

# Bilinear Exponential Family of MDPs: Frequentist Regret Bound with Tractable Exploration & Planning

Anonymous Author(s)

email

Affiliation

Address

## Abstract

We study the problem of episodic reinforcement learning in continuous state-action spaces with unknown rewards and transitions. Specifically, we consider the setting where the rewards and transitions are modeled using parametric bilinear exponential families. We propose an algorithm, BEF-RLSVI, that a) uses penalized maximum likelihood estimators to learn the unknown parameters, b) injects a calibrated Gaussian noise in the parameter of rewards to ensure exploration, and c) leverages linearity of the exponential family with respect to an underlying RKHS to perform tractable planning. We further provide a frequentist regret analysis of BEF-RLSVI that yields an upper bound of  $\tilde{O}(\sqrt{d^3 H^3 K})$ , where  $d$  is the dimension of the parameters,  $H$  is the episode length, and  $K$  is the number of episodes. Our analysis improves the existing bounds for the bilinear exponential family of MDPs by  $\sqrt{H}$  and removes the handcrafted clipping deployed in existing RLSVI-type algorithms. Our regret bound is order-optimal with respect to  $H$  and  $K$ .

**Keywords:** Episodic RL, Planning, Bilinear exponential family

## 1. Introduction

Reinforcement Learning (RL) is a well-studied and popular framework for sequential decision making, where an agent aims to compute a *policy* that allows her to maximize the accumulated reward over a horizon by interacting with an *unknown* environment (Sutton and Barto, 2018).

**Episodic RL.** In this paper, we consider the episodic finite-horizon MDP formulation of RL, in short *Episodic RL* Osband et al. (2013); Azar et al. (2017); Dann et al. (2017). Episodic RL is a tuple  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathbb{P}, r, K, H \rangle$ , where the state (resp. action) space  $\mathcal{S}$  (resp.  $\mathcal{A}$ ) might be continuous. The agent interacts with the environment in  $K$  episodes consisting of  $H$  steps. Episode  $k$  starts by observing state  $s_1^k$ . Then, for  $t = 1, \dots, H$ , the agent draws action  $a_t^k$  from a (possibly time-dependent) policy  $\pi_t(s_t^k)$ , observes the reward  $r(s_t^k, a_t^k) \in [0, 1]$ , and transits to a state  $s_{t+1}^k \sim \mathbb{P}(\cdot | s_t^k, a_t^k)$  following the transition function  $\mathbb{P}$ . The performance of a policy  $\pi$  is measured by the total expected reward  $V_1^\pi$  starting from a state  $s \in \mathcal{S}$ , the value function and the state-action value functions at step  $h \in [H]$  are defined as

$$V_h^\pi(s) \stackrel{\text{def}}{=} \mathbb{E} \left[ \sum_{t=h}^H r(s_t, a_t) \mid s_h = s \right], \quad \text{and} \quad Q_h^\pi(s, a) \stackrel{\text{def}}{=} \mathbb{E} \left[ \sum_{t=h}^H r(s_t, a_t) \mid s_h = s, a_h = a \right].$$

Here, computing the policy leading to maximization of cumulative reward requires the agent to strategically control the actions in order to learn the transition functions and reward functions as precisely as required. This tension between learning the unknown environment and reward maximization is quantified as *regret*: the typical performance measure of an episodic RL algorithm. *Regret* is defined as the difference between the *expected cumulative reward* or *value* collected by the optimal agent that knows the environment and the expected cumulative reward or value obtained by an agent that has to learn about the unknown environment. Formally, the regret over  $K$  episodes is

$$\mathcal{R}(K) \triangleq \sum_{k=1}^K \left( V_1^{\pi^*}(s_1^k) - V_1^{\pi_t}(s_1^k) \right).$$

**Key Challenges.** The first key challenge in episodic RL is to tackle the exploration–exploitation trade-off. This is traditionally addressed with the *optimism principle* that either carefully crafts optimistic upper bounds on the value (or state-action value) functions Azar et al. (2017), or maintains a posterior on the parameters to perform

18 posterior sampling Osband et al. (2013), or perturbs the value (or state-action value) function estimates with calibrated  
 19 noise Osband et al. (2016). Though the first two approaches induce theoretically optimal exploration, they might not  
 20 yield tractable algorithms for large/continuous state-action spaces as they either involve optimization in the optimistic  
 21 set or maintaining a high-dimensional posterior. Thus, *we focus on extending the third approach of Randomized*  
 22 *Least-Square Value Iteration (RLSVI) framework, and inject noise only in rewards to perform tractable exploration.*

23 *The second challenge, which emerges for continuous state-action spaces, is to learn a parametric functional approxima-*  
 24 *tion of either the value function or the rewards and transitions in order to perform planning and exploration. Different*  
 25 *functional representations (or models), such as linear Jin et al. (2020), bilinear Du et al. (2021), and bilinear exponential*  
 26 *families Chowdhury et al. (2021), are studied in literature to develop optimal algorithms for episodic RL with continuous*  
 27 *state-action spaces. Since the linear assumption is restrictive in real-life -where non-linear structures are abundant-,*  
 28 *generalized representations have obtained more attention recently Chowdhury et al. (2021); Li et al. (2021); Du et al.*  
 29 *(2021); Foster et al. (2021). The bilinear exponential family model is of special interest as it is expressive enough to*  
 30 *represent tabular MDPs (discrete state-action), factored MDPs Kearns and Koller (1999), linear MDPs Jin et al. (2020),*  
 31 *linearly controlled dynamical systems (such as Linear Quadratic Regulators Abbasi-Yadkori and Szepesvári (2011)) as*  
 32 *special cases Chowdhury et al. (2021). Thus, in this paper, we study the bilinear exponential family of MDPs, i.e. the*  
 33 *episodic RL setting where the rewards and transition functions can be modelled with bilinear exponential families.*

34 *The third challenge is to perform tractable planning<sup>1</sup> given the perturbation for exploration and the model class.*  
 35 Existing work (Osband and Van Roy, 2014; Chowdhury et al., 2021) assumes an oracle to perform planning and yield  
 36 policies that aren't explicit. The main difficulty in such planning approaches is that dynamic programming requires  
 37 calculating  $\int \mathbb{P}(s' | s, a) V_h(s)$  for all  $(s, a)$  pairs. This is not trivial unless the transition is assumed to be linear and  
 38 decouples  $s'$  from  $(s, a)$ , which is not known to hold except for tabular MDPs. Much ink has been spilled about this  
 39 challenge recently, e.g. Du et al. (2019) asks when misspecified linear representations are enough for a polynomial  
 40 sample complexity in several settings. Shariff and Szepesvári (2020); Lattimore et al. (2020); Van Roy and Dong (2019)  
 41 provide positive answers for specific linear settings. In this paper, *we aim to address this issue by designing a tractable*  
 42 *planner for the bilinear exponential family representation.*

43 In this paper, we aim to address the following question that encompasses the three challenges:

44 Can we design an algorithm that performs **tractable exploration** and **planning** for *bilinear exponential family of*  
 45 *MDPs* yielding a **near-optimal frequentist regret bound**?

46 **Our Contributions.** Our contributions to this question are three-fold.

47 1. *Formalism:* We assume that rewards and transitions are unknown, whereas existing efforts on the bilinear exponential  
 48 family of MDPs assume knowledge of rewards. This makes the addressed problem harder, practical, and more general.  
 49 We also observe that though the transition model can represent non-linear dynamics, it implies a linear behavior (see  
 50 Section 2) in a Reproducible Kernel Hilbert Space (RKHS). This observation contributes to the tractability of planning.

51 2. *Algorithm:* We propose an algorithm BEF-RLSVI that extends the RLSVI framework to bilinear exponential  
 52 families (see Section 3). BEF-RLSVI a) injects calibrated Gaussian noise in the rewards to perform exploration, b)  
 53 leverages the linearity of the transition model with respect to an underlying RKHS to perform tractable planning and c)  
 54 uses penalized maximum likelihood estimators to learn the parameters corresponding to rewards and transitions (see  
 55 Section 4). To the best of our knowledge, *BEF-RLSVI is the first algorithm for the bilinear exponential family of*  
 56 *MDPs with tractable exploration and planning under unknown rewards and transitions.*

57 3. *Analysis:* We carefully develop an analysis of BEF-RLSVI that yields  $\tilde{O}(\sqrt{d^3 H^3 K})$  regret which improves the  
 58 existing regret bound for bilinear exponential family of MDPs with known reward by a factor of  $\sqrt{H}$  (Section 3.2). Our  
 59 analysis (Section 5) builds on existing analyses of RLSVI-type algorithms Osband et al. (2016), but contrary to them,  
 60 we remove the need to handcraft a clipping of the value functions Zanette et al. (2020). We also do not need to *assume*  
 61 *anti-concentration bounds* as we can explicitly control it by the injected noise. This was not done previously except for  
 62 the linear MDPs. We illustrate this comparison in Table 1. We highlight three technical tools that we used to improve the  
 63 previous analyses: 1) Using transportation inequalities instead of the simulation lemma reduces a  $\sqrt{H}$  factor compared  
 64 to Ren et al. (2021), 2) Leveraging the observation that true value functions are bounded enables using an improved  
 65 elliptical lemma (compared to Chowdhury et al. (2021)), and 3) Noticing that the norm of features can only be large for a  
 66 finite amount of time allows us to forgo clipping and reduce a  $\sqrt{d}$  factor from the regret compared to Zanette et al. (2020).

1. By tractable planning, we mean having a planner with (pseudo-)polynomial complexity in the problem parameters, i.e. dimension of parameters, dimension of features, horizon, and number of episodes.

Table 1: A comparison of RL Algorithms for continuous state-actions with functional representations.

Algo	Regret	Tractable exploration	Tractable planning	Free of clipping	Model, assumptions
Thompson sampling Ren et al. (2021)	$\sqrt{d^2 H^3 K}$ (Bayesian)	✗	✓	N.A	Gaussian $\mathbb{P}$ Known rewards
LSVI-PHE Ishfaq et al. (2021)	$\sqrt{d^3 H^4 K}$ (Freq.)	✓	✓	✗	Generalized $V$ approx Tabular, anti-concentration
OPT-RLSVI Zanette et al. (2020)	$\sqrt{d^4 H^5 K}$ (Freq.)	✓	✓	✗	Linear $V$
EXP-UCRL Chowdhury et al. (2021)	$\sqrt{d^2 H^4 K}$ (Freq.)	✗	✗	N.A	Bilinear Exp family known rewards
BEF-RLSVI This work	$\sqrt{d^3 H^3 K}$ (Freq.)	✓	✓	✓	Bilinear Exp family

## 67 2. Bilinear exponential family of MDPs

68 In this section, we introduce the bilinear exponential family model coined in Chowdhury et al. (2021) and extend it to  
69 parametric rewards. Then, we state a novel observation about linearity of this representation.

**Bilinear exponential family model.** We consider both transition and reward kernels to be unknown and modeled with bilinear exponential families. Specifically,

$$\mathbb{P}(\tilde{s} | s, a) = \exp(\psi(\tilde{s})^\top M_{\theta^p} \varphi(s, a) - Z_{s,a}^p(\theta^p)), \quad (1)$$

$$\mathbb{P}(r | s, a) = \exp(r B^\top M_{\theta^r} \varphi(s, a) - Z_{s,a}^r(\theta^r)), \quad (2)$$

where  $\varphi \in (\mathbb{R}_+^q)^{\mathcal{S} \times \mathcal{A}}$  and  $\psi \in (\mathbb{R}_+^p)^{\mathcal{S}}$  are known functions, and  $B \in \mathbb{R}^p$  is a known scaling vector. The unknown reward and transition parameters are  $\theta^p, \theta^r \in \mathbb{R}^d$ .  $M_{\theta^p} \stackrel{\text{def}}{=} \sum_{i=1}^d \theta_i A_i$ , where  $A_i$  is a known  $p \times q$  matrix for each  $i$ . Finally,  $Z$  denotes the log partition function:  $Z_{s,a}^p(\theta^p) \stackrel{\text{def}}{=} \log \int_{\mathcal{S}} \exp(\psi(\tilde{s})^\top M_{\theta^p} \varphi(s, a)) d\tilde{s}$ , and  $Z^r$  is defined similarly. We denote  $V_{\theta^p, \theta^r, h}^\pi$  and  $Q_{\theta^p, \theta^r, h}^\pi$ , the value and state-action value function respectively, for policy  $\pi$  in the MDP parameterized by  $(\theta^p, \theta^r)$  at time  $h$ . A policy  $\pi^*$  is *optimal* if for all  $s \in \mathcal{S}$ ,  $V_{\theta^p, h}^{\pi^*}(s) = \max_{\pi \in \Pi} V_{\theta^p, h}^\pi(s)$ . A learning algorithm minimizes the (pseudo-)regret defined as:

$$\mathcal{R}(K) \triangleq \sum_{k=1}^K \left( V_{\theta^p, 1}^{\pi^*}(s_1^k) - V_{\theta^p, 1}^{\pi^k}(s_1^k) \right). \quad (3)$$

**Linearity of transitions.** Now, we state an observation about the bilinear exponential family and discuss how it helps with the challenge of planning in episodic RL. Specifically, the popular assumption of linearity of the transition kernel is a direct consequence of our model. Indeed,

$$2\psi(s')^\top M_{\theta^p} \varphi(s, a) = -\|(\psi(s') - M_{\theta^p} \varphi(s, a))\|^2 + \|\psi(s')\|^2 + \|M_{\theta^p} \varphi(s, a)\|^2.$$

The quadratic term resembles the Radial Basis Function (RBF). More precisely, for an RBF kernel with covariance  $\Sigma = I_p$  and  $k(x, y) \stackrel{\text{def}}{=} \exp(-\|x - y\|^2/2)$ , we find

$$\mathbb{P}(s' | s, a) = \langle \phi^p(s, a), \mu^p(s') \rangle_{\mathcal{H}} \quad (4)$$

70 where  $\mathcal{H}$  is the RKHS associated with the kernel,  $\mu^p(s') = (2\pi)^{-p/2} k(\psi(s'), \cdot) \exp(\|\psi(s')\|^2/2)$ , and  $\phi^p(s, a) =$   
71  $k(M_{\theta^p}^\top \varphi(s, a), \cdot) \exp(\|M_{\theta^p} \varphi(s, a)\|^2/2 - Z_{s,a}^p(\theta^p))$ . Equation (4) shows that  $s'$  is decoupled from  $(s, a)$ .

72 **Remark 1** Ren et al. (2021) is the only other work providing an example of linear transitions in continuous state-action  
73 spaces. It considers Gaussian transitions with an unknown mean ( $f^*(s, a)$ ) and known variance ( $\sigma^2$ ). Actually, linear  
74  $f^*$  is a special case of our model with  $\psi(s') = (s', \|s'\|^2)$  and  $M_{\theta^p} \varphi(s, a) = (f_\theta(s, a)/\sigma^2, -1/\sigma^2)$ .

**Importance of linearity.** To understand the planning challenge in RL, recall the Bellman equation:

$$Q_h^\pi(s, a) = r(s, a) + \int_{\tilde{s} \in \mathcal{S}} P(s' | s, a) V_{h+1}^\pi(\tilde{s}) d\tilde{s},$$

we must approximate the integral at the R.H.S for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . For a tabular MDP, we need to evaluate  $(Q_h^\pi)_{h \in [H]}$ , i.e. approximate  $|\mathcal{S}| \times |\mathcal{A}| \times H$  integrals per episode, which is very expensive. However, if Equation (4) holds, then

$$Q_{\hat{\theta}, h}^\pi(s, a) = r(s, a) + \left\langle \phi^\mathbb{P}(s, a), \int_{\mathcal{S}} \mu^\mathbb{P}(\tilde{s}) V_{\hat{\theta}, h+1}^\pi(\tilde{s}) d\tilde{s} \right\rangle. \quad (5)$$

75 When  $\phi^\mathbb{P}, \mu^\mathbb{P} \in \mathbb{R}^\tau$ , we can obtain  $Q_{\hat{\theta}^\mathbb{P}, \hat{\theta}^\mathbb{P}, h}$  by computing  $\tau$  integrals. For our model, although  $\phi^\mathbb{P}$  and  $\mu^\mathbb{P}$  are infinite  
76 dimensional, we show in Section 4 (§ planning) that the planning is still computationally tractable.

### 77 3. BEF-RLSVI: algorithm design and frequentist regret bound

78 In this section, we formally introduce the Bilinear Exponential Family Randomized Least-Squares Value Iteration  
79 (BEF-RLSVI) algorithm. Then, we present a high probability upper-bound on its regret.

#### 80 3.1 BEF-RLSVI: algorithm design

81 BEF-RLSVI is based on RLSVI (Osband et al., 2016) with the distinction that we only perturb the reward parameters  
82 and not all the parameters of the value function. RLSVI algorithms are reminiscent of Thompson Sampling, yet more  
83 tractable with better control over the probability to be optimistic.

