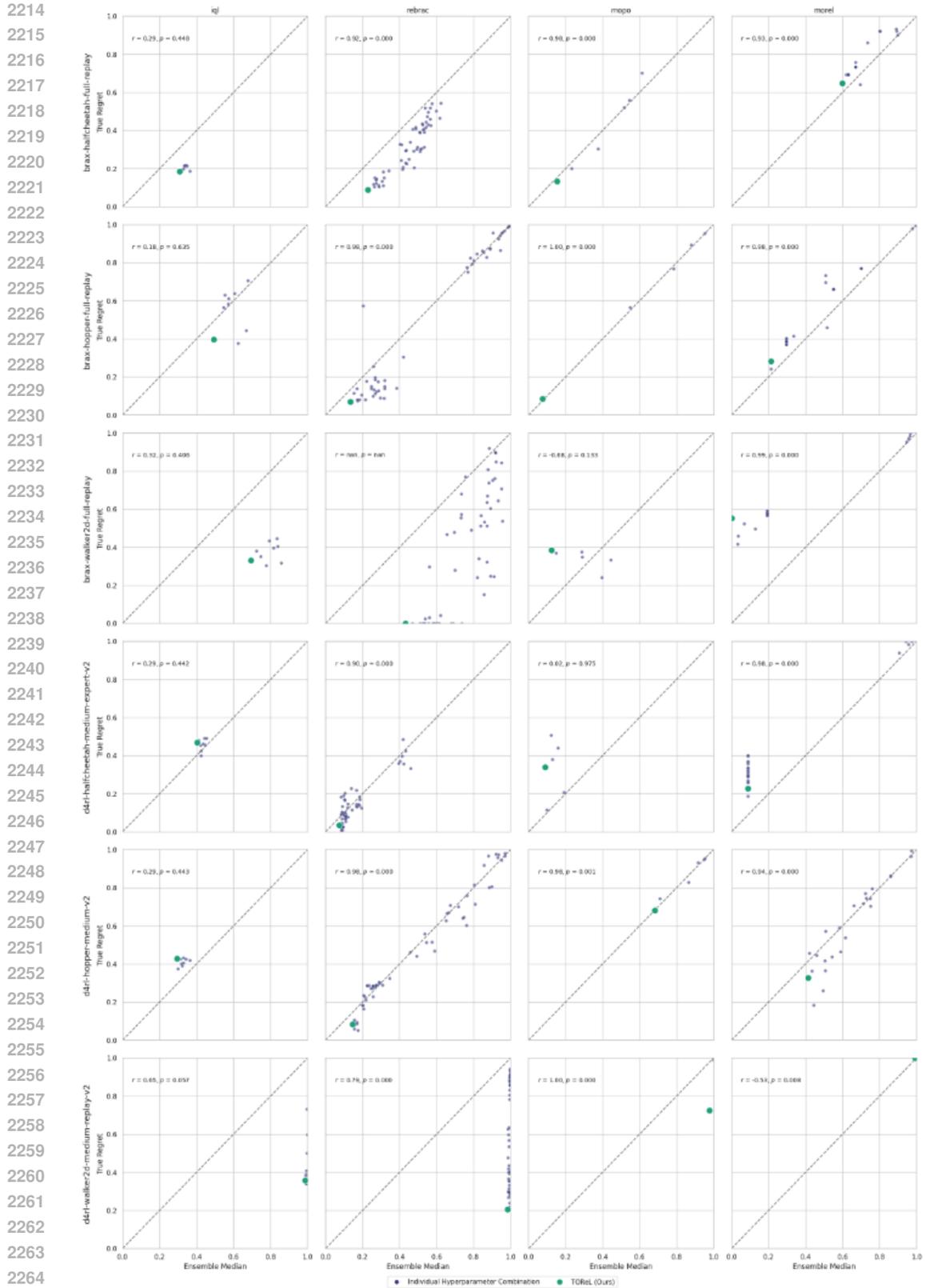


Figure 9: Scatter plots to visualise the positive correlation between the TOReL regret metric and the true regret.

2268
2269
2270
2271
2272
2273
2274

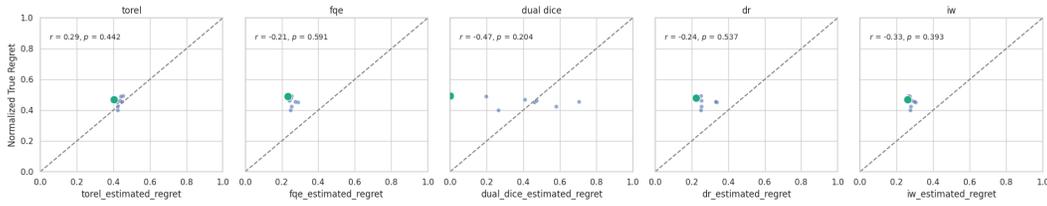


Figure 10: OPE metric scatter-plots: IQL-trained policies, d4rl-halfcheetah-medium-expert-v2.

2275
2276
2277
2278
2279

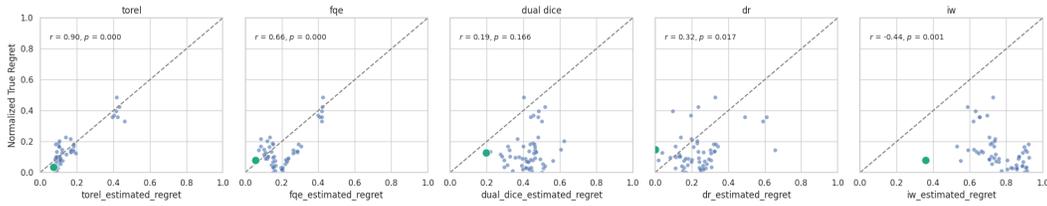


Figure 11: OPE metric scatter-plots: ReBRAC-trained policies, d4rl-halfcheetah-medium-expert-v2.