---

#### Algorithm 1 BEF-RLSVI

---

- 1: **Input:** failure rate  $\delta$ , constants  $\alpha^\mathbb{P}, \eta$  and  $(x_k)_{k \in [K]} \in \mathbb{R}^+$
  - 2: **for** episode  $k = 1, 2, \dots$  **do**
  - 3:   Observe initial state  $s_1^k$
  - 4:   Sample noise  $\xi_k \sim \mathcal{N}(0, x_k(\bar{G}_k^\mathbb{P})^{-1})$  such that
 
$$\bar{G}_k^\mathbb{P} = \frac{\eta}{\alpha^\mathbb{P}} \mathbb{A} + \sum_{\tau=1}^{k-1} \sum_{h=1}^H (\varphi(s_h^\tau, a_h^\tau)^\top A_i^\top A_j \varphi(s_h^\tau, a_h^\tau))_{i,j \in [d]}$$
  - 5:   Perturb reward parameter:  $\tilde{\theta}^x(k) = \hat{\theta}^x(k) + \xi_k$
  - 6:   Compute  $(Q_{\hat{\theta}^\mathbb{P}, \tilde{\theta}^x, h}^k)_{h \in [H]}$  via Bellman-backtracking, see Algorithm 2
  - 7:   **for**  $h = 1, \dots, H$  **do**
  - 8:     Pull action  $a_h^k = \arg \max_a Q_{\hat{\theta}^\mathbb{P}, \tilde{\theta}^x, h}(s_h^k, a)$
  - 9:     Observe reward  $r(s_h^k, a_h^k)$  and state  $s_{h+1}^k$ .
  - 10:   **end for**
  - 11:   Update the penalized ML estimators  $\hat{\theta}^\mathbb{P}(k), \hat{\theta}^x(k)$ , see Equation (6) and Equation (8)
  - 12: **end for**
- 

84 Algorithm 1 performs exploration by a Gaussian perturbation of the reward parameter (Line 4). Unlike optimistic  
85 approaches, this method is explicit and more efficient since it does not involve a high-dimensional optimization.

---

#### Algorithm 2 Bellman Backtracking

---

- 1: **Input** Parameters  $\hat{\theta}^\mathbb{P}, \tilde{\theta}^x$ , initialize  $\tilde{\theta} = (\tilde{\theta}^x, \hat{\theta}^\mathbb{P})$  and  $\forall s, V_{H+1}(s) = 0$
  - 2: **for** steps  $h = H - 1, H - 2, \dots, 0$  **do**
  - 3:   Calculate  $Q_{\tilde{\theta}, h}(s, a) = \mathbb{E}_{s,a}^{\tilde{\theta}^x} [r] + \langle \phi^\mathbb{P}(s, a), \int V_{\tilde{\theta}, h+1}(s') \mu^\mathbb{P}(s') ds' \rangle_{\mathcal{H}}$ .
  - 4: **end for**
- 

86 We can approximate Line 3 of Algorithm 2 with  $\mathcal{O}(pH^3K \log(HK))$  complexity without compromising regret  
87 guarantees (cf. § planning, Section 4). Therefore, Algorithm 2 provides tractable planning.

88 **3.2 BEF-RLSVI: regret upper-bound**

89 We state the standard smoothness assumptions on the model (Chowdhury et al., 2021; Jun et al., 2017; Lu et al., 2021).

**Assumption 2** *There exist constants  $\alpha^p, \alpha^r, \beta^p, \beta^r > 0$ , such that the representation model satisfies:*

$$\begin{aligned} \forall (s, a) \in \mathcal{S} \times \mathcal{A}, \forall \theta, x \in \mathbb{R}^d \quad \alpha^p \leq x^\top C_{s,a}^\theta [\psi] x \leq \beta^p \\ \forall (s, a) \in \mathcal{S} \times \mathcal{A}, \forall \theta, x \in \mathbb{R}^d \quad \alpha^r \leq \mathbb{V}\text{ar}_{s,a}^\theta(r) x^\top B^\top B x \leq \beta^r \end{aligned}$$

 90 where  $\mathbb{C}_{s,a}^\theta [\psi(s')] \triangleq \mathbb{E}_{s' \sim \mathbb{P}_\theta | s,a} [\psi(s') \psi(s')^\top] - \mathbb{E}_{s' \sim \mathbb{P}_\theta | s,a} [\psi(s')] \mathbb{E}_{s' \sim \mathbb{P}_\theta | s,a} [\psi(s')^\top]$  and  $\mathbb{V}\text{ar}_{s,a}^\theta(r)$  is the vari-  
 91 ance of the reward under  $\theta$  defined by  $\mathbb{V}\text{ar}_{s,a}^\theta(r) \triangleq \left( \mathbb{E}_{s,a}^\theta [r^2] - \mathbb{E}_{s,a}^\theta [r]^2 \right)$ .

 92 A closer look at the derivatives of the model (see Appendix D.3) tells us that previous inequalities directly imply a  
 93 control over the eigenvalues of the Hessian matrices of the log-normalizers.

94 We now state our main result, the regret upper-bound of BEF-RLSVI.

 95 **Theorem 3 (Regret bound)** *Let  $\mathbb{A} \triangleq (\text{tr}(A_i A_j^\top))_{i,j \in [d]}$  and  $G_{s,a} \triangleq (\varphi(s, a)^\top A_i^\top A_j \varphi(s, a))_{i,j \in [d]}$ . Under Assump-*  
 96 *tion 2 and further considering that*

- 97 1.
- $\max\{\|\theta^r\|_{\mathbb{A}}, \|\theta^p\|_{\mathbb{A}}\} \leq B_{\mathbb{A}}, \|\mathbb{A}^{-1} G_{s,a}\| \leq B_{\varphi, \mathbb{A}}$
- and
- $\mathbb{E}_{\theta^r}[r(s, a)] \in [0, 1]$
- for all
- $(s, a)$
- .
- 
- 98 2. noise
- $\xi_k \sim \mathcal{N}(0, x_k (\bar{G}_k^p)^{-1})$
- satisfies
- $x_k \geq \left( H \sqrt{\frac{\beta^p \beta^p (K, \delta)}{\alpha^p \alpha^r}} + \frac{\sqrt{\beta^r \beta^r (K, \delta) \min\{1, \frac{\alpha^p}{\alpha^r}\}}}{2\alpha^r} \right)^2 \propto dH^2$
- ,

 then for all  $\delta \in (0, 1]$ , with probability at least  $1 - 7\delta$ ,

$$\begin{aligned} \mathcal{R}(K) \leq \sqrt{KH} \left[ \underbrace{2H \left( \sqrt{\frac{2\beta^p}{\alpha^p}} \beta^p (K, \delta) \gamma_K^p + (1 + \sqrt{\gamma_K^r}) \sqrt{\log(1/\delta^2)} \right)}_{\text{Transition concentration} \approx dH} + \underbrace{\beta^r \sqrt{\frac{\beta^r (n, \delta) \gamma_K^r}{2\alpha^r}}}_{\text{Reward concentration} \approx d} \right. \\ \left. + \underbrace{c\beta^r \sqrt{x_K d \gamma_K^r \log(dK/\delta)} + \frac{\beta^r \sqrt{x_K d \gamma_K^r \log(e/\delta^2)}}{\Phi(-1)} (1 + \sqrt{\log(d/\delta)})}_{\text{Noise concentration} \approx d^{3/2} H} \right] \\ + \sqrt{H \gamma_K^r} \left[ \underbrace{\beta^r C_d \left( \sqrt{\frac{\beta^r (K, \delta)}{2\alpha^r}} + c \sqrt{x_K d \log(dK/\delta)} \right)}_{\text{Estimation error for no clipping} \approx dH} + \underbrace{\frac{\beta^r d \sqrt{x_K}}{\Phi(-1)} (1 + \sqrt{\log(d/\delta)}) \sqrt{C_d \left( 1 + \frac{\alpha^r B_{\varphi, \mathbb{A}} H}{\eta} \right)}}_{\text{Learning error for no clipping} \approx (dH)^{3/2}} \right], \end{aligned}$$

 99 where for  $i \in [p, r]$ ,  $\beta^i(K, \delta) \triangleq \frac{\eta}{2} B_{\mathbb{A}}^2 + \gamma_K^i + \log(1/\delta)$ , and  $\gamma_K^i \triangleq d \log(1 + \frac{\beta^i}{\eta} B_{\varphi, \mathbb{A}} H K)$ . Also,  $C_d \triangleq$   
 100  $\frac{3d}{\log(2)} \log \left( 1 + \frac{\alpha^r \|\mathbb{A}\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)} \right)$ ,  $\Phi$  is the Gaussian CDF, and  $c$  is a universal constant.

 101 Theorem 3 entails a regret  $\mathcal{R}(K) = \mathcal{O}(\sqrt{d^3 H^3 K})$  for BEF-RLSVI, where  $d$  is the number of parameters of the  
 102 bilinear exponential family model,  $K$  is the number of episodes, and  $H$  is the horizon of an episode.

 103 *Comparison with other bounds.* The closest work to ours is Chowdhury et al. (2021) as it considers the same model  
 104 for transitions but with known rewards. They propose a UCRL-type and PSRL-type algorithm, which achieve a regret  
 105 of order  $\tilde{\mathcal{O}}(\sqrt{d^2 H^4 K})$ . There are two notable algorithmic differences with our work. First, they do exploration using  
 106 intractable-optimistic upper bounds or high-dimensional posteriors, while we do it with explicit perturbation. The second  
 107 difference is in planning. While they assume access to a planning oracle, we do it explicitly with pseudo-polynomial  
 108 complexity (Section 4). Moreover, we improve the regret bound by a  $\sqrt{H}$  factor thanks to an improved analysis, (cf.  
 109 Lemma 24). But similar to all RLSVI-type algorithms, we pick up an extra  $\sqrt{d}$  (cf. (Abeille and Lazaric, 2017)).

110 [Zanette et al. \(2020\)](#) adapts RLSVI for continuous state-action spaces. Assuming low-rank models of transitions and  
 111 rewards, it shows a regret bound  $R(K) = \tilde{O}(\sqrt{d^4 H^5 K})$ , which is larger than ours by  $O(\sqrt{dH^2})$ . In algorithm design,  
 112 we improve on their work by removing the need to carefully clip the value function. Analytically, our model allows us  
 113 to use transportation inequalities (*cf.* Lemma 19) instead of the simulation lemma, which saves us a  $\sqrt{H}$  factor.

114 [Ren et al. \(2021\)](#) considers Gaussian transitions, i.e.  $s' = f^*(s, a) + \epsilon$  such that  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ . This is a particular case  
 115 of our model. They propose to use Thompson Sampling, and have the merit of being the first to have observed linearity  
 116 of the value function from this transition structure. But they do not connect it to the finite dimensional approximation  
 117 of [Rahimi and Recht \(2007\)](#) unlike us (Section 4). Finally, they show a Bayesian regret bound of  $O(\sqrt{d^2 H^3 K})$ . This  
 118 notion of regret is weaker than frequentist regret, hence this result is not directly comparable with Theorem 3.

119 *Tightness of regret bound.* A lower bound for episodic RL with continuous state-action spaces is still missing. However,  
 120 for tabular RL, ([Domingues et al., 2021](#)) proves a lower bound of order  $\Omega(\sqrt{H^3 S A K})$ . To represent a tabular MDP in  
 121 our model, we need  $d = S^2 \times A$  parameters (Section 4.3, ([Chowdhury et al., 2021](#))). In this case, our bound becomes  
 122  $R(K) = O(\sqrt{(S^2 A)^3 H^3 K})$ , which is clearly not tight in  $S$  and  $A$ . This is understandable due to the relative generality  
 123 of our setting. We are however positively surprised that **our bound is tight in terms of its dependence on  $H$  and  $K$ .**

#### 124 4. Algorithm design: building blocks of BEF-RLSVI

125 We present necessary details about BEF-RLSVI and discuss the key algorithm design techniques.

**Estimation of parameters.** We estimate transitions and rewards from observations similar to EXP-UCRL [Chowdhury et al. \(2021\)](#), *i.e.* by using a penalized maximum likelihood estimator

$$\hat{\theta}^p(k) \in \arg \min_{\theta \in \mathbb{R}^d} \sum_{t=1}^k \sum_{h=1}^H -\log \mathbb{P}_\theta (s_{h+1}^t | s_h^t, a_h^t) + \eta \text{pen}(\theta).$$

Here,  $\text{pen}(\theta)$  is a trace-norm penalty:  $\text{pen}(\theta) = \frac{1}{2} \|\theta\|_{\mathbb{A}}$  and  $\mathbb{A} = (\text{tr}(A_i A_j^\top))_{i,j}$ . By properties of the exponential family, the penalized maximum likelihood estimator verifies, for all  $i \leq d$ :

$$\sum_{t=1}^k \sum_{h=1}^H \left( \psi(s_{h+1}^t) - \mathbb{E}_{s_h^t, a_h^t}^{\hat{\theta}_k^p} [\psi(s')] \right)^\top A_i \varphi(s_h^t, a_h^t) = \eta \nabla_i \text{pen}(\hat{\theta}_k^p). \quad (6)$$

Equation (6) can be solved in closed form for simple distributions, like Gaussian, but it can involve integral approximations for other distribution. We estimate the parameter for reward, *i.e.*  $\theta_r$ , similarly

$$\hat{\theta}^r(k) \in \arg \min_{\theta \in \mathbb{R}^d} \sum_{t=1}^k \sum_{h=1}^H -\log \mathbb{P}_\theta (r_t | s_h^t, a_h^t) + \eta \text{pen}(\theta), \quad (7)$$

$$\implies \sum_{t=1}^k \sum_{h=1}^H \left( r_t - \mathbb{E}_{s_h^t, a_h^t}^{\hat{\theta}_k^r} [r] \right) B^\top A_i \varphi(s_h^t, a_h^t) = \eta \nabla_i \text{pen}(\hat{\theta}_k^r) \quad \forall i \in [d]. \quad (8)$$

**Exploration.** A significant challenge in RL is handling exploration in continuous spaces. The majority of the literature is split between intractable, upper confidence bound-style optimism or Thompson sampling algorithms with high-dimensional posterior and guarantees only in terms of Bayesian regret. In BEF-RLSVI, we adopt the approach of reward perturbation motivated by the RLSVI-framework [Zanette et al. \(2020\)](#); [Osband et al. \(2016\)](#). We show that perturbing the reward estimation can guarantee optimism with a constant probability, *i.e.* there exists  $\nu \in (0, 1]$  such that for all  $k \in [K]$  and  $s_1^k \in \mathcal{S}$ ,

$$\mathbb{P} \left( \tilde{V}_1(s_1^k) - V_1^*(s_1^k) \geq 0 \right) \geq \nu.$$

126 [Zanette et al. \(2020\)](#) proves that this suffices to bound the learning error. However, their method clashes with not  
 127 clipping the value function, as it modifies the probability of optimism. Thus, [Zanette et al. \(2020\)](#) proposes an involved  
 128 clipping procedure to handle the issue of unstable values. Instead, by careful geometric analysis (*cf.* Lemma 26),  
 129 we bound the occurrences of the unstable values, and in turn, upper bound the regret without clipping. Note that

130 unlike (Ishfaq et al., 2021), BEF-RLSVI does not guarantee that the estimated value function is optimistic but still is  
 131 able to control the learning error (cf. Section 5).

132 **Planning.** Recall that with our model assumptions, we can write the state-action value function linearly (Equation (5)).  
 133 Using BEF-RLSVI, we have at step  $h$ :

$$Q_{\hat{\theta}^v, \hat{\theta}^r, h}^\pi(s, a) = \mathbb{E}_{\tilde{\theta}^r} [r(s, a)] + \left\langle \phi^p(s, a), \int_{\mathcal{S}} \mu^p(\tilde{s}) V_{\hat{\theta}^v, \hat{\theta}^r, h+1}^\pi(\tilde{s}) d\tilde{s} \right\rangle.$$

134 Then, we select the best action greedily using dynamic programming to compute  $Q_h(s, a)$ . Although our model  
 135 yields infinite dimensional  $\phi^p$  and  $\psi^p$ , approximating them (cf. next paragraph) with linear features of dimension  
 136  $\mathcal{O}(pH^2K \log(HK))$  is possible without increasing the regret. Thus, the planning is done in  $\mathcal{O}(pH^3K \log(HK))$ ,  
 137 which is pseudo-polynomial in  $p, H$  and  $K$ , i.e. tractable.