2280
2281
2282
2283
2284
2285
2286

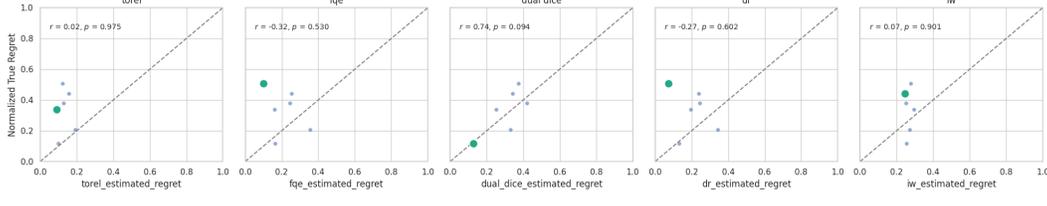


Figure 12: OPE metric scatter-plots: MOPO-trained policies, d4rl-halfcheetah-medium-expert-v2.

2290
2291
2292
2293
2294
2295
2296
2297

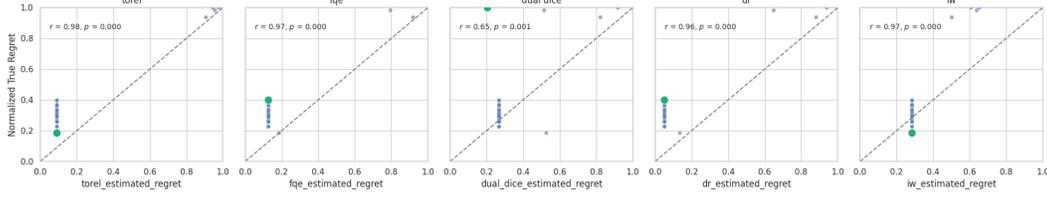


Figure 13: OPE metric scatter-plots: MOREL-trained policies, d4rl-halfcheetah-medium-expert-v2.

2300
2301
2302
2303
2304
2305
2306
2307
2308
2309

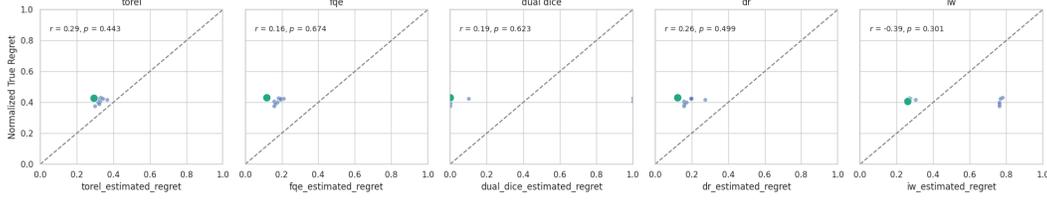


Figure 14: OPE metric scatter-plots: IQL-trained policies, d4rl-hopper-medium-v2.

2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321

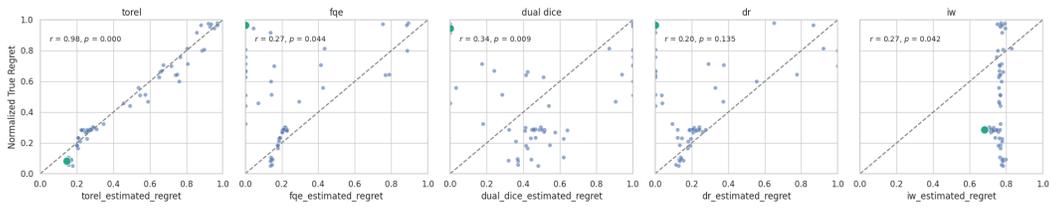


Figure 15: OPE metric scatter-plots: ReBRAC-trained policies, d4rl-hopper-medium-v2.

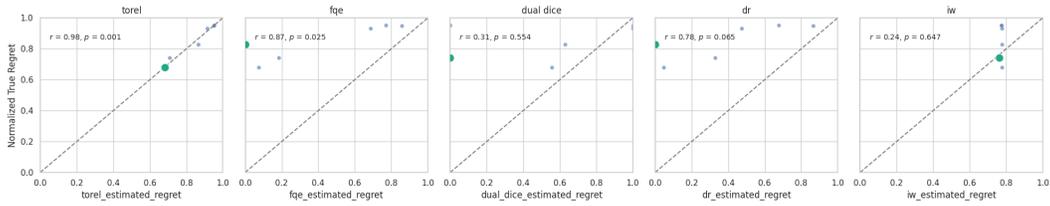


Figure 16: OPE metric scatter-plots: MOPO-trained policies, d4rl-hopper-medium-v2.

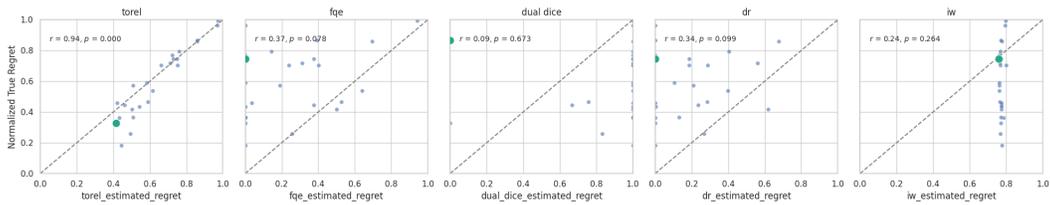


Figure 17: OPE metric scatter-plots: MOREL-trained policies, d4rl-hopper-medium-v2.

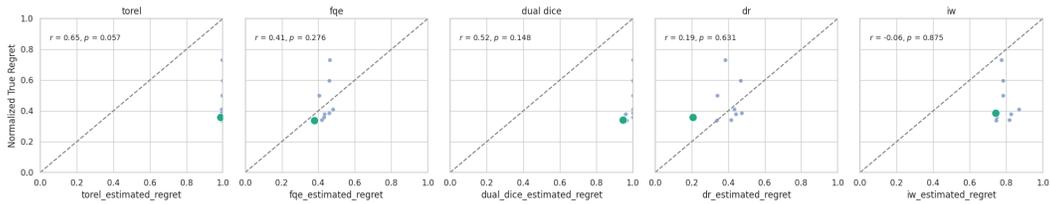


Figure 18: OPE metric scatter-plots: IQL-trained policies, d4rl-walker-medium-replay-v2.

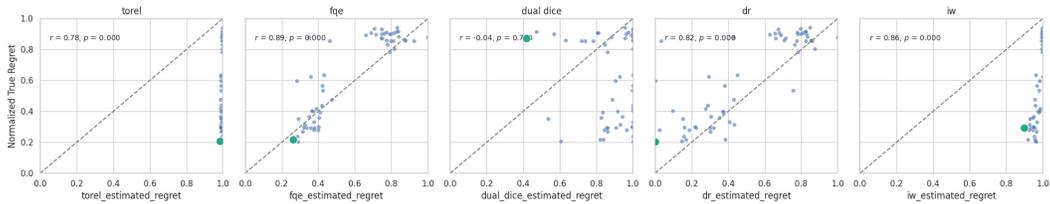


Figure 19: OPE metric scatter-plots: ReBRAC-trained policies, d4rl-walker-medium-replay-v2.

2376
 2377
 2378
 2379
 2380
 2381
 2382
 2383
 2384
 2385
 2386
 2387
 2388
 2389
 2390
 2391
 2392
 2393
 2394
 2395
 2396
 2397
 2398
 2399
 2400
 2401
 2402
 2403
 2404
 2405
 2406
 2407
 2408
 2409
 2410
 2411
 2412
 2413
 2414
 2415
 2416
 2417
 2418
 2419
 2420
 2421
 2422
 2423
 2424
 2425
 2426
 2427
 2428
 2429

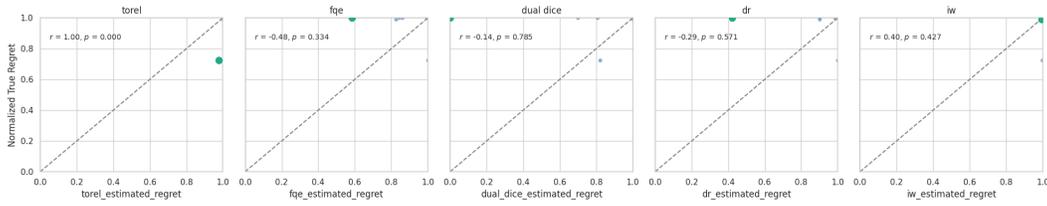


Figure 20: OPE metric scatter-plots: MOPO-trained policies, d4rl-walker-medium-replay-v2.

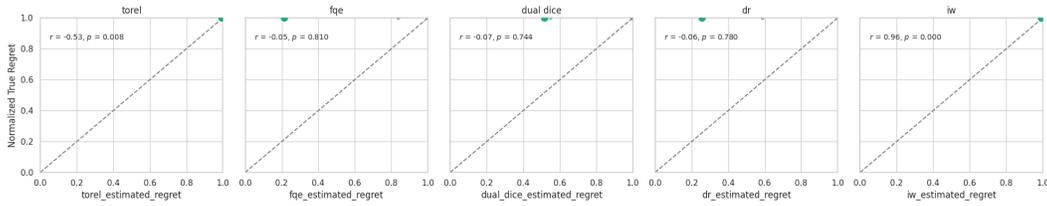


Figure 21: OPE metric scatter-plots: MOREL-trained policies, d4rl-walker-medium-replay-v2.

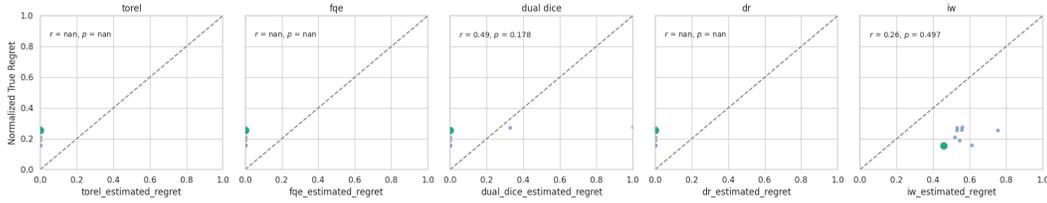


Figure 22: OPE metric scatter-plots: IQL-trained policies, d4rl-pen-expert-v1.

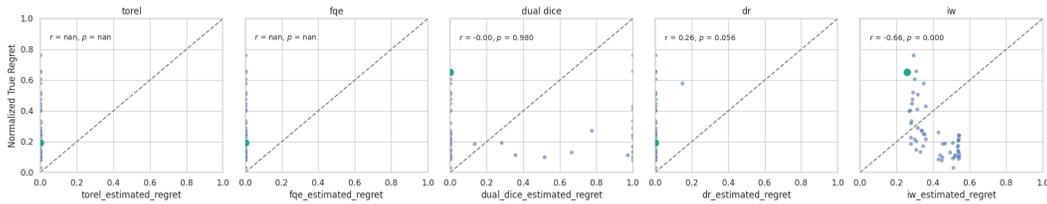


Figure 23: OPE metric scatter-plots: ReBRAC-trained policies, d4rl-pen-expert-v1.