138 For details about the finite-dimensional approximation of our transition kernel, refer to Appendix E. Now, we highlight  
 139 the schematic of a finite-dimensional approximation of  $\phi^p$  and  $\psi^p$ . We proceed in three steps. **1)** We have with high  
 140 probability  $\mathbb{S}(V_{\hat{\theta}^v, \hat{\theta}^r, h}) \leq dH^{3/2}$  (Section 5). **2)** If we have a uniform  $\epsilon$ -approximation of  $\mathbb{P}_{\theta^v}$ , we show that using it  
 141 incurs at most an extra  $\mathcal{O}(\epsilon dH^{5/2}K)$  regret. **3)** Finally, following Rahimi and Recht (2007), we approximate uniformly  
 142 the shift invariant kernels, here the RBF in Equation (4), within  $\epsilon$  error and with features of dimensions  $\mathcal{O}(p\epsilon^{-2} \log \frac{1}{\epsilon^2})$ ,  
 143 where  $p$  is dimension of  $\psi$ . Associating these three elements and choosing  $\epsilon = 1/\sqrt{(H^2K)}$ , we establish our claim.

## 144 5. Theoretical analysis: proof outline

To convey the novelties in our analysis, we provide a proof sketch for Theorem 3. We start by decomposing the regret  
 into an estimation loss and a learning error, as given below

$$R(K) = \sum_{k=1}^K (V_{\hat{\theta}^v, \theta^r, 1}^* - V_{\hat{\theta}^v, \theta^r, 1}^{\pi_k})(s_{1k}) = \sum_{k=1}^K \underbrace{(V_{\hat{\theta}^v, \theta^r, 1}^* - V_{\hat{\theta}^v, \hat{\theta}^r, 1}^{\pi_k})}_{\text{learning}} + \underbrace{(V_{\hat{\theta}^v, \hat{\theta}^r, 1}^{\pi_k} - V_{\hat{\theta}^v, \theta^r, 1}^{\pi_k})}_{\text{Estimation}}(s_{1k}). \quad (9)$$

145 For the **estimation error**, we use smoothness arguments with concentrations of parameters up to some novelties.  
 146 Regarding the **learning error**, we show that the injected noise ensures a constant probability of anti-concentration.  
 147 Applying Assumption 2 and Lemma 24 leads to the upper-bound.

### 148 5.1 Bounding the estimation error

We further decompose the estimation error into the errors in estimating transitions and rewards.

$$V_{\hat{\theta}^v, \hat{\theta}^r}^\pi(s_{1k}) - V_{\hat{\theta}^v, \theta^r}^\pi(s_{1k}) = \underbrace{V_{\hat{\theta}^v, \theta^r}^\pi(s_{1k}) - V_{\hat{\theta}^v, \theta^r}^\pi(s_{1k})}_{\text{transition estimation}} + \underbrace{V_{\hat{\theta}^v, \hat{\theta}^r}^\pi(s_{1k}) - V_{\hat{\theta}^v, \theta^r}^\pi(s_{1k})}_{\text{reward estimation}} \quad (10)$$

149 **Transition estimation** Since the reward parameter is exact, the value function’s span is  $\leq H$ . Then, using the  
 150 transportation of Lemma 19 we obtain the bound  $H \sum_{h=1}^H \sqrt{2 \text{KL}_{s_{hk}, a_{hk}}(\theta^v, \hat{\theta}^v)}$ . We notice that since the reward  
 151 parameter is exact, the bound is actually  $H \min\{1, \sum_{h=1}^H \sqrt{2 \text{KL}_{s_{hk}, a_{hk}}(\theta^v, \hat{\theta}^v)}\}$ . Using Lemma 24 under Assumption 2,  
 152 we win a  $\sqrt{H}$  factor compared to the analysis of Chowdhury and Gopalan (2019).

153 **Reward estimation** Previous work uses clipping to help control this error, but in this case it can hinder the opti-  
 154 mism probability by biasing the noise. Zanette et al. (2020) proposes an involved clipping depending on the norms  
 155  $\|(A_i \varphi(s_h^k, a_h^k))_{i \in [d]}\|_{(\tilde{\mathcal{G}}_k^r)^{-1}}$ , which is somewhat delicate to analyze and deploy. We remedy the situation acting solely  
 156 in the proof. First let’s define what we call the set of “bad rounds”:  $\left\{k \in [K], \exists h : \|(A_i \varphi(s_h^k, a_h^k))_{i \in [d]}\|_{(\tilde{\mathcal{G}}_k^r)^{-1}} \geq 1\right\}$ ,  
 157 these rounds are why clipping is necessary. Thanks to Lemma 26, we know that the number of such rounds is at most  
 158  $\mathcal{O}(d)$ . Surprisingly, it depends neither on  $H$  nor on  $K$ . We show that the “bad rounds” incur at most  $\mathcal{O}(d^{3/2}H^2)$  regret,  
 159 independent of  $K$ . Therefore, our algorithm can forgo clipping for free.

160 **Remark 4** *If it wasn't for the episodic nature of our setting, we could have used the forward algorithm to eliminate the*  
 161 *span control issue. We refer to* [Vovk \(2001\)](#); [Azouy and Warmuth \(2001\)](#) *for a description of this algorithm,* [Ouhamma](#)  
 162 *et al. (2021)* *for a stochastic analysis, and Section 4 therein for an application to linear bandits.*

## 163 5.2 Bounding the learning error

164 To upper-bound this term of the regret, we first show that the estimated value function is optimistic with a constant  
 165 probability. Then, we show that this is enough to control the learning error.

**Stochastic optimism.** The perturbation ensures a constant probability of optimism. Specifically,

$$\begin{aligned} (V_{\hat{\theta}^p, \tilde{\theta}^x, 1} - V_{\hat{\theta}^p, \theta^x, 1}^*)(s_1) &\geq (Q_{\hat{\theta}^p, \tilde{\theta}^x, 1}^* - Q_1^*)(s_1, \pi^*(s_1)) \\ &\geq \underbrace{V_{\hat{\theta}^p, \theta^x}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \theta^x}^{\pi^*}(s_1)}_{\text{first term}} + \underbrace{V_{\hat{\theta}^p, \tilde{\theta}^x}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \theta^x}^{\pi^*}(s_1)}_{\text{second term}} + \underbrace{V_{\hat{\theta}^p, \tilde{\theta}^x}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \tilde{\theta}^x}^{\pi^*}(s_1)}_{\text{third term}} \end{aligned}$$

The first and second terms are perturbation free, we handle them similarly to the estimation error, *i.e.* using concentration arguments for  $\hat{\theta}^p$  and  $\hat{\theta}^x$ . For the third term, we use transportation of rewards (Lemma 23) and anti-concentration of  $\xi_k$  (Lemma 18). We find that with probability at least  $1 - 2\delta$

$$\begin{aligned} (V_{\hat{\theta}^p, \tilde{\theta}^x, 1} - V_{\hat{\theta}^p, \theta^x, 1}^*)(s_1) &\geq \xi_k^\top \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} \left[ \sum_{t=1}^H \frac{\text{Var}^{\theta^x}(r)}{2} (A_t \varphi(\tilde{s}_t, \pi^*(\tilde{s}_t)))_{i \in [d]} \right] B \\ &\quad - Hc(n, \delta) \left\| \sum_{h=1}^H \mathbb{E}_{(\tilde{s}_h)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} [(A_i \varphi(\tilde{s}_h, \pi^*(\tilde{s}_h)))_{i \in [d]}] \right\|_{(\bar{G}_k^p)^{-1}}, \end{aligned}$$

166 where  $c(n, \delta) = \left( \sqrt{\beta^p \beta^x(n, \delta) / \alpha^p} + \sqrt{\beta^x \beta^x(n, \delta) \min\{1, \alpha^p / \alpha^x\} / (2\alpha^x)} \right)$ . Since  $\xi_k \sim \mathcal{N}(0, x_k (\bar{G}_k^p)^{-1})$  and  $x_k \geq H^2 c(n, \delta)^2$ ,  
 167 we get  $\mathbb{P} \left( V_{\hat{\theta}^p, \tilde{\theta}^x, 1}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \theta^x, 1}^*(s_1) \geq 0 \right) \geq \Phi(-1)$ , where  $\Phi$  is the normal CDF. This is ensured by the anti-  
 168 concentration property of Gaussian random variables, see Lemma 18.

169 **From stochastic optimism to error control:** Existing algorithms require the value function to be optimistic (*i.e.*  
 170 negative learning error) with large probability. Contrary to them, BEF-RLSVI only requires the estimated value to be  
 171 optimistic with a constant probability. When it is, the learning happens. Otherwise, the policy is still close to a good  
 172 one thanks to the decreasing estimation error, and the learning still happens. This part of the proof is similar in spirit to  
 173 that of [Zanette et al. \(2020\)](#).

174 Upper bound on  $V_1^*$ : Draw  $(\bar{\xi}_k)_{k \in [K]}$  i.i.d copies of  $(\xi_k)_{k \in [K]}$  and define the event where optimism holds as  $\bar{O}_k \triangleq$   
 175  $\{V_{\hat{\theta}^p, \tilde{\theta}^x, 1}^k(s_1^k) - V_1^*(s_1^k) \geq 0\}$ . This implies that  $V_1^*(s_1^k) \leq \mathbb{E}_{\bar{\xi}_k | \bar{O}_k} [V_{\hat{\theta}^p, \tilde{\theta}^x + \bar{\xi}_k, 1}(s_1^k)]$ .

Lower bound on  $V_{\hat{\theta}^p, \tilde{\theta}^x}$ : Consider  $\underline{V}_1(s_1^k)$  to be a solution of the optimization problem

$$\min_{\xi_k} V_{\hat{\theta}^p, \tilde{\theta}^x + \xi_k, 1}(s_1^k) \quad \text{subject to: } \|\xi_k\|_{\bar{G}_k} \leq \sqrt{x_k d \log(d/\delta)},$$

176 As the injected noise concentrates, we obtain  $\underline{V}_1(s_1^k) \leq V_{\hat{\theta}^p, \tilde{\theta}^x}(s_1^k)$ .

Combination: Using these upper and lower bounds, we show that with probability at least  $1 - \delta$ ,

$$\begin{aligned} V_1^*(s_1^k) - V_{\hat{\theta}^p, \tilde{\theta}^x + \bar{\xi}_k, 1}(s_1^k) &\leq \mathbb{E}_{\bar{\xi}_k | \bar{O}_k} [V_{\hat{\theta}^p, \tilde{\theta}^x + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)] \\ &\leq \left( \mathbb{E}_{\bar{\xi}_k} [V_{\hat{\theta}^p, \tilde{\theta}^x + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)] - \mathbb{E}_{\bar{\xi}_k | \bar{O}_k} [V_{\hat{\theta}^p, \tilde{\theta}^x + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)] \right) / \mathbb{P}(\bar{O}_k), \end{aligned}$$

177 The last step follows from the tower rule. Note that the term inside the expectations is positive with high probability but  
 178 not necessarily in expectation. We follow the lines of the estimation error analysis to complete the proof of Theorem 3.  
 179 We refer to Appendix B.2 for the detailed proof.



## 180 6. Related works: functional representations with regret and tractability

181 Our work extends the endeavor of using functional representations to perform optimal regret minimization in continuous  
182 state-action MDPs. We now provide a few complementary details.

183 *General functional representation.* (Dai et al., 2018) provides the first convergence guarantee for general nonlinear  
184 function representations in the Maximum Entropy RL setting, where entropy of a policy is used as a regularizer to  
185 induce exploration. Thus, the analysis cannot address episodic RL, where we have to explicitly ensure exploration with  
186 optimism. (Wang et al., 2020) proposes a framework that leverages the optimism with confidence bound approach for  
187 general functional representations with bounded Eluder dimensions, which is a complexity measure in RL. However,  
188 knowing the Eluder dimension is crucial for the optimistic confidence bound in their algorithm. Eluder dimension is  
189 not known for MDPs except linear and tabular MDPs. *To concretize our design, we focus on the general but explicit*  
190 *bilinear exponential family of MDPs than any abstract representation.*

191 *Bilinear exponential family of MDPs.* Exponential families are studied widely in RL theory, from bandits to MDPs  
192 Lu et al. (2021); Korda et al. (2013); Filippi et al. (2010); Kveton and Hauskrecht (2006), as an expressive parametric  
193 family to design theoretically-grounded model-based algorithms. Chowdhury et al. (2021) first studies episodic RL  
194 with Bilinear Exponential Family (BEF) of transitions, which is linear in both state-action pairs and the next-state. It  
195 proposes a regularized log-likelihood method to estimate the model parameters, and two optimistic algorithms with  
196 upper confidence bounds and posterior sampling. Due to its generality to unifiedly model tabular MDPs, factored MDPs,  
197 linear MDPs, and linearly controlled dynamical systems, the BEF-family of MDPs has received increasing attention (Li  
198 et al., 2021). Li et al. (2021) estimates the model parameters based on score matching that enables them to replace  
199 regularity assumption on the log-partition function with Fisher-information and assumption on the parameters. Both  
200 (Chowdhury et al., 2021; Li et al., 2021) achieve a worst-case regret of order  $\tilde{O}(\sqrt{d^2 H^4 K})$  for known reward. On a  
201 different note, (Du et al., 2021; Foster et al., 2021) also introduces a new structural framework for generalization in RL,  
202 called bilinear classes as it requires the Bellman error to be upper bounded by a bilinear form. Instead of using bilinear  
203 forms to capture non-linear structures, this class is not identical to BEF class of MDPs, and studying the connection is  
204 out of the scope of this paper. Specifically, *we address the shortcomings of the existing works on BEF-family of MDPs*  
205 *that assume known rewards, absence of RLSVI-type algorithms, and access to oracle planners.*

206 *Tractable planning and linearity.* Planning is a major byproduct of the chosen functional representation. In general,  
207 planning can incur high computational complexity if done naively. Specially, Du et al. (2019) shows that for some  
208 settings, even with a linear  $\epsilon$ -approximation of the  $Q$ -function, a planning procedure able to produce an  $\epsilon$ -optimal policy  
209 has a complexity at least  $2^H$ . Thus, different works (Shariff and Szepesvári, 2020; Lattimore et al., 2020; Van Roy  
210 and Dong, 2019) propose to leverage different low-dimensional representations of value functions or transitions to  
211 perform efficient planning. Here, we take note from (Ren et al., 2021) that Gaussian transitions induce an explicit  
212 linear value function in an RKHS. And generalize this observation with the bilinear exponential. Moreover, using  
213 uniformly good features Rahimi and Recht (2007) to approximate transition dynamics from our model enables us  
214 to design a tractable planner. We provide a detailed discussion of this approximation in Section 4. More practically,  
215 Ren et al. (2021); Nachum and Yang (2021) use representations given by random Fourier features (Rahimi and Recht,  
216 2007) to approximate the transition dynamics and provide experiments validating the benefits of this approach for  
217 high-dimensional Atari-games.

## 218 7. Conclusion and future work

219 We propose the BEF-RLSVI algorithm for the bilinear exponential family of MDPs in the setting of episodic-RL.  
220 BEF-RLSVI explores using a Gaussian perturbation of rewards, and plans tractably ( $\mathcal{O}(pH^3 K \log(HK))$  complexity)  
221 thanks to properties of the RBF kernel. Our proof shows that clipping can be forwent for similar RLSVI-type algorithms.  
222 Moreover, we prove a  $\sqrt{d^3 H^3 K}$  frequentist regret bound, which improves over existing work, accommodates unknown  
223 rewards, and matches the lower bound in terms of  $H$  and  $K$ . Regarding future work, we believe that our proof approach  
224 can be extended to rewards with bounded variance. We also believe that the extra  $\sqrt{d}$  in our bound is an artefact of  
225 the proof, and specifically, the anti-concentration. We will investigate it further. Finally, we plan to study the practical  
226 efficiency of BEF-RLSVI through experiments on tasks with continuous state-action spaces in an extended version of  
227 this work.