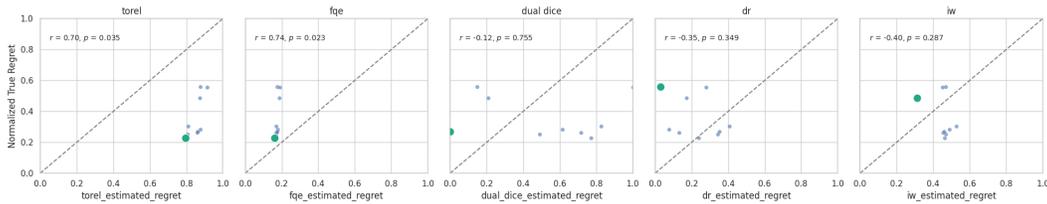


Figure 24: OPE metric scatter-plots: IQL-trained policies, d4rl-hammer-expert-v1.

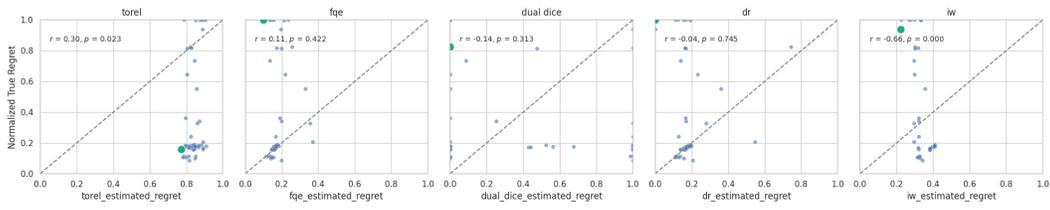


Figure 25: OPE metric scatter-plots: ReBRAC-trained policies, d4rl-hammer-expert-v1.

2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483

E.2 SOReL

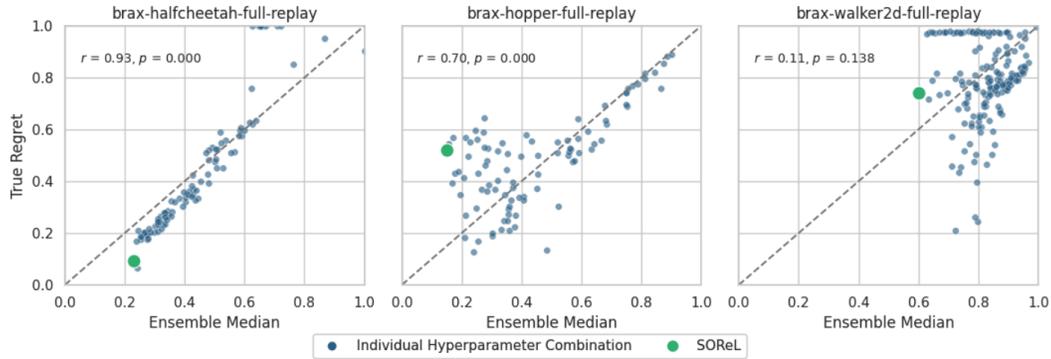


Figure 26: SOReL BAMDP hyperparameter sweeps (tuning set ϕ_{III}) for 200,000 randomly sampled transitions of the brax datasets. The plots correspond to $\phi_{III} \leftarrow \arg \min_{\phi_{III}} \text{RegretMetric}(\phi_I, \phi_{II}, \phi_{III}, \mathcal{D}_N)$ in Algorithm 1. SOReL selects the BAMDP hyperparameters that yield the lowest approximate regret (green). For Walker2d the high approximate regret for all hyperparameter combinations ($R_N > 0.6$) suggests that in Algorithm 1 $R_N > \mathcal{R}_{\text{Deploy}}$ has not been satisfied: the practitioner should change the model or approximate inference method to obtain a lower PIL before re-tuning the BAMDP hyperparameters. Alternatively, the practitioner could consider collecting more data - though the high approximate regret highlights the associated risk. While we report the true regret to validate the approach, in practice only the policy with the lowest approximate regret would be deployed.

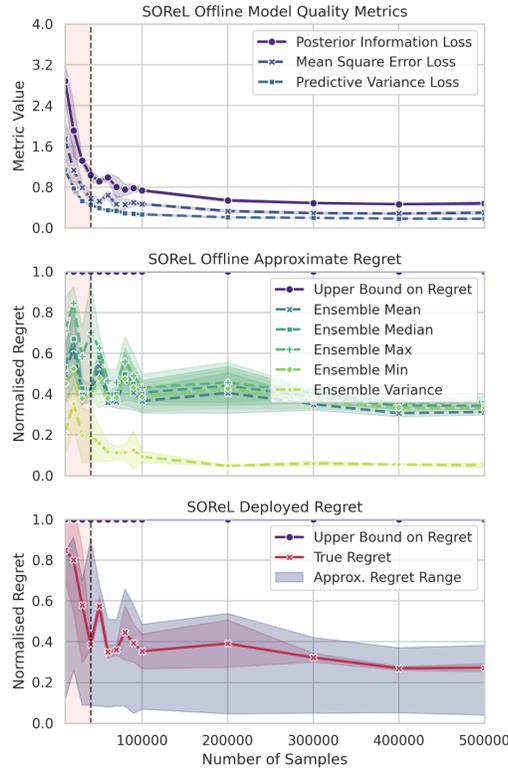


Figure 27: Simplified version of SOReL on brax-halfcheetah-full-replay. The plot showing only the ensemble median as the approximate regret is given in the main body of the paper. Shaded purple shows the approximate regret range across all the metrics (with varying degrees of conservatism). Shaded red indicates where $\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) \not\approx \mathcal{V}(\mathcal{D}_N)$ (for a threshold of 0.25), and hence the approximate regret may be unreliable. Mean and standard deviation given over 3 seeds.

The only environment for which the world model is sufficiently accurate relative to its discount factor, resulting in a non-trivial upper bound, is pendulum-v1 (Fig. 28b). This is due to two factors: (i) the other environments use a lower discount factor (0.998 vs. 0.995), and (ii) learning accurate world models with low MSE is inherently more challenging in high-dimensional settings. We note that this is a supervised-learning problem, and orthogonal to our line of work.