228 **References**

- 229 Yasin Abbasi-Yadkori and Csaba Szepesvári. Regret bounds for the adaptive control of linear quadratic systems. In  
 230 *Proceedings of the 24th Annual Conference on Learning Theory*, pages 1–26. JMLR Workshop and Conference  
 231 Proceedings, 2011.
- 232 Marc Abeille and Alessandro Lazaric. Linear thompson sampling revisited. In *Artificial Intelligence and Statistics*,  
 233 pages 176–184. PMLR, 2017.
- 234 Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In  
 235 *International Conference on Machine Learning*, pages 263–272. PMLR, 2017.
- 236 Katy S Azoury and Manfred K Warmuth. Relative loss bounds for on-line density estimation with the exponential  
 237 family of distributions. *Machine Learning*, 43(3):211–246, 2001.
- 238 Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of*  
 239 *independence*. Oxford university press, 2013.
- 240 Sayak Ray Chowdhury and Aditya Gopalan. Online learning in kernelized markov decision processes. In *The 22nd*  
 241 *International Conference on Artificial Intelligence and Statistics*, pages 3197–3205. PMLR, 2019.
- 242 Sayak Ray Chowdhury, Aditya Gopalan, and Odalric-Ambrym Maillard. Reinforcement learning in parametric mdps  
 243 with exponential families. In *International Conference on Artificial Intelligence and Statistics*, pages 1855–1863.  
 244 PMLR, 2021.
- 245 Bo Dai, Albert Shaw, Lihong Li, Lin Xiao, Niao He, Zhen Liu, Jianshu Chen, and Le Song. Sbeed: Convergent  
 246 reinforcement learning with nonlinear function approximation. In *International Conference on Machine Learning*,  
 247 pages 1125–1134. PMLR, 2018.
- 248 Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying pac and regret: Uniform pac bounds for episodic  
 249 reinforcement learning. *Advances in Neural Information Processing Systems*, 30, 2017.
- 250 Omar Darwiche Domingues, Pierre Ménard, Emilie Kaufmann, and Michal Valko. Episodic reinforcement learning in  
 251 finite mdps: Minimax lower bounds revisited. In *Algorithmic Learning Theory*, pages 578–598. PMLR, 2021.
- 252 Simon Du, Sham Kakade, Jason Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong Wang. Bilinear classes:  
 253 A structural framework for provable generalization in rl. In *International Conference on Machine Learning*, pages  
 254 2826–2836. PMLR, 2021.
- 255 Simon S Du, Sham M Kakade, Ruosong Wang, and Lin F Yang. Is a good representation sufficient for sample efficient  
 256 reinforcement learning? *arXiv preprint arXiv:1910.03016*, 2019.
- 257 Sarah Filippi, Olivier Cappe, Aurélien Garivier, and Csaba Szepesvári. Parametric bandits: The generalized linear case.  
 258 *Advances in Neural Information Processing Systems*, 23, 2010.
- 259 Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of interactive decision  
 260 making. *arXiv preprint arXiv:2112.13487*, 2021.
- 261 Haque Ishfaq, Qiwen Cui, Viet Nguyen, Alex Ayoub, Zhuoran Yang, Zhaoran Wang, Doina Precup, and Lin F Yang.  
 262 Randomized exploration for reinforcement learning with general value function approximation. *arXiv preprint*  
 263 *arXiv:2106.07841*, 2021.
- 264 Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear  
 265 function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR, 2020.
- 266 Kwang-Sung Jun, Aniruddha Bhargava, Robert Nowak, and Rebecca Willett. Scalable generalized linear bandits:  
 267 Online computation and hashing. *arXiv preprint arXiv:1706.00136*, 2017.
- 268 Michael Kearns and Daphne Koller. Efficient reinforcement learning in factored mdps. In *IJCAI*, volume 16, pages  
 269 740–747, 1999.

- 270 Nathaniel Korda, Emilie Kaufmann, and Remi Munos. Thompson sampling for 1-dimensional exponential family  
271 bandits. *Advances in neural information processing systems*, 26, 2013.
- 272 Branislav Kveton and Milos Hauskrecht. Solving factored mdps with exponential-family transition models. In *ICAPS*,  
273 pages 114–120, 2006.
- 274 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 275 Tor Lattimore, Csaba Szepesvari, and Gellert Weisz. Learning with good feature representations in bandits and in rl  
276 with a generative model. In *International Conference on Machine Learning*, pages 5662–5670. PMLR, 2020.
- 277 Gene Li, Junbo Li, Nathan Srebro, Zhaoran Wang, and Zhuoran Yang. Exponential family model-based reinforcement  
278 learning via score matching. *arXiv preprint arXiv:2112.14195*, 2021.
- 279 Yangyi Lu, Amirhossein Meisami, and Ambuj Tewari. Low-rank generalized linear bandit problems. In *International  
280 Conference on Artificial Intelligence and Statistics*, pages 460–468. PMLR, 2021.
- 281 Ofir Nachum and Mengjiao Yang. Provable representation learning for imitation with contrastive fourier features.  
282 *Advances in Neural Information Processing Systems*, 34, 2021.
- 283 Ian Osband and Benjamin Van Roy. Model-based reinforcement learning and the eluder dimension. *Advances in Neural  
284 Information Processing Systems*, 27, 2014.
- 285 Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via posterior sampling.  
286 *Advances in Neural Information Processing Systems*, 26, 2013.
- 287 Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value functions. In  
288 *International Conference on Machine Learning*, pages 2377–2386. PMLR, 2016.
- 289 Reda Ouhamma, Odalric-Ambrym Maillard, and Vianney Perchet. Stochastic online linear regression: the forward  
290 algorithm to replace ridge. *Advances in Neural Information Processing Systems*, 34:24430–24441, 2021.
- 291 Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Advances in neural information  
292 processing systems*, 20, 2007.
- 293 Tongzheng Ren, Tianjun Zhang, Csaba Szepesvári, and Bo Dai. A free lunch from the noise: Provable and practical  
294 exploration for representation learning. *arXiv preprint arXiv:2111.11485*, 2021.
- 295 Roshan Shariff and Csaba Szepesvári. Efficient planning in large mdps with weak linear function approximation. *arXiv  
296 preprint arXiv:2007.06184*, 2020.
- 297 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 298 Benjamin Van Roy and Shi Dong. Comments on the du-kakade-wang-yang lower bounds. *arXiv preprint  
299 arXiv:1911.07910*, 2019.
- 300 Volodya Vovk. Competitive on-line statistics. *International Statistical Review*, 69(2):213–248, 2001.
- 301 Ruosong Wang, Ruslan Salakhutdinov, and Lin F Yang. Reinforcement learning with general value function approxima-  
302 tion: Provably efficient approach via bounded eluder dimension. *arXiv preprint arXiv:2005.10804*, 2020.
- 303 Andrea Zanette, David Brandfonbrener, Emma Brunskill, Matteo Pirootta, and Alessandro Lazaric. Frequentist regret  
304 bounds for randomized least-squares value iteration. In *International Conference on Artificial Intelligence and  
305 Statistics*, pages 1954–1964. PMLR, 2020.

306 **Appendix A. Notations**

 307 We dedicate this section to index all the notations used in this paper. Note that every notation is defined when it is  
 308 introduced as well.

Table 2: Notations

$H$	$\stackrel{\text{def}}{=}$	number of steps in a given episode
$K$	$\stackrel{\text{def}}{=}$	number of episodes
$T$	$\stackrel{\text{def}}{=}$	$KH$ , total number of steps
$s_h^k$	$\stackrel{\text{def}}{=}$	state at time $h$ of episode $k$ , denoted $s_h$ when $k$ is clear from context
$a_h^k$	$\stackrel{\text{def}}{=}$	action at time $h$ of episode $k$ , denoted $a_h$ when $k$ is clear from context
$r(s, a)$	$\stackrel{\text{def}}{=}$	realization of the reward in state $s$ under action $a$
$\theta^p$	$\stackrel{\text{def}}{=}$	parameter of the transition distribution, $\in \mathbb{R}^d$
$\theta^r$	$\stackrel{\text{def}}{=}$	parameter of the reward distribution, $\in \mathbb{R}^d$
$\theta$	$\stackrel{\text{def}}{=}$	$\in \mathbb{R}^d$ denotes either $\theta^r$ or $\theta^p$ , unless stated otherwise
$\hat{\theta}$	$\stackrel{\text{def}}{=}$	$\theta$ estimator with Maximum Likelihood unless stated otherwise
$\tilde{\theta}$	$\stackrel{\text{def}}{=}$	$\hat{\theta} + \xi$ where $\xi$ is a chosen noise. Perturbed estimation of $\theta$ .
$[\theta_1, \theta_2]$	$\stackrel{\text{def}}{=}$	the $d$ -dimensional $\ell_\infty$ hypercube joining $\theta_1$ and $\theta_2$
$\mathbb{P}_{\theta^p}$	$\stackrel{\text{def}}{=}$	transition under the exponential family model with parameter $\theta^p$
$\psi$	$\stackrel{\text{def}}{=}$	feature function, $\in (\mathbb{R}_+^p)^{\mathcal{S}}$
$\varphi$	$\stackrel{\text{def}}{=}$	feature function, $\in (\mathbb{R}_+^q)^{\mathcal{S} \times \mathcal{A}}$
$B$	$\stackrel{\text{def}}{=}$	$p$ -dimensional vector
$M_\theta$	$\stackrel{\text{def}}{=}$	$\sum_{i=1}^d \theta_i A_i$ , where $A_i$ are $p \times q$ matrices.
$Z^r$	$\stackrel{\text{def}}{=}$	the rewards' log partition function
$Z^p$	$\stackrel{\text{def}}{=}$	the transitions' log partition function
$\mathcal{H}$	$\stackrel{\text{def}}{=}$	Hilbert space where we decompose transitions
$\mu^p$	$\stackrel{\text{def}}{=}$	feature function after decomposition, $\in (\mathbb{R}_+)^{\mathcal{S} \times \mathcal{H}}$
$\phi^p$	$\stackrel{\text{def}}{=}$	feature function after decomposition, $\in (\mathbb{R}_+)^{\mathcal{S} \times \mathcal{A} \times \mathcal{H}}$
$G_{s,a}$	$\stackrel{\text{def}}{=}$	$(\varphi(s, a)^\top A_i^\top A_j \varphi(s, a))_{i,j \in [d]}$
$\bar{G}_k^r$	$\stackrel{\text{def}}{=}$	$\bar{G}_{(k-1)h}^r = \frac{\eta}{\alpha^r} \mathbb{A} + \sum_{\tau=1}^{k-1} \sum_{h=1}^H G_{s_h^\tau, a_h^\tau}$
$\bar{G}_k^p$	$\stackrel{\text{def}}{=}$	$\bar{G}_{(k-1)h}^p = \frac{\eta}{\alpha^p} \mathbb{A} + \sum_{\tau=1}^{k-1} \sum_{h=1}^H G_{s_h^\tau, a_h^\tau}$
$\mathbb{C}_{s,a}^\theta [\psi(s')]$	$\stackrel{\text{def}}{=}$	$\mathbb{E}_{s,a}^\theta [\psi(s') \psi(s')^\top] - \mathbb{E}_{s,a}^\theta [\psi(s')] \mathbb{E}_{s,a}^\theta [\psi(s')^\top]$
$\beta^p$	$\stackrel{\text{def}}{=}$	$\sup_{\theta, s, a} \lambda_{\max} (\mathbb{C}_{s,a}^\theta [\psi(s')])$ linked to the maximum eigenvalue of $\nabla^2 Z^p$
$\alpha^p$	$\stackrel{\text{def}}{=}$	$\inf_{\theta, s, a} \lambda_{\min} (\mathbb{C}_{s,a}^\theta [\psi(s')])$ linked to the minimum eigenvalue of $\nabla^2 Z^p$
$\beta^r$	$\stackrel{\text{def}}{=}$	$\lambda_{\max} (BB^\top) \sup_{\theta, s, a} \text{Var}_{s,a}^\theta (r)$ , linked to the maximum eigenvalue of $\nabla^2 Z^r$
$\alpha^r$	$\stackrel{\text{def}}{=}$	$\lambda_{\min} (BB^\top) \inf_{\theta, s, a} \text{Var}_{s,a}^\theta (r)$ , linked to the minimum eigenvalue of $\nabla^2 Z^r$

## 309 Appendix B. Regret analysis

We provide a high probability analysis of the regret of BEF-RLSVI under standard regularity assumptions of the representation. First we recall the regret definition then we separate the perturbation error from the statistical estimation:

$$\mathcal{R}(K) = \sum_{k=1}^K (V_{\theta^p, \theta^x, 1}^* - V_{\theta^p, \theta^x, 1}^{\pi_k})(s_1^k) = \sum_{k=1}^K \left( \underbrace{V_{\theta^p, \theta^x, 1}^* - V_{\hat{\theta}^p, \hat{\theta}^x, 1}^{\pi_k}}_{\text{learning}} + \underbrace{V_{\hat{\theta}^p, \hat{\theta}^x, 1}^{\pi_k} - V_{\theta^p, \theta^x, 1}^{\pi_k}}_{\text{Estimation}} \right) (s_1^k)$$

### 310 B.1 Estimation error

311 To show that the estimation error ( $\sum_{k=1}^K V_{\hat{\theta}^p, \hat{\theta}^x, 1}^{\pi_k} - V_{\theta^p, \theta^x, 1}^{\pi_k}$ ) can be controlled, we decompose it to an error that comes  
 312 from the estimation of the transition parameter and one that comes from the estimation of the reward parameter:

$$V_{\hat{\theta}^p, \hat{\theta}^x, 1}^{\pi_k}(s_1^k) - V_{\theta^p, \theta^x, 1}^{\pi_k}(s_1^k) = \underbrace{V_{\hat{\theta}^p, \theta^x, 1}^{\pi_k}(s_1^k) - V_{\theta^p, \theta^x, 1}^{\pi_k}(s_1^k)}_{\text{transition estimation}} + \underbrace{V_{\hat{\theta}^p, \hat{\theta}^x, 1}^{\pi_k}(s_1^k) - V_{\hat{\theta}^p, \theta^x, 1}^{\pi_k}(s_1^k)}_{\text{reward estimation}},$$

313 we control each term separately in Section B.1.1 and Section B.1.2. Therefore, we obtain the following lemma  
 314 controlling the estimation error.

**Lemma 5** *The estimation error satisfies, with probability at least  $1 - 5\delta$*

$$\begin{aligned} \sum_{k=1}^K V_{\hat{\theta}^p, \hat{\theta}^x, 1}^{\pi_k}(s_1^k) - V_{\theta^p, \theta^x, 1}^{\pi_k}(s_1^k) &\leq 2H \sqrt{\frac{2\beta^p}{\alpha^p} \beta^p(N, \delta) N \gamma_K^p} + 2H \sqrt{2N \log(1/\delta)} \\ &+ \left[ \sqrt{KHd \log(1 + \alpha^x \eta^{-1} B_{\varphi, \mathbb{A}} n)} + C_d \sqrt{Hd \log(1 + \alpha \eta^{-1} B_{\varphi, \mathbb{A}} H)} \right] \times \left( \sqrt{\frac{\beta^x(n, \delta)}{2\alpha^x}} \right. \\ &\left. + c \sqrt{(\max_k x_k) d \log(dK/\delta)} \right) \beta^x + \sqrt{2KHd \log(1 + \alpha^x \eta^{-1} B_{\varphi, \mathbb{A}} n) \log(1/\delta)} \end{aligned}$$

315 where for  $i \in [p, x]$ ,  $\beta^i(K, \delta) \triangleq \frac{\eta}{2} B_{\mathbb{A}}^2 + \gamma_K^i + \log(1/\delta)$ , and  $\gamma_K^i \triangleq d \log(1 + \frac{\beta^i}{\eta} B_{\varphi, \mathbb{A}} HK)$ . Also,  $C_d \triangleq$   
 316  $\frac{3d}{\log(2)} \log\left(1 + \frac{\alpha^x \|\mathbb{A}\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)}\right)$ , and  $c$  is a universal constant.

317 **Proof** It follows directly by combining Lemma 6 and Lemma 9 using a union bound. ■

318

#### 319 B.1.1 TRANSITION ESTIMATION

320 The goal of this section is to prove the following lemma which bounds the regret due to transition estimation.

**Lemma 6** *We have, with probability at least  $1 - 2\delta$*

$$\sum_{k=1}^K V_{\hat{\theta}^p, \theta^x, 1}^{\pi_k}(s_1^k) - V_{\theta^p, \theta^x, 1}^{\pi_k}(s_1^k) \leq 2H \sqrt{\frac{2\beta^p}{\alpha^p} \beta^p(N, \delta) N \gamma_K^p} + 2H \sqrt{2N \log(1/\delta)}$$

321 where  $\gamma_K^p := d \log(1 + \beta^p \eta^{-1} B_{\varphi, \mathbb{A}} HK)$ , and  $\beta^p(K, \delta) \triangleq \frac{\eta}{2} B_{\mathbb{A}}^2 + \gamma_K^p + \log(1/\delta)$ .

322 **Proof** The proof proceeds in two parts. First, we will reveal a bound in terms of the induced local geometry, *i.e.* a  
 323 bound in terms of KL-divergence. Second, we explicit the bound by transferring the induced local geometry to the  
 324 euclidean one.

325 **1) Bound in terms of local geometry.** We provide a bound on the estimation error of the transition in terms of KL  
 326 divergences, for that end we show that the estimation error can be decomposed and well controlled. We start by writing  
 327 the one-step decomposition:

$$\begin{aligned}
 & V_{\hat{\theta}^p, \theta^x, 1}^\pi(s_1^k) - V_{\theta^p, \theta^x, 1}^\pi(s_1^k) \\
 &= \mathbb{E}_{s_1^k, a_1^k}^{\hat{\theta}^p} \left[ V_{\hat{\theta}^p, \theta^x, 2}^\pi \right] - \mathbb{E}_{s_1^k, a_1^k}^{\theta^p} \left[ V_{\hat{\theta}^p, \theta^x, 2}^\pi \right] + \mathbb{E}_{s_1^k, a_1^k}^{\theta^p} \left[ V_{\hat{\theta}^p, \theta^x, 2}^\pi - V_{\theta^p, \theta^x, 2}^\pi \right] \\
 &= \mathbb{E}_{s_1^k, a_1^k}^{\hat{\theta}^p} \left[ V_{\hat{\theta}^p, \theta^x, 2}^\pi \right] - \mathbb{E}_{s_1^k, a_1^k}^{\theta^p} \left[ V_{\hat{\theta}^p, \theta^x, 2}^\pi \right] + V_{\hat{\theta}^p, \theta^x, 2}^\pi(s_{2k}) - V_{\theta^p, \theta^x, 2}^\pi(s_{2k}) + \zeta_1^k \\
 &= \sum_{h=1}^H \mathbb{E}_{s_{hk}, a_{hk}}^{\hat{\theta}^p} \left[ V_{\hat{\theta}^p, \theta^x, h+1}^\pi \right] - \mathbb{E}_{s_{hk}, a_{hk}}^{\theta^p} \left[ V_{\hat{\theta}^p, \theta^x, h+1}^\pi \right] + \zeta_{hk}
 \end{aligned}$$

where  $\zeta_{hk} = \mathbb{E}_{s_{hk}, a_{hk}}^{\theta^p} \left[ V_{\hat{\theta}^p, \theta^x, h+1}^\pi - V_{\theta^p, \theta^x, h+1}^\pi \right] - \left( V_{\hat{\theta}^p, \theta^x, h+1}^\pi(s_{h+1k}) - V_{\theta^p, \theta^x, h+1}^\pi(s_{h+1k}) \right)$  is a martingale sequence, and the last equality comes by induction. Here we consider the true reward parameter which verifies  $|\mathbb{E}_{\theta^x}[r(s, a)]| \leq 1$  by assumption, therefore  $|\zeta_{hk}| \leq 2H$ . Using the Azuma-Hoeffding inequality [Boucheron et al. \(2013\)](#), with probability at least  $1 - \delta$

$$\sum_{k=1}^K \sum_{h=1}^H \zeta_{hk} \leq 2H \sqrt{2KH \log(1/\delta)}$$

We finish bounding the first term using Lemma 19, indeed

$$\begin{aligned}
 \mathbb{E}_{s_{hk}, a_{hk}}^{\hat{\theta}^p} \left[ V_{\hat{\theta}^p, \theta^x, h+1}^\pi \right] - \mathbb{E}_{s_{hk}, a_{hk}}^{\theta^p} \left[ V_{\hat{\theta}^p, \theta^x, h+1}^\pi \right] &\leq H \sqrt{2 \text{KL}_{s_{hk}, a_{hk}}(\theta^p, \hat{\theta}^p)} \\
 &\leq H \min \left\{ 1, \sqrt{2 \text{KL}_{s_{hk}, a_{hk}}(\theta^p, \hat{\theta}^p)} \right\},
 \end{aligned}$$

328 the last inequality follows because  $\forall h, \mathbb{S}(V_{\hat{\theta}^p, \theta^x, h+1}) \leq H$ .

329 **Remark 7** Traditionally, the expected value difference bound follows from the simulation lemma [Ren et al. \(2021\)](#). The  
 330 simulation lemma incurs an extra  $\sqrt{H}$  factor compared to our bound.

We deduce that with probability at least  $1 - \delta$ :

$$\begin{aligned}
 & \sum_{k=1}^K V_{\hat{\theta}^p, \theta^x}(s_1^k) - V_{\theta^p, \theta^x}(s_1^k) \\
 & \leq H \sum_{k=1}^K \min \left\{ 1, \sum_{h=1}^H \sqrt{2 \text{KL}_{s_{hk}, a_{hk}}(\theta^p, \hat{\theta}^p)} \right\} + 2H \sqrt{2KH \log(1/\delta)} \tag{11}
 \end{aligned}$$

**2) Bounding the sum of KL divergences.** we explicit the bound of inequality (11) using Assumption 2 along with properties of the exponential family (*cf.* Section D.3). We have for all  $(s, a)$ ,

$$\forall \theta^p, \theta^{p'}, \quad \frac{\alpha^p}{2} \|\theta^{p'} - \theta^p\|_{G_{s,a}}^2 \leq \text{KL}_{s,a}(\theta^p, \theta^{p'}) \leq \frac{\beta^p}{2} \|\theta^{p'} - \theta^p\|_{G_{s,a}}^2. \tag{12}$$

This implies that

$$\text{KL}_{s,a}(\hat{\theta}^p(k), \theta^p) \leq \frac{\beta^p}{2} \left\| \theta^p - \hat{\theta}^p(k) \right\|_{G_{s,a}}^2 \leq \beta^p \left\| (\bar{G}_k^p)^{-1/2} G_{s,a} (\bar{G}_k^p)^{-1/2} \right\| \frac{1}{2} \left\| \theta^p - \hat{\theta}^p(k) \right\|_{\bar{G}_k^p}^2,$$

331 where  $\bar{G}_k^{\mathbb{P}} \equiv \bar{G}_{(k-1)H}^{\mathbb{P}} := G_k + (\alpha^{\mathbb{P}})^{-1}\eta\mathbb{A}$  and  $G_k \equiv \sum_{\tau=1}^{k-1} \sum_{h=1}^H G_{s_s^{\tau}, a_h^{\tau}}$ .

From Corollary 14, with probability at least  $1 - \delta$  and for all  $k \in \mathbb{N}$

$$\left\| \theta^{\mathbb{P}} - \hat{\theta}^{\mathbb{P}}(k) \right\|_{\bar{G}_k^{\mathbb{P}}}^2 \leq 2\beta^{\mathbb{P}}(k, \delta)/\alpha^{\mathbb{P}}.$$

Also, using Lemma 24, we have

$$\sum_{t=1}^T \sum_{h=1}^H \min \left\{ 1, \left\| (\bar{G}_k^{\mathbb{P}})^{-1/2} G_{s, a} (\bar{G}_k^{\mathbb{P}})^{-1/2} \right\| \right\} \leq 2d \log(1 + \alpha^{\mathbb{P}} \eta^{-1} B_{\varphi, \mathbb{A}} H K).$$

332 Combining these two results we obtain, with probability at least  $1 - \delta$ :

$$\sum_{t=1}^T \sum_{h=1}^H \min \left\{ 1, \text{KL}_{s_h^t, a_h^t} \left( \hat{\theta}^{\mathbb{P}}(k), \theta^{\mathbb{P}} \right) \right\} \leq \frac{2\beta^{\mathbb{P}}}{\alpha^{\mathbb{P}}} \beta^{\mathbb{P}}(K, \delta) \gamma_K^{\mathbb{P}}. \quad (13)$$

333 **Remark 8** Notice that the minimum with 1 is crucial, indeed, without it the bound deteriorates by a factor  $H$  as was  
334 the case in Chowdhury et al. (2021).

3) *Combining the bounds.* By applying Cauchy-Schwarz in inequality (11), we obtain, with probability at least  $1 - \delta$ , and for all  $K \in \mathbb{N}$

$$\sum_{k=1}^K V_{\hat{\theta}^{\mathbb{P}}, \theta^{\mathbb{P}}}(s_1^k) - V_{\theta^{\mathbb{P}}, \theta^{\mathbb{P}}}(s_1^k) \leq H \sqrt{2 \sum_{k=1}^K \sum_{h=1}^H \text{KL}_{s_{hk}, a_{hk}}(\theta^{\mathbb{P}}, \hat{\theta}^{\mathbb{P}}) + 2H \sqrt{2KH \log(1/\delta)}}.$$

335 Injecting inequality (13) proves the desired result with probability at least  $1 - 2\delta$ . ■

336

### 337 B.1.2 REWARD ESTIMATION

338 Now, we provide the bound over the regret due to estimating the reward parameter.

**Lemma 9** With probability at least  $1 - 3\delta$ , the following result holds true.

$$\begin{aligned} \sum_{k=1}^K V_{\hat{\theta}^{\mathbb{P}}, \tilde{\theta}^{\mathbb{P}}, 1}(s_1^k) - V_{\hat{\theta}^{\mathbb{P}}, \theta^{\mathbb{P}}, 1}(s_1^k) &\leq \left( \sqrt{\frac{\beta^{\mathbb{P}}(K, \delta)}{2\alpha^{\mathbb{P}}}} + c \sqrt{(\max_{k \leq K} x_k) d \log(dK/\delta)} \right) \beta^{\mathbb{P}} \\ &\times \left( \sqrt{C_d \left( 1 + \frac{\alpha^{\mathbb{P}} B_{\varphi, \mathbb{A}} H}{\eta} \right)} + \sqrt{K \log(e/\delta^2)} \right) \sqrt{H d \log(1 + \alpha^{\mathbb{P}} \eta^{-1} B_{\varphi, \mathbb{A}} H K)}, \end{aligned}$$

339 where  $\beta^{\mathbb{P}}(K, \delta) \triangleq \frac{\eta}{2} B_{\mathbb{A}}^2 + \gamma_K^{\mathbb{P}} + \log(1/\delta)$ , and  $\gamma_K^{\mathbb{P}} \triangleq d \log(1 + \frac{\beta^{\mathbb{P}}}{\eta} B_{\varphi, \mathbb{A}} H K)$ . Also,  $C_d \triangleq \frac{3d}{\log(2)} \log \left( 1 + \frac{\alpha^{\mathbb{P}} \|\mathbb{A}\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)} \right)$ ,

340 and  $c$  is a universal constant.

341 **Proof** The reward estimation error in Equation (10) can be written explicitly. Indeed, using Lemma 23

$$\begin{aligned}
 V_{\hat{\theta}^r, \tilde{\theta}^r, 1}^\pi(s_1^k) - V_{\hat{\theta}^r, \theta^r, 1}^\pi(s_1^k) &= \mathbb{E}_{(\tilde{s}_h)_{1 \leq h \leq H} \sim \pi | \hat{\theta}^r, s_1^k} \left[ \sum_{h=1}^H \frac{\text{Var}_{\tilde{s}_h, \pi(\tilde{s}_h)}(r)}{2} B^\top M_{\tilde{\theta}^r - \theta^r} \varphi(\tilde{s}_h, \pi(\tilde{s}_h)) \right] \\
 &\leq \mathbb{E} \left[ \sum_{h=1}^H \frac{\text{Var}_{\tilde{s}_h, \pi(\tilde{s}_h)}(r)}{2} \|\tilde{\theta}^r - \theta^r\|_{\tilde{G}_k^r} \|(B^\top A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\tilde{G}_k^r)^{-1}} \right] \\
 &\leq \|\tilde{\theta}^r - \theta^r\|_{\tilde{G}_k^r} \mathbb{E} \left[ \sum_{h=1}^H \frac{\text{Var}_{\tilde{s}_h, \pi(\tilde{s}_h)}(r)}{2} \|(B^\top A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\tilde{G}_k^r)^{-1}} \right] \\
 &\leq \|\tilde{\theta}^r - \theta^r\|_{\tilde{G}_k^r} \frac{\beta^r}{2} \mathbb{E} \left[ \underbrace{\sum_{h=1}^H \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\tilde{G}_k^r)^{-1}}}_{\stackrel{\text{def}}{=} \text{traj}_k} \right],
 \end{aligned}$$

342 where  $\text{traj}_k \stackrel{\text{def}}{=} \sum_{h=1}^H \|(A_i \varphi(s_h, \pi(s_h)))_{1 \leq i \leq d}\|_{(G_k^r)^{-1}}$ .

**Bad rounds.** We separate the analysis of this estimation error into bad and good rounds. Here we analyze the bad rounds, which are define by the following set:

$$\mathcal{T} = \{k \in \mathbb{N}^*, \exists h \in [H], \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\tilde{G}_k^r)^{-1}} \geq 1\}$$

1) We know that  $\|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\tilde{G}_k^r)^{-1}} \leq \|\mathbb{A}\|_2^2 B_{\varphi, \mathbb{A}}^2$ . Consequently, according to Lemma 26

$$|\mathcal{T}| \leq \frac{3d}{\log(2)} \log \left( 1 + \frac{\alpha \|\mathbb{A}\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)} \right).$$

343 2) Since  $G_k$  is positive semi-definite, we have  $\tilde{G}_k^r \succeq (\alpha^r)^{-1} \eta \mathbb{A}$ , and in turn, for all state-action couples  $(s, a)$ ,  
 344  $\|(\tilde{G}_k^r)^{-1} G_{s,a}\| \leq \frac{\alpha^r}{\eta} \|\mathbb{A}^{-1} G_{s,a}\| \leq \frac{\alpha^r B_{\varphi, \mathbb{A}}}{\eta}$ .

This further yields

$$\left\| I + (\tilde{G}_k^r)^{-1} \sum_{h=1}^H G_{s_h^t, a_h^t} \right\| \leq 1 + \sum_{h=1}^H \left\| (\tilde{G}_k^r)^{-1} G_{s_h^t, a_h^t} \right\| \leq 1 + \frac{\alpha^r B_{\varphi, \mathbb{A}} H}{\eta}.$$

Let us define  $\tilde{G}_{k+H}^r := \tilde{G}_k^r + \sum_{h=1}^H G_{s_h^k, a_h^k}$ . Then,

$$\tilde{G}_{k+H}^{-1} G_{s,a} = \left( I + (\tilde{G}_k^r)^{-1} \sum_{h=1}^H G_{s_h^t, a_h^t} \right)^{-1} (\tilde{G}_k^r)^{-1} G_{s,a}.$$

Therefore, for all pairs  $(s, a)$ ,

$$\begin{aligned}
 \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\tilde{G}_k^r)^{-1}} &= \sqrt{\text{tr}((A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}^\top (\tilde{G}_k^r)^{-1} (A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d})} \\
 &= \sqrt{\text{tr} \left( \left( 1 + \frac{\alpha^r B_{\varphi, \mathbb{A}} H}{\eta} \right) (\tilde{G}_{k+H}^r)^{-1} G_{s,a} \right)} \\
 &\leq \sqrt{\left( 1 + \frac{\alpha^r B_{\varphi, \mathbb{A}} H}{\eta} \right) \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\tilde{G}_{k+H}^r)^{-1}}}
 \end{aligned}$$

Since  $\|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\tilde{G}_{k+H}^r)^{-1}} \leq 1$ , we have  $\|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\tilde{G}_k^r)^{-1}} \leq \min \left\{ 1, \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\tilde{G}_k^r)^{-1}} \right\}$ .  
 Consequently

$$\sum_{h=1}^H \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\tilde{G}_{k+H}^r)^{-1}} \leq \sqrt{H d \log(1 + \alpha^r \eta^{-1} B_{\varphi, \mathbb{A}} H)}.$$



3) From 1) and 2), we deduce that the total regret induced by rounds from  $\mathcal{T}$  is bounded.

$$\sum_{k \in \mathcal{T}} \sum_{h \in [H]} V_{\hat{\theta}^p, \tilde{\theta}^r, 1}^\pi(s_1^k) - V_{\hat{\theta}^p, \theta^r, 1}^\pi(s_1^k) \leq \|\tilde{\theta}^r - \theta^r\|_{G_k^r} \frac{\beta^x}{2} \sqrt{\frac{3d}{\log(2)} \log\left(1 + \frac{\alpha^x \|\mathbb{A}\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)}\right) \left(1 + \frac{\alpha^x B_{\varphi, \mathbb{A}} H}{\eta}\right) H d \log(1 + \alpha^x \eta^{-1} B_{\varphi, \mathbb{A}} H)} \quad (14)$$

345 **Remark 10** *The bad rounds analysis is one of our most important contributions as it enables us to forgo clipping*  
 346 *without consequences. Consequently, this is a novel method to control the reward estimation error that improves on*  
 347 *existing work for whom clipping was essential.*

**Good rounds.** Going forward we consider rounds from  $\bar{\mathcal{T}}$ . Let us define

$$\zeta'_k \stackrel{\text{def}}{=} \text{traj}_k - \mathbb{E}_{(\bar{s}_h)_{1 \leq h \leq H} \sim \pi|\hat{\theta}^p, s_1^k} [\widetilde{\text{traj}}_k].$$

where  $\widetilde{\text{traj}}_k$  is the same quantity as  $\text{traj}_k$  but with a random realization of state transitions.