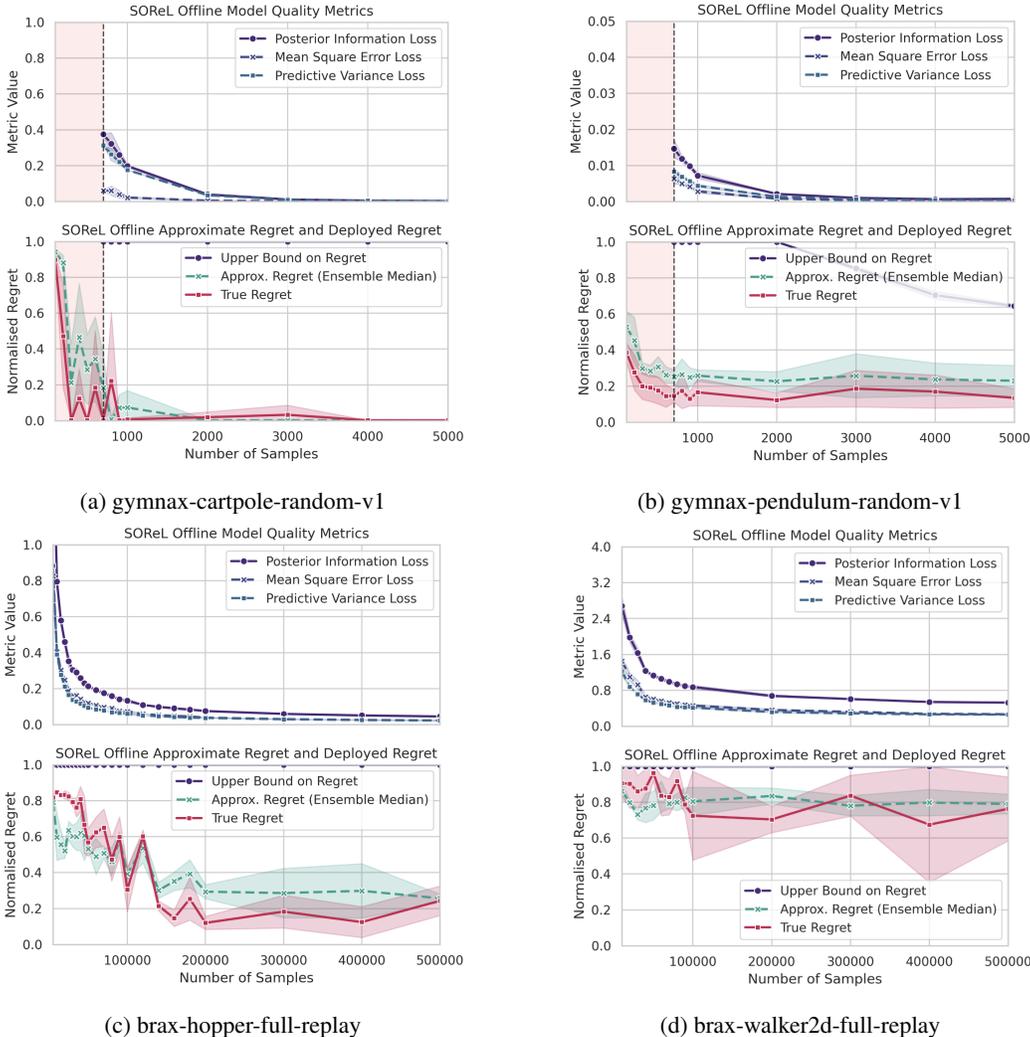


Figure 28: Simplified version of SOReL applied to various tasks. For the gym-nax environments, $N < 700$ is shaded red because the PIL is undefined in this region (validation set $<$ batch size): without being able to ensure $\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) < \approx \mathcal{V}(\mathcal{D}_N)$ the practitioner would have been unable to determine whether to trust the approximate regret. Mean and standard deviation given over 3 seeds.

F IMPLEMENTATION DETAILS

Our implementations of SOReL and TOReL, along with all of the code to reproduce the experiments, are made publicly available with this work.

F.1 DIVERSE FULL-REPLAY DATASETS

As mentioned in Section 6.2, without a model prior, the offline dataset must include transitions from poor, medium and expert regions of performance. For the brax environments, we collect our own full-replay datasets to ensure that this is the case. We arbitrarily choose the hyperparameters, simply requiring that the agent spends sufficient time in all three regions of performance. The training curves

2592
 2593
 2594
 2595
 2596
 2597
 2598
 2599
 2600
 2601
 2602
 2603
 2604
 2605
 2606
 2607
 2608
 2609
 2610
 2611
 2612
 2613
 2614
 2615
 2616
 2617
 2618
 2619
 2620
 2621
 2622
 2623
 2624
 2625
 2626
 2627
 2628
 2629
 2630
 2631
 2632
 2633
 2634
 2635
 2636
 2637
 2638
 2639
 2640
 2641
 2642
 2643
 2644
 2645