Since all feature norms are smaller than one,  $(\zeta'_k)_k$  is a martingale sequence with  $|\zeta'_k| \leq \sqrt{H d \log(1 + \alpha^x \eta^{-1} B_{\varphi, \mathbb{A}} H K)}$ . We deduce that with probability at least  $1 - \delta$ :

$$\sum_{k=1}^K \zeta'_k \leq \sqrt{2K H d \log(1 + \alpha^x \eta^{-1} B_{\varphi, \mathbb{A}} H K) \log(1/\delta)}$$

Therefore, we have with probability at least  $1 - 3\delta$ :

$$\sum_{k \in \mathcal{T}^c} V_{\hat{\theta}^p, \tilde{\theta}^r, 1}^\pi(s_1^k) - V_{\hat{\theta}^p, \theta^r, 1}^\pi(s_1^k) \leq \left( \sqrt{\frac{\beta^x(K, \delta)}{2\alpha^x}} + c \sqrt{(\max_k x_k) d \log(dK/\delta)} \right) \times \beta^x \sqrt{K H d \log(1 + \alpha^x \eta^{-1} B_{\varphi, \mathbb{A}} H K) \log(e/\delta^2)}.$$

348 The last inequality follows from controlling the concentration of the reward parameter. First we observe that (Corol-  
 349 lary 16) with probability at least  $1 - \delta$ , uniformly over  $k \in \mathbb{N}$ ,  $\left\| \theta^r - \hat{\theta}^r(k) \right\|_{G_k^r}^2 \leq \frac{2}{\alpha^x} \beta^x(k, \delta)$ . Second, we also have  
 350 that for all  $k \geq 1$ , with probability at least  $1 - \delta$ ,  $\|\xi_k\|_{G_k^r} \leq c \sqrt{x_k d \log(d/\delta)}$ , we then use a union bound. Combining  
 351 with Equation (14) we find

$$\sum_{k=1}^K V_{\hat{\theta}^p, \tilde{\theta}^r, 1}^\pi(s_1^k) - V_{\hat{\theta}^p, \theta^r, 1}^\pi(s_1^k) \leq \left( \sqrt{\frac{\beta^x(K, \delta)}{2\alpha^x}} + c \sqrt{(\max_k x_k) d \log(dK/\delta)} \right) \times \beta^x \sqrt{K H d \log(1 + \alpha^x \eta^{-1} B_{\varphi, \mathbb{A}} H K) \log(e/\delta^2)}.$$

352 This concludes the proof. ■

353

354 **Remark 11** *If we use Lemma 23 without the martingale difference sequence, it will lead to a linear regret. Indeed, the*  
 355 *span of the sum of norms over an episode is of order  $\sqrt{H}$ . Using the martingale technique instead allows us to retrieve*  
 356 *a telescopic sum controlled using the elliptical lemma, this is essential to obtaining a sub-linear regret bound.*

357 **B.2 Learning error**

 358 We now start the control of an important regret term, due to the distance between the estimated value function and the  
 359 optimal value function.

**Lemma 12** *If the variance parameter of the injected noise  $(\xi_k)_k$  satisfies*

$$x_k \geq \left( H \sqrt{\frac{\beta^p \beta^r(k, \delta)}{\alpha^p \alpha^r}} + \frac{\sqrt{\beta^r \beta^r(k, \delta) \min\{1, \frac{\alpha^p}{\alpha^r}\}}}{2\alpha^r} \right),$$

*then the learning error is controlled with probability at least  $1 - 2\delta$  as*

$$\begin{aligned} \sum_{k=1}^K V_1^*(s_1^k) - V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_{k,1}}^{\pi}(s_1^k) &\leq \frac{d\beta^r \sqrt{x_k} \left(1 + \sqrt{\log(d/\delta)}\right)}{\Phi(-1)} \sqrt{H \log(1 + \alpha^r \eta^{-1} B_{\varphi, \Delta} H K)} \\ &\times \left( \sqrt{C_d \left(1 + \frac{\alpha^r B_{\varphi, A} H}{\eta}\right)} + \sqrt{K \log(e/\delta^2)} \right), \end{aligned}$$

 360 where for  $i \in [p, r]$ ,  $\beta^i(K, \delta) \triangleq \frac{\eta}{2} B_{\Delta}^2 + \gamma_K^i + \log(1/\delta)$ , and  $\gamma_K^i \triangleq d \log(1 + \frac{\beta^i}{\eta} B_{\varphi, \Delta} H K)$ . Also  $C_d \stackrel{\text{def}}{=} \frac{3d}{\log(2)} \log\left(1 + \frac{\alpha^r \|\Delta\|_2^2 B_{\varphi, \Delta}^2}{\eta \log(2)}\right)$ , and  $\Phi$  is the normal CDF.  
 361

 362 This result basically means that we are no longer obliged to follow optimistic value functions, the perturbed estimation  
 363 is enough to have a tight bound on the learning error.

 364 **B.2.1 STOCHASTIC OPTIMISM**

 The goal here is to show that by injecting our carefully designed noise in the rewards we can ensure optimism with a constant probability. Consider the optimal policy  $\pi^*$ , we have:

$$\begin{aligned} (V_{\hat{\theta}^p, \hat{\theta}^r, 1} - V_{\theta^p, \theta^r, 1}^*)(s_1) &\geq (Q_{\hat{\theta}^p, \hat{\theta}^r, 1}^* - Q_1^*)(s_1, \pi^*(s_1)) \\ &\geq \underbrace{V_{\hat{\theta}^p, \theta^r}^{\pi^*}(s_1) - V_{\theta^p, \theta^r}^{\pi^*}(s_1)}_{\text{first term}} + \underbrace{V_{\hat{\theta}^p, \hat{\theta}^r}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \theta^r}^{\pi^*}(s_1)}_{\text{second term}} + \underbrace{V_{\hat{\theta}^p, \hat{\theta}^r}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \hat{\theta}^r}^{\pi^*}(s_1)}_{\text{third term}} \end{aligned}$$

**First term.** By assumption, the expected reward under the true parameter satisfies  $\mathbb{E}_{\theta^r}[r(s, a)] \in [0, 1]$ , then  $\mathbb{S}\left(\sum_{t=1}^H \mathbb{E}_{\theta^r}[r(s_t, \pi(s_t))]\right) \leq H$ . Consequently, the first term can be controlled using Lemma 19

$$\begin{aligned} V_{\hat{\theta}^p, \theta^r}^{\pi^*}(s_1) - V_{\theta^p, \theta^r}^{\pi^*}(s_1) &\leq H \sqrt{\text{KL}(P_{\hat{\theta}^p}(s_2, \dots, s_H), P_{\theta^p}(s_2, \dots, s_H))} \\ &\leq H \sqrt{\mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} \left[ \sum_{t=1}^H \psi(\tilde{s}_{t+1})^\top M_{\hat{\theta}^p - \theta^p} \varphi(\tilde{s}_t, \pi^*(\tilde{s}_t)) + Z_{\hat{\theta}^p}^p(\tilde{s}_t, \pi^*(\tilde{s}_t)) - Z_{\theta^p}^p(\tilde{s}_t, \pi^*(\tilde{s}_t)) \right]} \end{aligned}$$

 Using Taylor's expansion, for all  $h \in [H]$ ,  $\exists \hat{\theta}_h \in [\theta^p, \hat{\theta}^p]$  such that:

$$\begin{aligned} \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} &\left[ \psi(\tilde{s}_{t+1})^\top M_{\hat{\theta}^p - \theta^p} \varphi(\tilde{s}_t, \pi^*(\tilde{s}_t)) + Z_{\hat{\theta}^p}^p(\tilde{s}_t, \pi^*(\tilde{s}_t)) - Z_{\theta^p}^p(\tilde{s}_t, \pi^*(\tilde{s}_t)) \right] \\ &= \frac{1}{2} (\hat{\theta}^p - \theta^p)^\top \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} \left[ \nabla_{s_h, \pi^*(s_h)}^2 Z^p(\theta_h) \right] (\hat{\theta}^p - \theta^p) \\ &\leq \frac{\beta^p}{2} \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} \left[ \|\hat{\theta}^p - \theta^p\|_{G_{\tilde{s}_h, \pi^*(\tilde{s}_h)}}^2 \right]. \end{aligned}$$

Define  $u_k \stackrel{\text{def}}{=} \sum_{h=1}^H \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} [(A_i \varphi(\tilde{s}_h, \pi^*(\tilde{s}_h)))_{i \in [d]}]$ , then

$$\begin{aligned} V_{\hat{\theta}^p, \hat{\theta}^x}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \theta^x}^{\pi^*}(s_1) &\leq H \sqrt{\frac{\beta^p}{2} \sum_{h=1}^H \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} [\|\hat{\theta}^p - \theta^p\|_{G_{\tilde{s}_h, \pi^*(\tilde{s}_h)}}^2]} \\ &\leq H \sqrt{\frac{\beta^p}{2}} \|\hat{\theta}^p - \theta^p\|_{\sum_{h=1}^H \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} [G_{\tilde{s}_h, \pi^*(\tilde{s}_h)}]} \\ &\leq H \sqrt{\frac{\beta^p}{2}} \|\hat{\theta}^p - \theta^p\|_{u_k u_k^\top} \\ &\leq H \sqrt{\frac{\beta^p}{2}} \|(\bar{G}_k^p)^{-1/2} u_k u_k^\top (\bar{G}_k^p)^{-1/2}\| \|\hat{\theta}^p - \theta^p\|_{\bar{G}_k^p} \\ &\leq H \sqrt{\frac{\beta^p}{2}} \|u_k\|_{(\bar{G}_k^p)^{-1}} \|\hat{\theta}^p - \theta^p\|_{\bar{G}_k^p} \end{aligned}$$

The third line follows because  $\forall x \in \mathbb{R}^d$ ,  $\|x\|_{\sum_{i=1}^d a_i a_i^\top} \leq \|x\|_{(\sum_{i=1}^d a_i)(\sum_{i=1}^d a_i)^\top}$ , and the last one follows because  $\text{tr}(AB) \leq \text{tr}(A) \text{tr}(B)$  for any two real positive semi-definite matrices  $A$  and  $B$ .

We deduce, with probability at least  $1 - \delta$ :

$$V_{\hat{\theta}^p, \hat{\theta}^x}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \theta^x}^{\pi^*}(s_1) \leq H \sqrt{\frac{\beta^p \beta^p(k, \delta)}{\alpha^p}} \left\| \sum_{h=1}^H \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} [(A_i \varphi(\tilde{s}_h, \pi^*(\tilde{s}_h)))_{i \in [d]}] \right\|_{(\bar{G}_k^p)^{-1}}$$

**Second term.** We have

$$\begin{aligned} V_{\hat{\theta}^p, \hat{\theta}^x}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \theta^x}^{\pi^*}(s_1) &= \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} \left[ \sum_{t=1}^H \frac{\text{Var}^{\theta^x}(r)}{2} B^\top M_{\hat{\theta}^x - \theta^x} \varphi(\tilde{s}_t, \pi^*(\tilde{s}_t)) \right] \\ &= (\hat{\theta}^x - \theta^x)^\top \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} \left[ \sum_{t=1}^H \frac{\text{Var}^{\theta^x}(r)}{2} (A_i \varphi(\tilde{s}_t, \pi^*(\tilde{s}_t)))_{i \in [d]} \right] B \\ &\leq \frac{\sqrt{\beta^x}}{2} \|\hat{\theta}^x - \theta^x\|_{\bar{G}_k^x} \left\| \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} \left[ \sum_{t=1}^H (A_i \varphi(\tilde{s}_t, \pi^*(\tilde{s}_t)))_{i \in [d]} \right] \right\|_{(\bar{G}_k^x)^{-1}} \end{aligned}$$

The last inequality comes from Cauchy-Schwarz. Applying that the norm (sum) makes appear only symmetric matrices times the variances so that we can bound the latter by  $\beta^x$ .

We conclude that with probability at least  $1 - \delta$ ,

$$V_{\hat{\theta}^p, \hat{\theta}^x}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \theta^x}^{\pi^*}(s_1) \leq \frac{\beta^x \sqrt{\beta^x(k, \delta)}}{\sqrt{2\alpha^x}} \left\| \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} \left[ \sum_{t=1}^H (A_i \varphi(\tilde{s}_t, \pi^*(\tilde{s}_t)))_{i \in [d]} \right] \right\|_{(\bar{G}_k^x)^{-1}}$$

We want to write all the norms in the same matrix. Therefore, with probability at least  $1 - \delta$ ,

$$\begin{aligned} V_{\hat{\theta}^p, \hat{\theta}^x}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \theta^x}^{\pi^*}(s_1) &\leq \sqrt{\frac{\beta^x \beta^x(k, \delta) \min\{1, \frac{\alpha^p}{\alpha^x}\}}{2\alpha^x}} \\ &\quad \times \left\| \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} \left[ \sum_{t=1}^H (A_i \varphi(\tilde{s}_t, \pi^*(\tilde{s}_t)))_{i \in [d]} \right] \right\|_{(\bar{G}_k^p)^{-1}} \end{aligned}$$

**Third term.** We have

$$\begin{aligned} V_{\hat{\theta}^p, \hat{\theta}^x, 1}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \tilde{\theta}^x, 1}^{\pi^*}(s_1) &= \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} \left[ \sum_{t=1}^H \frac{\text{Var}^{\theta_j^x}(r)}{2} B^\top M_{\hat{\theta}^x - \tilde{\theta}^x} \varphi(\tilde{s}_t, \pi^*(\tilde{s}_t)) \right] \\ &= \xi_k^\top \mathbb{E}_{(\tilde{s}_t)_{t \in [H]} \sim \hat{\theta}^p | s_1^k} \left[ \sum_{t=1}^H \frac{\text{Var}^{\theta_j^x}(r)}{2} (A_i \varphi(\tilde{s}_t, \pi^*(\tilde{s}_t)))_{i \in [d]} \right] B \end{aligned}$$

Given the normal CDF  $\Phi$ , we obtain that with probability at least  $\Phi(-1)$

$$V_{\hat{\theta}^p, \hat{\theta}^r}^{\pi^*}(s_1) - V_{\hat{\theta}^p, \hat{\theta}^r}^{\pi^*}(s_1) \geq \sqrt{x_k \alpha^r} \left\| \left[ \sum_{t=1}^H \frac{\text{Var}^{\theta^r}(r)}{2} (A_i \varphi(\tilde{s}_t, \pi^*(\tilde{s}_t)))_{i \in [d]} \right] \right\|_{(\bar{G}_k^r)^{-1}}$$

365 Choosing  $x_k \geq \left( H \sqrt{\frac{\beta^p \beta^r(k, \delta)}{\alpha^p \alpha^r}} + \frac{\sqrt{\beta^r \beta^r(k, \delta) \min\{1, \frac{\alpha^p}{\alpha^r}\}}}{2\alpha^r} \right)$  and using Lemma 18, we find that the perturbed value  
 366 function is optimistic with probability at least  $\Phi(-1)$ .

### 367 B.2.2 CONTROLLING THE LEARNING ERROR

368 In this section we see the core difference with optimistic algorithms. On the one hand, optimistic approaches require the  
 369 value function generating the agent's policy to be larger than the optimal one with large probability, and can therefore  
 370 ensure that the learning error is negative. On the other hand, BEF-RLSVI only ensures that the value function is  
 371 optimistic with a constant probability: intuitively when this event holds the learning happens, and if it does not then the  
 372 policy is still close to a good one thanks to the decreasing estimation error.

**Upper bound on  $V_1^*$ .** Let us draw  $(\bar{\xi}_k)_{k \in [K]}$  i.i.d copies of  $(\xi_k)_{k \in [K]}$ . Define the optimism event at episode  $k$ :

$$\bar{O}_k = \{V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) - V_1^*(s_1^k) \geq 0\} \quad (15)$$

we know that  $\mathbb{P}(\bar{O}_k) \geq \Phi(-1)$ . This event provides the upper bound:

$$V_1^*(s_1^k) \leq \mathbb{E}_{\bar{\xi}_k | \bar{O}_k} [V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k)] \quad (16)$$

**Lower bound on  $V_{\hat{\theta}^p, \hat{\theta}^r}$ .** We define this bound with an optimization problem under concentration of the noise. Consider  $\underline{V}_1(s_1^k)$  is the solution of

$$\begin{aligned} \min_{\xi_k} V_{\hat{\theta}^p, \hat{\theta}^r + \xi_k, 1}(s_1^k) \\ \|\xi_k\|_{\bar{G}_k^p} \leq \sqrt{x_k d \log(d/\delta)}, \quad \forall t \in [H] \end{aligned} \quad (17)$$

Under the concentration of our injected noise, we obtain

$$\underline{V}_1(s_1^k) \leq V_{\hat{\theta}^p, \hat{\theta}^r}(s_1^k) \quad (18)$$

**Combining the error bounds.** Combining the upper bound of Equation (16) with the lower bound of Equation (18), we get, with probability at least  $1 - \delta$ :

$$V_1^*(s_1^k) - V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) \leq \mathbb{E}_{\bar{\xi}_k | \bar{O}_k} [V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)]$$

Also, using the tower rule,

$$\begin{aligned} \mathbb{E}_{\bar{\xi}_k} [V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)] \\ = \mathbb{E}_{\bar{\xi}_k | \bar{O}_k} [V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)] \mathbb{P}(\bar{O}_k) + \mathbb{E}_{\bar{\xi}_k | \bar{O}_k^c} [V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)] \mathbb{P}(\bar{O}_k^c) \end{aligned}$$

Therefore,

$$\begin{aligned} V_1^*(s_1^k) - V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) \\ \leq \left( \mathbb{E}_{\bar{\xi}_k} [V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)] - \mathbb{E}_{\bar{\xi}_k | \bar{O}_k^c} [V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)] \mathbb{P}(\bar{O}_k^c) \right) / \mathbb{P}(\bar{O}_k) \\ = \left( \mathbb{E}_{\bar{\xi}_k} [V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)] - \mathbb{E}_{\bar{\xi}_k | \bar{O}_k^c} [V_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)] \mathbb{P}(\bar{O}_k^c) \right) / \mathbb{P}(\bar{O}_k). \end{aligned}$$

373 The last line follows since  $\xi_k$  and  $\bar{\xi}_k$  are i.i.d.

374 The rest of the analysis proceeds similarly to the proof of the reward estimation.

Let us call the argument of the minimum in Equation (17) as  $\underline{\xi}_k$ . Using Lemma 23, we find

$$\begin{aligned}
 & V_{\hat{\theta}^v, \hat{\theta}^r, 1}^\pi(s_1^k) - V_{\hat{\theta}^v, \hat{\theta}^r + \underline{\xi}_k, 1}^\pi(s_1^k) \\
 &= \mathbb{E}_{(\tilde{s}_h)_{1 \leq h \leq H} \sim \pi|\hat{\theta}^v, s_1^k} \left[ \sum_{h=1}^H \frac{\text{Var}_{\tilde{s}_h, \pi}(\tilde{s}_h)(r)}{2} B^\top M_{\hat{\theta}^r - \hat{\theta}^r - \underline{\xi}_k} \varphi(\tilde{s}_h, \pi(\tilde{s}_h)) \right] \\
 &\leq \mathbb{E} \left[ \sum_{h=1}^H \frac{\text{Var}_{\tilde{s}_h, \pi}(\tilde{s}_h)(r)}{2} \|\hat{\theta}^r - \hat{\theta}^r - \underline{\xi}_k\|_{\bar{G}_k^v} \|(B^\top A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^v)^{-1}} \right] \\
 &\leq \|\hat{\theta}^r - \hat{\theta}^r - \underline{\xi}_k\|_{\bar{G}_k^v} \mathbb{E} \left[ \sum_{h=1}^H \frac{\text{Var}_{\tilde{s}_h, \pi}(\tilde{s}_h)(r)}{2} \|(B^\top A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^v)^{-1}} \right] \\
 &\leq \|\tilde{\xi}_k - \underline{\xi}_k\|_{\bar{G}_k^v} \frac{\beta^r}{2} \mathbb{E} \left[ \sum_{h=1}^H \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^v)^{-1}} \right]
 \end{aligned}$$

Then,

$$\begin{aligned}
 & \mathbb{E}_{\tilde{\xi}_k} \left[ V_{\hat{\theta}^v, \hat{\theta}^r, 1}^\pi(s_1^k) - V_{\hat{\theta}^v, \hat{\theta}^r + \underline{\xi}_k, 1}^\pi(s_1^k) \right] \\
 &\leq \frac{\beta^r}{2} \mathbb{E}_{\tilde{\xi}_k} [\|\tilde{\xi}_k - \underline{\xi}_k\|_{\bar{G}_k^v}] \mathbb{E}_{(\tilde{s}_h) \sim \pi|\hat{\theta}^v} \left[ \sum_{h=1}^H \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^v)^{-1}} \right].
 \end{aligned}$$

Also,

$$\begin{aligned}
 & \left| \mathbb{E}_{\xi_k | \bar{O}_k^v} [V_{\hat{\theta}^v, \hat{\theta}^r + \xi_k, 1}^\pi(s_1^k) - V_1(s_1^k)] \right| \\
 &\leq \frac{\beta^r}{2} \mathbb{E}_{\tilde{\xi}_k | \bar{O}_k^v} [\|\tilde{\xi}_k - \underline{\xi}_k\|_{\bar{G}_k^v}] \mathbb{E}_{(\tilde{s}_h) \sim \pi|\hat{\theta}^v} \left[ \sum_{h=1}^H \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^v)^{-1}} \right] \\
 &\leq \frac{\beta^r}{2} \mathbb{E}_{\tilde{\xi}_k} [\|\tilde{\xi}_k - \underline{\xi}_k\|_{\bar{G}_k^v}] \mathbb{E}_{(\tilde{s}_h) \sim \pi|\hat{\theta}^v} \left[ \sum_{h=1}^H \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^v)^{-1}} \right].
 \end{aligned}$$

We have a bound on the expected value of the sum of feature norms in the proof of Lemma 9. Also,

$$\begin{aligned}
 \mathbb{E}_{\tilde{\xi}_k} [\|\tilde{\xi}_k - \underline{\xi}_k\|_{\bar{G}_k^v}] &\leq \mathbb{E}_{\tilde{\xi}_k} [\|\tilde{\xi}_k\|_{\bar{G}_k^v}] + \mathbb{E}_{\tilde{\xi}_k} [\|\underline{\xi}_k\|_{\bar{G}_k^v}] \\
 &\leq \sqrt{\mathbb{E}_{\tilde{\xi}_k} [\|\tilde{\xi}_k\|_{\bar{G}_k^v}^2]} + \sqrt{x_k d \log(d/\delta)} \\
 &\leq \sqrt{x_k d} + \sqrt{x_k d \log(d/\delta)}
 \end{aligned}$$

375 The second line follows from Cauchy-Schwarz and by definition of  $\underline{\xi}_k$ . The last line is due to the fact that  $x_k (\bar{G}_k^v)^{-1} \sim$   
 376  $\mathcal{N}(0, x_k I_d)$ , which implies  $\|\tilde{\xi}_k\|_{\bar{G}_k^v}^2 \sim \mathcal{N}(0, dx_k)$ . We conclude the proof by taking the sum of feature norms from the  
 377 proof of Lemma 9.

We conclude that with probability at least  $1 - 2\delta$ :

$$\begin{aligned}
 & \sum_{k=1}^K V_1^*(s_1^k) - V_{\hat{\theta}^v, \hat{\theta}^r + \tilde{\xi}_k, 1}^\pi(s_1^k) \leq \frac{\beta^r}{\Phi(-1)} (\sqrt{x_k d} + \sqrt{x_k d \log(d/\delta)}) \\
 & \left[ \sqrt{\frac{3d}{\log(2)} \log \left( 1 + \frac{\alpha^r \|\mathbb{A}\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)} \right) \left( 1 + \frac{\alpha^r B_{\varphi, \mathbb{A}} H}{\eta} \right) H d \log(1 + \alpha^r \eta^{-1} B_{\varphi, \mathbb{A}} H)} \right. \\
 & \left. + \sqrt{K H d \log(1 + \alpha^r \eta^{-1} B_{\varphi, \mathbb{A}} H K) \log(e/\delta^2)} \right]
 \end{aligned}$$

378 **Appendix C. Concentrations**

 379 **C.1 Concentration of the transition parameter**

 380 We recall the important concentration of the maximum likelihood estimator for general bilinear exponential families (cf.  
 381 Theorem 1 of Chowdhury et al. (2021)).

**Theorem 13** Suppose  $\{\mathcal{F}_t\}_{t=0}^\infty$  is a filtration such that for each  $t$ , (i)  $s_{t+1}$  is  $\mathcal{F}_t$ -measurable, (ii)  $(s_t, a_t)$  is  $\mathcal{F}_{t-1}$  measurable, and (iii) given  $(s_t, a_t)$ ,  $s_{t+1} \sim P_{\theta^p}^p(\cdot | s_t, a_t)$  according to the exponential family defined by Equation (1). Let  $\hat{\theta}^p(k)$  be the penalized MLE defined by Equation (6), and let  $Z_{s,a}^p(\theta)$  be strictly convex in  $\theta$  for all  $(s, a)$ . Then, for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$ , the following holds uniformly over all  $n \in \mathbb{N}$ :

$$\sum_{t=1}^k \text{KL}_{s_t, a_t}(\hat{\theta}^p(k), \theta^p) + \frac{\eta}{2} \left\| \theta^p - \hat{\theta}^p(k) \right\|_{\mathbb{A}}^2 - \frac{\eta}{2} \|\theta^p\|_{\mathbb{A}}^2 \leq \log \left( \frac{C_{\mathbb{A},k}^p}{\delta} \right),$$

where  $C_{\mathbb{A},k}^p = \left( \int_{\mathbb{R}^d} \exp\left(-\frac{\eta}{2} \|\theta'\|_{\mathbb{A}}^2\right) d\theta' \right) / \left( \int_{\mathbb{R}^d} \exp\left(-\sum_{t=1}^k \text{KL}_{s_t, a_t}(\theta_k, \theta') - \frac{\eta}{2} \|\theta' - \theta_k\|_{\mathbb{A}}^2\right) d\theta' \right)$ . Define  $G_{s,a} \stackrel{\text{def}}{=} (\varphi(s, a)^\top A_i^\top A_j \varphi(s, a))_{i,j \in [d]}$ , we have

$$C_{\mathbb{A},k}^p \leq \det \left( I + \beta^p \eta^{-1} \mathbb{A}^{-1} \sum_{t=1}^k G_{s_t, a_t} \right),$$

 382 where  $\beta^p = \sup_{\theta, s, a} \lambda_{\max}(\mathbb{C}_{s,a}^\theta[\psi(s')])$ .

 383 A proof of this result can be found in the work Chowdhury et al. (2021). We provide an almost similar proof for the  
 384 concentration of rewards in the next section.

**Corollary 14** The previous theorem implies a simple euclidean confidence region. Indeed, with probability at least  $1 - \delta$ , for all  $k \in \mathbb{N}$

$$\left\| \theta^p - \hat{\theta}^p(k) \right\|_{\mathbb{C}_n^p}^2 \leq \frac{2}{\alpha^p} \beta^p(k, \delta),$$

 385 where  $\beta^p(k, \delta) \stackrel{\text{def}}{=} \beta_{(k-1)H}^p(\delta) = \frac{2}{2} B_A^2 + \log(2C_{A,k}^p/\delta)$ .

**Proof** The result follows from the following simple calculations:

$$\begin{aligned} \frac{1}{2} \left\| \theta^p - \hat{\theta}^p(k) \right\|_{\mathbb{C}_k}^2 &= \frac{(\alpha^p)^{-1} \eta}{2} \left\| \theta^p - \hat{\theta}^p(k) \right\|_{\mathbb{A}}^2 + \sum_{\tau=1}^{k-1} \sum_{h=1}^H \frac{1}{2} \left\| \theta^p - \hat{\theta}^p(k) \right\|_{G_{s_h^\tau, a_h^\tau}}^2 \\ &\leq (\alpha^p)^{-1} \left( \frac{\eta}{2} \left\| \theta^p - \hat{\theta}^p(k) \right\|_{\mathbb{A}}^2 + \sum_{\tau=1}^{k-1} \sum_{h=1}^H \text{KL}_{s_h^\tau, a_h^\tau}(\theta_k, \theta) \right). \end{aligned}$$

386

387

 388 **C.2 Concentration of the reward parameter (contribution)**

**Theorem 15** Suppose  $\{\mathcal{F}_t\}_{t=0}^\infty$  is a filtration such that for each  $t$ , (i)  $r(s_t, a_t)$  is  $\mathcal{F}_t$ -measurable, (ii)  $(s_t, a_t)$  is  $\mathcal{F}_{t-1}$  measurable, and (iii) given  $(s_t, a_t)$ ,  $r(s_t, a_t) \sim P_{\theta^x}^x(\cdot | s_t, a_t)$  according to the exponential family defined by (2). Let  $\hat{\theta}^x(k)$  be the penalized MLE defined by Equation (8), and let  $Z_{s,a}^x(\theta)$  be strictly convex in  $\theta$  for all  $(s, a)$ . Then, for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$ , the following holds uniformly over all  $k \in \mathbb{N}$ :

$$\sum_{t=1}^k \text{KL}_{s_t, a_t}(\hat{\theta}^x(k), \theta^x) + \frac{\eta}{2} \left\| \theta^x - \hat{\theta}^x(k) \right\|_{\mathbb{A}}^2 - \frac{\eta}{2} \|\theta^x\|_{\mathbb{A}}^2 \leq \log \left( \frac{C_{\mathbb{A},k}^x}{\delta} \right),$$

where  $C_{\mathbb{A},k}^x = \left( \int_{\mathbb{R}^d} \exp\left(-\frac{\eta}{2} \|\theta'\|_{\mathbb{A}}^2\right) d\theta' \right) / \left( \int_{\mathbb{R}^d} \exp\left(-\sum_{t=1}^k \text{KL}_{s_t, a_t}(\theta_k, \theta') - \frac{\eta}{2} \|\theta' - \theta_k\|_{\mathbb{A}}^2\right) d\theta' \right)$ . Define  $G_{s,a} \stackrel{\text{def}}{=} (\varphi(s, a)^\top A_i^\top A_j \varphi(s, a))_{i,j \in [d]}$ , we have

$$C_{\mathbb{A},k} \leq \det \left( I + \beta^x \eta^{-1} \mathbb{A}^{-1} \sum_{t=1}^k G_{s_t, a_t} \right),$$

389 where  $\beta^x := \|B\|_2^2 \sup_{\theta, s, a} \text{Var}_{s,a}^\theta(r)$ .

390 **Proof** We proceed similar to the proof of Theorem 1 in Chowdhury and Gopalan (2019).

**Step 1: Martingale construction.** First, observe that by assuming strict convexity, the log-partition function  $Z_{s,a}^x$  becomes a Legendre function. Now for the conditional exponential family model, the KL divergence between  $\mathbb{P}_{\theta^x}^r(\cdot | s, a)$  and  $\mathbb{P}_{\theta^{x'}}^r(\cdot | s, a)$  can be expressed as a Bregman divergence associated to  $Z_{s,a}^x$  with the parameters reversed, i.e.

$$\text{KL}_{s,a}(\theta^x, \theta^{x'}) := \text{KL}(P_{\theta^x}(\cdot | s, a), P_{\theta^{x'}}(\cdot | s, a)) = B_{Z_{s,a}^x}(\theta^{x'}, \theta^x).$$

Now, for any  $\lambda \in \mathbb{R}^d$ , we introduce the function  $B_{Z_{n,\alpha}, \theta^x}(\lambda) = B_{Z_{n,\alpha}}(\theta^x + \lambda, \lambda)$  and define

$$M_n^\lambda = \exp \left( \lambda^\top S_n - \sum_{t=1}^n B_{Z_{n_t, a_t}, \theta^x}(\lambda) \right)$$

where  $\forall i \leq d$ , we denote  $(S_n)_i = \sum_{t=1}^n (r(s_t, a_t) - \mathbb{E}_{s_t, a_t}^{\theta^x}[r]) B^\top A_i \varphi(s_t, a_t)$ . Note that  $M_n^\lambda > 0$  and it is  $\mathcal{F}_n$ -measurable. Furthermore, we have for all  $(s, a)$ ,

$$\begin{aligned} & \mathbb{E}_{s,a}^{\theta^x} \left[ \exp \left( \sum_{i=1}^d \lambda_i \left( r(s_t, a_t) - \mathbb{E}_{s_t, a_t}^{\theta^x}[r] \right) B^\top A_i \varphi(s_t, a_t) \right) \right] \\ &= \exp(-\lambda^\top \nabla Z_{s,a}^x(\theta^x)) \int_{\mathcal{S}} \exp \left( \sum_{i=1}^d (\theta_i^x + \lambda_i) B^\top A_i \varphi(s, a) - Z_{s,a}^x(\theta^x) \right) dr \\ &= \exp(Z_{s,a}^x(\theta^x + \lambda) - Z_{s,a}^x(\theta^x) - \lambda^\top \nabla Z_{s,a}^x(\theta^x)) = \exp(B_{Z_{s,a}^x}(\theta^x)) \end{aligned}$$

This implies  $\mathbb{E}[\exp(\lambda^\top S_n) | \mathcal{F}_{n-1}] = \exp(\lambda^\top S_{n-1} + B_{Z_{n_n, a_n}, \theta^x}(\lambda))$  thus  $\mathbb{E}[M_n^\lambda | \mathcal{F}_{n-1}] = M_{n-1}^\lambda$ . Therefore  $\{M_n^\lambda\}_{n=0}^\infty$  is a non-negative martingale adapted to the filtration  $\{\mathcal{F}_n\}_{n=0}^\infty$  and actually satisfies  $\mathbb{E}[M_n^\lambda] = 1$ . For any prior density  $q(\theta)$  for  $\theta$ , we now define a mixture of martingales

$$M_n = \int_{\mathbb{R}^d} M_n^\lambda q(\theta^x + \lambda) d\lambda \quad (19)$$

391 Then  $\{M_n\}_{n=0}^\infty$  is also a non-negative martingale adapted to  $\{\mathcal{F}_n\}_{n=0}^\infty$  and in fact,  $\mathbb{E}[M_n] = 1$ .