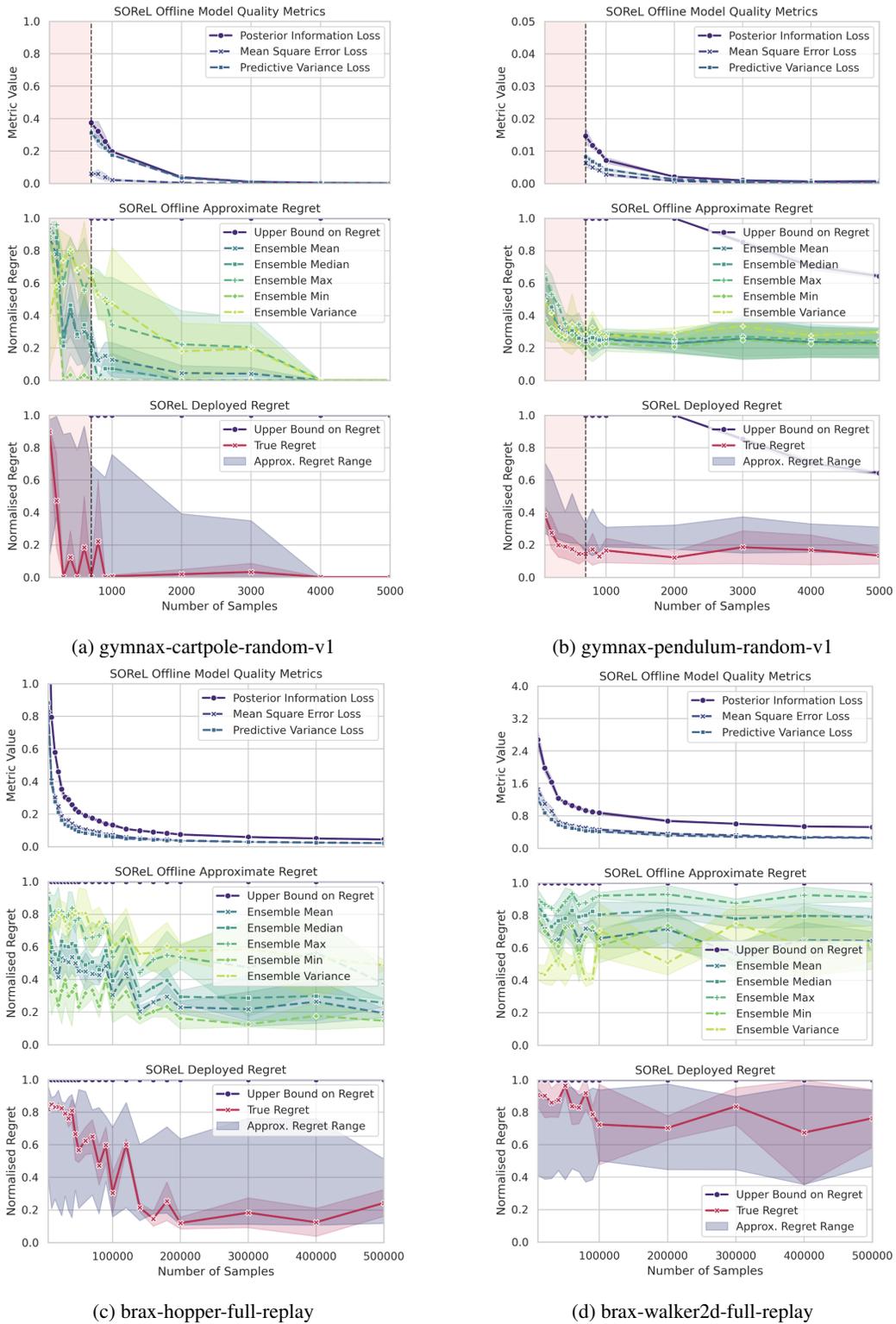


Figure 29: Simplified version of SOReL applied to various tasks. Shaded purple shows the approximate regret range across all the metrics (with varying degrees of conservatism). For the gymnax environments, $N < 700$ is shaded red because the PIL is undefined in this region (validation set $<$ batch size): without being able to ensure $\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) < \approx \mathcal{V}(\mathcal{D}_N)$ the practitioner would have been unable to determine whether to trust the approximate regret. Mean and standard deviation given over 3 seeds.

obtained while collecting the offline datasets are given in Figure 30. The gymnax environments are simple enough that collecting a dataset using an ensemble of randomly initialised policies leads to sufficient coverage across all three regions of performance.

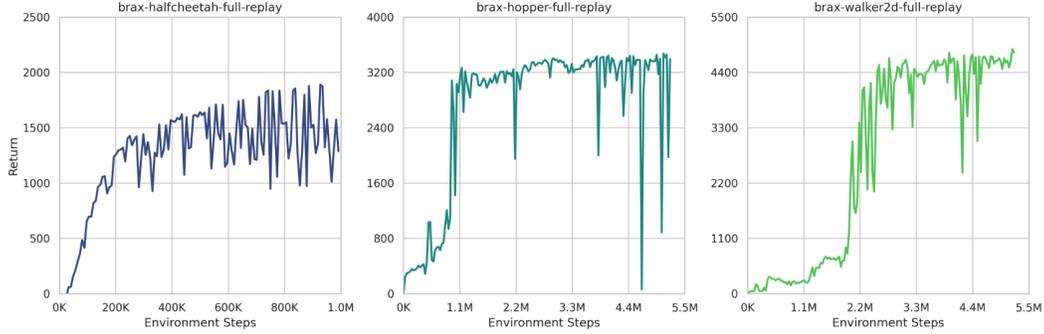


Figure 30: Training curves while collecting the brax full-replay offline datasets. We ensure that the agent spend sufficient time in poor, medium and expert regions of performance such that the offline dataset captures diverse transitions.

F.2 WORLD MODEL AND RP APPROXIMATE INFERENCE

To ensure compatibility with Uniflora implementations (Jackson et al., 2025), our world model is a variation of the Gaussian World Model presented in 4.4, but amended to predict the change in state $\Delta := s' - s$ rather than the absolute next state s' . We also allow the model to characterise its uncertainty with variance functions $\sigma_{r,\theta}^2(s, a)$ and $\sigma_{\Delta,\theta}^2(s, a)$. The Gaussian reward and state transition models then have the form:

$$P_R(s, a, \theta) = \mathcal{N}(r_\theta(s, a), \sigma_{r,\theta}^2(s, a)), \quad P_\Delta(s, a, \theta) = \mathcal{N}(\Delta_\theta(s, a), \sigma_{\Delta,\theta}^2(s, a)),$$

with mean reward function $r_\theta(s, a)$ and mean state transition function $\Delta_\theta(s, a)$, as before. Let $r(s, a, \mathcal{D}_N) := \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [r_\theta(s, a)]$ and $\Delta(s, a, \mathcal{D}_N) := \mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} [\Delta_\theta(s, a)]$ denote the Bayesian mean reward and state transition functions and $r^*(s, a)$ and $\Delta^*(s, a)$ denote the true mean reward and state transition functions. We define the normalised mean squared error between the true and Bayesian mean functions as:

$$\mathcal{E}(\mathcal{D}_N, \mathcal{M}^*) := \mathbb{E}_{(s,a) \sim \rho_\pi^*} \left[\frac{\|r(s, a, \mathcal{D}_N) - r^*(s, a)\|_2^2}{2\sigma_r^2(\mathcal{D}_N)} + \frac{\|\Delta(s, a, \mathcal{D}_N) - \Delta^*(s, a)\|_2^2}{2\sigma_\Delta^2(\mathcal{D}_N)} \right],$$

and the normalised predictive variance using the law of total variance as:

$$\mathcal{V}(\mathcal{D}_N) := \mathbb{E}_{(s,a) \sim \rho_\pi^*} \left[\mathbb{E}_{\theta \sim P_\Theta(\mathcal{D}_N)} \left[\frac{\|r(s, a, \mathcal{D}_N) - r_\theta(s, a)\|_2^2}{2\sigma_r^2(\mathcal{D}_N)} + \|\sigma_{r,\theta}^2(s, a)\|_2^2 \right. \right. \\ \left. \left. + \frac{\|\Delta(s, a, \mathcal{D}_N) - \Delta_\theta(s, a)\|_2^2}{2\sigma_\Delta^2(\mathcal{D}_N)} + \|\sigma_{\Delta,\theta}^2(s, a)\|_2^2 \right] \right],$$

where $\sigma_r^2(\mathcal{D}_N)$ and $\sigma_\Delta^2(\mathcal{D}_N)$ denote the variance of the reward and change in state over the offline dataset.

When rolling out sequences of trajectories on which to train our (Bayes-Optimal) policy, we uniformly sample a model from the ensemble of elite models and then sample the transition from the corresponding Gaussian output distribution.

Our ensemble consists of multilayer perceptrons (MLPs) with ReLU activation, which we train using negative log-likelihood loss derived in Section C.1. Training the models in parallel allows us to simultaneously optimise maximum and minimum (log) variance parameters for each dimension across the model ensemble, which we use to soft-clamp the (log) variances output by the individual models. This prevents any individual model becoming overly confident or too uncertain in one dimension. All models in our ensemble have identical structure, but are initialised differently using LeCun

LeCun et al. (1998) initialisation. The maximum and minimum log-variance terms are initialised at constants. The exact loss function and ensemble dynamics model are the same as the one implemented by Jackson et al. (2025), but we use an Adam optimiser Kingma and Ba (2014) with cosine learning rate schedule rather than constant learning rate. A percentage of the available offline dataset is used as a validation set to calculate the PIL. At the end of training, only a subset of elite models are retained, based on their validation MSE. Although the current implementation uses hard-coded reset and termination conditions during model rollouts, the dynamics model could naturally be extended to learn reset and termination heads. When sampling transitions, we conservatively clip the rewards to remain within the support of the offline dataset distribution.

Hyperparameter	gymnax- cartpole- random-v1	gymnax- pendulum- random-v1	brax- halfcheetah- full-replay	brax- hopper- full-replay	brax- walker2d- full-replay
Num. layers	3	3	3	3	3
Layer Size	200	200	200	200	200
Activation	ReLU	ReLU	ReLU	ReLU	ReLU
Num. Ensemble Models	7	7	7	10	10
Num. Elite Models	5	5	5	8	8
Log Var. Diff. Coeff.	0.01	0.01	0.01	0.01	0.01
Batch Size	64	64	256	256	256
Num. Epochs	400	400	400	400	400
Learning Rate	0.001	0.001	0.001	0.001	0.001
Learning Rate Schedule	cosine	cosine	cosine	cosine	cosine
Final Learning Rate %	10	10	10	10	10
Weight Decay	2.5e-05	2.5e-05	2.5e-05	2.5e-05	2.5e-05
Validation Split	0.1	0.1	0.1	0.1	0.1

Table 6: World model ensemble dynamics hyperparameters.

2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807

F.3 BAMDP SOLVER

We use the RNN-PPO implementation of Lu et al. (2022a), which we amend to be compatible with continuous action spaces. We sweep over the hyperparameters given in Table 7.

Hyperparameter	Value / Sweep Values
Learning rate	[0.0001, 0.0003]
Anneal learning rate	True
Number of environments	[4, 64, 128, 256, 512]
Steps per environment	[32, 64, 128]
Total timesteps	Set to 500,000, 1,000,000 or 50,000,000
Update epochs	[2, 4, 8]
Number of minibatches	[2, 4, 8, 16]
Discount factor (γ)	[0.99, 0.995, 0.998]
GAE lambda	[0.8, 0.9, 0.95]
Clip ϵ	[0.2, 0.3]
Entropy coefficient	[0.000, 0.001, 0.010]
Value function coefficient	0.5
Max gradient norm	[0.5, 1.0]
Layer Size	256
Activation function	tanh
RNN size	Set to 64, 128 or 256
Burn-in Percentage	25

Table 7: RNN-PPO hyperparameters swept over.

Hyperparameter	gymnax- cartpole- random-v1	gymnax- pendulum- random-v1	brax- halfcheetah- full-replay	brax- hopper- full-replay	brax- walker2d- full-replay
Learning rate	0.0003	0.0003	0.0003	0.0003	0.0003
Anneal learning rate	True	True	True	True	True
Number of environments	4	128	512	512	512
Steps per environment	128	64	64	32	64
Total timesteps	500,000	1,000,000	50,000,000	50,000,000	50,000,000
Update epochs	4	8	8	2	4
Number of minibatches	4	16	16	8	8
Gamma	0.99	0.99	0.99	0.998	0.995
GAE lambda	0.95	0.95	0.95	0.8	0.95
Clip ϵ	0.2	0.2	0.2	0.3	0.2
Entropy coefficient	0.01	0.003	0.003	0.001	0.001
Value function coefficient	0.5	0.5	0.5	0.5	0.5
Max gradient norm	0.5	0.5	0.5	1.0	0.5
Layer Size	256	256	256	256	256
Activation function	tanh	tanh	tanh	tanh	tanh
RNN size	64	128	256	256	256
Burn-in Percentage	25	25	25	25	25

Table 8: RNN-PPO hyperparameters for gymnax and brax environments. For computational reasons, we sweep over the hyperparameters of each task once, for a fixed dataset size (1000 datapoints for the gymnax tasks and 200,000 datapoints for brax tasks) and choose the hyperparameters corresponding to the lowest approximate regret, which we then use to train the policy of all other dataset sizes. Ideally, we would sweep over all hyperparameters for each dataset size.

2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861

F.4 ORL IMPLEMENTATIONS

We use Jackson et al. (2025)’s implementations of the ORL algorithms. We use their default hyperparameters and sweep over their suggested hyperparameters, which we summarise in Table 9.

	Hyperparameter	Value / Sweep Values
Generic Optimisation	Discount factor γ	0.99
	Polyak averaging coefficient	0.005
IQL	Learning rate	0.0003
	Batch size	256
	Beta	[0.5, 3.0, 10.0]
	τ (expectile)	[0.5, 0.7, 0.9]
	Advantage clip	100.0
ReBRAC	Learning rate	0.001
	Batch size	1024
	Critic BC coefficient	[0, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.1]
	Actor BC coefficient	[0.0005, 0.001, 0.002, 0.003, 0.03, 0.1, 0.3, 1.0]
	Critic layer norm	true
	Actor layer norm	false
	Observation normalization	false
	Noise clip	0.5
	Policy noise	0.2
	Num critic updates per step	2
MOPO	Learning rate	0.0001
	Batch size	256
	Model retain epochs	5
	Number of critics	10
	Rollout batch size	50000
	Rollout interval	1000
	Rollout length	[1, 3, 5]
	Dataset sample ratio	0.05
	Step penalty coefficient	[1.0, 5.0]
MOReL	Learning rate	0.0001
	Batch size	256
	Model retain epochs	5
	Number of critics	10
	Rollout batch size	50000
	Rollout interval	1000
	Rollout length	5
	Dataset sample ratio	0.01
	Threshold coefficient	[0, 5, 10, 15, 20, 25]
Termination penalty offset	[-30, -50, -100, -200]	

Table 9: Hyperparameters and sweep ranges for IQL, ReBRAC, MOPO, and MOReL.

For the sample efficiency experiments in Fig. 3, we use the default UCB bandit-based hyperparameter-tuning algorithm hyperparameters. We rescale the regret into a score to be maximised, normalised between 0 and 100 ($100 \cdot (1 - \text{regret})$), enabling fair comparison and compatibility with the algorithm.

F.5 REGRET NORMALISATION

We normalise the regret using the (known) minimum and maximum returns (R_{min} and R_{max}) that an online-learnt policy would achieve in the true environment. In the absence of access to these values, we suggest estimating R_{max} or R_{min} as $\frac{r_{min}}{1-\gamma}$ and $\frac{r_{max}}{1-\gamma}$ respectively, where r_{min} and r_{max} denote the 2.5th and 97.5th percentiles of episode rewards in the offline dataset. These thresholds are

suggested to avoid unrealistic assumptions, such as that the best possible policy consistently receives the 100th percentile reward at every time step. Such unrealistic assumptions may significantly distort results, especially in reward distributions with heavy tails where an inflated $R_{max} - R_{min}$ would compress the true and expected regrets.

	gymnax- cart pole- random- v1	gymnax- pend ulum- random- v1	brax- half cheetah- full- replay	brax- hop per- full- replay	brax- walker 2d full- replay	d4rl- halfcheetah medium- expert- v2	d4rl- hopper med ium- v2	d4rl- walker2d medium- replay- v2
r_{min}	0	-13.3	-0.50	0.00	0.00	-0.28	-0.20	0.00
r_{max}	1	-0.2	3.50	3.50	3.50	12.14	3.23	4.59
$P_{2.5}$	1	-13.2	-0.84	1.08	0.01	1.99	1.25	-0.19
$P_{97.5}$	1	-0.16	3.44	4.34	5.65	12.39	4.73	4.92

Table 10: Top half of table: for the gymnax and brax tasks, we define r_{min} and r_{max} using approximate known minimum and maximum returns Jesson et al. (2024), and dividing by the episode length. For the D4RL tasks, we divide the given D4RL minimum and maximum reference scores by the episode length to find r_{min} and r_{max} . Bottom half of table: the suggested normalisation values if expert and random scores were unknown, and determined using the 95th percentile of the offline dataset. Based on the above datasets, for all datasets apart from cartpole-v1, normalisation using the 95th percentile approximation would lead to a more conservative approximate regret. We note that as long as we use the same normalisation constants for both the approximate and true regrets, the absolute value of the normalisation constants is arbitrary. cartpole-v1 is an exception, as the reward is constant for each step (episode returns vary only due to early termination).

We calculate the infinite horizon discounted return from the finite horizon discounted return as follows:

$$R_{inf} = R_{fin} \cdot \left(1 + \frac{\gamma^s}{1 - \gamma^s}\right),$$

where s represents the maximum number of episode steps. The normalised regret is then calculated from the infinite horizon discounted return using:

$$\text{Regret} = \frac{R_{max} - R_{inf}}{R_{max} - R_{min}}.$$

F.6 EXPERIMENT COMPUTE RESOURCES

All of our experiments were run within a week using four L40S GPUs.