**Step 2: Method of mixtures.** Considering the prior density  $\mathcal{N}(0, (\eta\mathbb{A})^{-1})$ , we obtain from (19) that

$$M_n = c_0 \int_{\mathbb{R}^d} \exp \left( \lambda^\top S_n - \sum_{t=1}^n B_{Z_{s_t, a_t}, \theta^x}(\lambda) - \frac{\eta}{2} \|\theta^x + \lambda\|_{\mathbb{A}}^2 \right) d\lambda, \quad (20)$$

where  $c_0 = \frac{1}{\int_{\mathbb{R}^d} \exp(-\frac{\eta}{2} \|\theta'\|_{\mathbb{A}}^2) d\theta'}$ . We now introduce the function  $Z_n^x(\theta) = \sum_{t=1}^n Z_{s_t, a_t}^x(\theta)$ . Note that  $Z_n^x$  is also Legendre function and its associated Bregman divergence satisfies

$$B_{Z_n^x}(\theta', \theta) = \sum_{t=1}^n \left( Z_{s_t, a_t}^x(\theta') - Z_{s_t, a_t}^x(\theta) - (\theta' - \theta)^\top \nabla Z_{s_t, a_t}^x(\theta) \right) = \sum_{t=1}^n B_{Z_{s_t, a_t}^x}(\theta', \theta)$$

Furthermore, we have  $\sum_{t=1}^n B_{Z_{s_t, \alpha_t}^r, \theta^r}(\lambda) = B_{Z_n^r, \theta^r}(\lambda)$ . From the penalized likelihood formula (8), recall that

$$\forall i \leq d, \quad \sum_{t=1}^n \nabla_i Z_{s_t, \alpha_t}^r(\hat{\theta}^r(k)) + \frac{\eta}{2} \nabla_i \|\hat{\theta}^r(k)\|_{\mathbb{A}}^2 = \sum_{t=1}^k r_t B^\top A_i \varphi(s_t, a_t).$$

This yields

$$S_k = \sum_{t=1}^k \left( \nabla Z_{s_t, \alpha_t}^r(\hat{\theta}^r(k)) - \nabla Z_{s_t, \alpha_t}^r(\theta^r) \right) + \eta \mathbb{A} \hat{\theta}^r(k) = \nabla Z_k^r(\hat{\theta}^r(k)) - \nabla Z_k^r(\theta^r) + \eta \mathbb{A} \hat{\theta}^r(k) \quad (21)$$

We now obtain from (20) and (21) that

$$M_k = c_0 \cdot \exp\left(-\frac{\eta}{2} \|\theta^r\|_{\mathbb{A}}^2\right) \int_{\mathbb{R}^d} \exp(\lambda^\top x_k - B_{Z_k, \theta^*}(\lambda) + g_k(\lambda)) d\lambda, \quad (22)$$

where we introduced  $g_k(\lambda) = \frac{\eta}{2} \left( 2\lambda^\top \mathbb{A} \hat{\theta}^r(k) + \|\theta^r\|_{\mathbb{A}}^2 - \|\theta^r + \lambda\|_{\mathbb{A}}^2 \right)$  and  $x_k = \nabla Z_k^r(\hat{\theta}^r(k)) - \nabla Z_k^r(\theta^r)$ .

Now, note that  $\sup_{\lambda \in \mathbb{R}^d} g_k(\lambda) = \frac{\eta}{2} \left\| \theta^r - \hat{\theta}^r(k) \right\|_{\mathbb{A}}^2$ , where the supremum is attained at  $\lambda^* = \hat{\theta}^r(k) - \theta^r$ . We then have

$$\begin{aligned} g_k(\lambda) &= g_n(\lambda) + \sup_{\lambda \in \mathbb{R}^*} g_k(\lambda) - g_k(\lambda^*) \\ &= \frac{\eta}{2} \left\| \hat{\theta}^r(k) - \theta^r \right\|_{\mathbb{A}}^2 + \eta (\lambda - \lambda^*)^\top \mathbb{A} (\theta^r + \lambda^*) + \frac{\eta}{2} \|\theta^r + \lambda^*\|_{\mathbb{A}}^2 - \frac{\eta}{2} \|\theta^r + \lambda\|_{\mathbb{A}}^2 \\ &= B_{Z_0^r}(\theta^r, \hat{\theta}^r(k)) + (\lambda - \lambda^*)^\top \nabla Z_0^r(\theta^r + \lambda^*) + Z_0^r(\theta^r + \lambda^*) - Z_0^r(\theta^r + \lambda) \end{aligned} \quad (23)$$

where we have introduced the Legendre function  $Z_0^r(\theta) = \frac{\eta}{2} \|\theta\|_{\mathbb{A}}^2$ . We now have from (27) that

$$\begin{aligned} &\sup_{\lambda \in \mathbb{R}^d} (\lambda^\top x_n - B_{Z_n^r, \theta^r}(\lambda)) \\ &= B_{Z_n^r, \theta^r}^*(x_n) = B_{Z_n^r, \theta^r}^* \left( \nabla Z_n^r(\hat{\theta}^r(n)) - \nabla Z_n^r(\theta^r) \right) = B_{Z_n^r}(\theta^r, \hat{\theta}^r(n)). \end{aligned}$$

Further, any optimal  $\lambda$  must satisfy

$$\nabla Z_n^r(\theta^r + \lambda) - \nabla Z_n^r(\theta^r) = x_n \implies \nabla Z_n^r(\theta^r + \lambda) = \nabla Z_n^r(\hat{\theta}^r(n)).$$

One possible solution is  $\lambda = \lambda^*$ . Now, since  $Z_n^r$  is strictly convex, the supremum is indeed attained at  $\lambda = \lambda^*$ . We then have

$$\begin{aligned} &\lambda^\top x_n - B_{Z_n^r, \theta^r}(\lambda) \\ &= \lambda^\top x_n - B_{Z_n^r, \theta^r}(\lambda) + B_{Z_n^r}(\theta^r, \hat{\theta}^r(n)) - (\lambda^* x_n - B_{Z_n^r, \theta^r}(\lambda^*)) \\ &= B_{Z_n^r}(\theta^r, \hat{\theta}^r(n)) + (\lambda - \lambda^*)^\top \nabla Z_n^r(\theta^r + \lambda^*) + B_{Z_n^r, \theta^*}(\lambda^*) - B_{Z_n^r, \theta^*}(\lambda) \\ &\quad - (\lambda - \lambda^*)^\top \nabla Z_n^r(\theta^r) \\ &= B_{Z_n^r}(\theta^r, \hat{\theta}^r(n)) + (\lambda - \lambda^*)^\top \nabla Z_n^r(\theta^r + \lambda^*) + Z_n^r(\theta^r + \lambda^*) - Z_n^r(\theta^r + \lambda) \end{aligned} \quad (24)$$



Plugging Equation (23) and Equation (24) in Equation (22), we obtain

$$\begin{aligned}
 M_n &= c_0 \cdot \exp \left( \sum_{j \in \{0, n\}} B_{Z_j^x}(\theta^x, \theta_j) - \frac{\eta}{2} \|\theta^x\|_A^2 \right) \\
 &\quad \times \int_{\mathbb{R}^d} \exp \left( \sum_{j \in \{0, n\}} \left( (\lambda - \lambda^*)^\top \nabla Z_j^x(\theta^x + \lambda^*) + Z_j^x(\theta^x + \lambda^*) - Z_j^x(\theta^x + \lambda) \right) \right) d\lambda \\
 &= c_0 \cdot \exp \left( \sum_{j \in \{0, n\}} B_{Z_j^x}(\theta^x, \hat{\theta}^x(n)) - \frac{\eta}{2} \|\theta^x\|^2 \right) \\
 &\quad \times \exp \left( - \sum_{j \in \{0, n\}} \left( (\theta^x + \lambda^*)^\top \nabla Z_j^x(\theta^x + \lambda^*) - Z_j^x(\theta^x + \lambda^*) \right) \right) \\
 &\quad \times \int_{\mathbb{R}^d} \exp \left( \sum_{j \in \{0, n\}} \left( (\theta^x + \lambda)^\top \nabla Z_j^x(\theta^x + \lambda^*) - Z_j^x(\theta^x + \lambda) \right) \right) d\lambda \\
 &= \frac{c_0}{c_n} \exp \left( \sum_{j \in \{0, n\}} B_{Z_j^x}(\theta^x, \hat{\theta}^x(n)) - \frac{\eta}{2} \|\theta^x\|_A^2 \right) \\
 &\quad \times \frac{\int_{\mathbb{R}^d} \exp \left( \sum_{j \in \{0, n\}} \left( (\theta^x + \lambda)^\top \nabla Z_j^x(\theta^x + \lambda^*) - Z_j^x(\theta^x + \lambda) \right) \right) d\lambda}{\int_{\mathbb{R}^d} \exp \left( \sum_{j \in \{0, n\}} \left( (\theta')^\top \nabla Z_j^x(\theta^x + \lambda^*) - Z_j^x(\theta') \right) \right) d\theta'} \\
 &= \frac{c_0}{c_n} \cdot \exp \left( B_{Z_n}(\theta^x, \hat{\theta}^x(n)) + B_{Z_0}(\theta^x, \hat{\theta}^x(n)) - \frac{\eta}{2} \|\theta^x\|_A^2 \right),
 \end{aligned}$$

where we introduced  $c_n = \frac{\exp(\sum_{j \in \{0, n\}} ((\theta^x + \lambda^*)^\top \nabla Z_j^x(\theta^x + \lambda^*) - Z_j^x(\theta^x + \lambda^*)))}{\int_{\mathbb{R}^d} \exp(\sum_{j \in \{0, n\}} ((\theta')^\top \nabla Z_j^x(\theta^x + \lambda^*) - Z_j^x(\theta'))) d\theta'}$ . Since  $\lambda^* = \hat{\theta}^x(n) - \theta^x$ , we have

$$c_n = \frac{1}{\int_{\mathbb{R}^d} \exp \left( - \sum_{j \in \{0, n\}} B_{Z_j^x}(\theta', \theta^x + \lambda^*) \right) d\theta'} = \frac{1}{\int_{\mathbb{R}^d} \exp \left( - \sum_{t=1}^n B_{Z_{s_t, a_t}}(\theta', \hat{\theta}^x(n)) - \frac{\eta}{2} \|\theta' - \hat{\theta}^x(n)\|_{\mathbb{A}'}^2 \right) d\theta'}$$

Therefore, we have from (5) that

$$C_{A, n} := \frac{c_n}{c_0} = \frac{\int_{\mathbb{R}^d} \exp \left( - \frac{\eta}{2} \|\theta'\|_A^2 \right) d\theta'}{\int_{\mathbb{R}^d} \exp \left( - \sum_{t=1}^n \text{KL}_{s_t, a_t}(\hat{\theta}^x(n), \theta') - \frac{\eta}{2} \|\theta' - \hat{\theta}^x(n)\|_{\mathbb{A}}^2 \right) d\theta'}$$

An application of Markov's inequality now yields

$$\mathbb{P} \left[ \sum_{t=1}^n \text{KL}_{s_t, a_t}(\hat{\theta}^x(n), \theta^x) + \frac{\eta}{2} \|\theta^x - \hat{\theta}^x(n)\|_{\mathbb{A}}^2 - \frac{\eta}{2} \|\theta^x\|_{\mathbb{A}}^2 \geq \log \left( \frac{C_{A, n}}{\delta} \right) \right] = \mathbb{P} \left[ M_n \geq \frac{1}{\delta} \right] \leq \delta \mathbb{E} [M_n] = \delta$$

**Step 3: A stopped martingale and its control.** Let  $N$  be a stopping time with respect to the filtration  $\{\mathcal{F}_n\}_{n=0}^\infty$ . Now, by the martingale convergence theorem,  $M_\infty = \lim_{n \rightarrow \infty} M_n$  is almost surely well-defined, and thus  $M_N$  is well-defined as well irrespective of whether  $N < \infty$  or not. Let  $Q_n = M_{\min\{N, n\}}$  be a stopped version of  $\{M_n\}_n$ . Then an application of Fatou's lemma yields

$$\mathbb{E} [M_N] = \mathbb{E} \left[ \liminf_{n \rightarrow \infty} Q_n \right] \leq \liminf_{n \rightarrow \infty} \mathbb{E} [Q_n] = \liminf_{n \rightarrow \infty} \mathbb{E} [M_{\min\{N, n\}}] \leq 1,$$

392 since the stopped martingale  $\{M_{\min\{N,n\}}\}_{n \geq 1}$  is also a martingale. Therefore, by the properties of  $M_n$ , (12) also holds  
 393 for any random stopping time  $N < \infty$ . To complete the proof, we now employ a random stopping time construction as  
 394 in Abbasi-Yadkori et al. (2011)

We define a random stopping time  $N$  by

$$N = \min \left\{ n \geq 1 : \sum_{t=1}^n \text{KL}_{s_t, a_t}(\hat{\theta}^x(n), \theta^x) + \frac{\eta}{2} \|\theta^x - \hat{\theta}^x(n)\|_A^2 - \frac{\eta}{2} \|\theta^x\|_A^2 \geq \log \left( \frac{C_{A,n}}{\delta} \right) \right\}$$

with  $\min\{\emptyset\} := \infty$  by convention. We then have

$$\mathbb{P} \left[ \exists n \geq 1, \sum_{t=1}^n \text{KL}_{s_t, a_t}(\hat{\theta}^x(n), \theta^x) + \frac{\eta}{2} \|\theta^x - \hat{\theta}^x(n)\|_A^2 - \frac{\eta}{2} \|\theta^x\|_A^2 \geq \log \left( \frac{C_{A,n}}{\delta} \right) \right] = \mathbb{P}[N < \infty] \leq \delta,$$

395 which concludes the proof of the first part.

396

**Proof of second part: upper bound on  $C_{A,n}$ .** First, we have for some  $\tilde{\theta} \in [\hat{\theta}^x(n), \theta']_\infty$  that

$$\text{KL}_{s,a}(\hat{\theta}^x(n), \theta') = \frac{1}{2} \sum_{i,j=1}^d (\theta' - \hat{\theta}^x(n))_i \text{Var}_{s,a}^\theta(r) \times \varphi(s,a)^\top A_i^\top B B^\top A_j \varphi(s,a) (\theta' - \hat{\theta}^x(n))_j \quad (25)$$

Now (25) implies that

$$\begin{aligned} \sum_{t=1}^n \text{KL}_{s_t, a_t}(\hat{\theta}^x(n), \theta') &\leq \frac{\beta}{2} \sum_{t=1}^n \sum_{i,j=1}^d (\theta' - \hat{\theta}^x(n))_i \varphi(s_t, a_t)^\top A_i^\top A_j \varphi(s_t, a_t) (\theta' - \hat{\theta}^x(n))_j \\ &= \frac{\beta^x}{2} \|\theta' - \hat{\theta}^x(n)\|_{\sum_{t=1}^n G_{s_t, a_t}}, \end{aligned}$$

where  $\beta^x := \lambda_{\max}(B B^\top) \times \sup_{\theta, s, a} \text{Var}_{s,a}^\theta(r)$  and  $\forall i, j \leq d$ ,  $(G_{s,a})_{i,j} := \varphi(s,a)^\top A_i^\top A_j \varphi(s,a)$ . Therefore, we obtain

$$\begin{aligned} C_{A,n} &\leq \frac{\int_{\mathbb{R}^d} \exp\left(-\frac{\eta}{2} \|\theta'\|_A^2\right) d\theta'}{\int_{\mathbb{R}^d} \exp\left(-\frac{1}{2} \|\theta' - \hat{\theta}^x(n)\|_{(\beta^x \sum_{t=1}^n G_{s_t, a_t} + \eta A)}^2\right) d\theta'} \\ &= \frac{(2\pi)^{d/2}}{\det(\eta A)^{1/2}} \times \frac{\det(\beta^x \sum_{t=1}^n G_{s_t, a_t} + \eta A)^{1/2}}{(2\pi)^{d/2}} = \det\left(I + \beta^x \eta^{-1} A^{-1} \sum_{t=1}^n G_{s_t, a_t}\right), \end{aligned}$$

397 which completes the proof of the second part.

398

399

**Corollary 16** Here also, the theorem implies a euclidean control. With probability at least  $1 - \delta$  uniformly over  $k \in \mathbb{N}$

$$\|\theta^x - \hat{\theta}^x(k)\|_{G_k^x}^2 \leq \frac{2}{\alpha^x} \beta^x(k, \delta),$$

400 where  $\beta^x(k, \delta) \stackrel{\text{def}}{=} \beta_{(k-1)H}^x(\delta) = \frac{2}{2} B_A^2 + \log(2C_{A,k}^x/\delta)$ .

### 401 C.3 Gaussian concentration and anti-concentration

**Lemma 17 (Gaussian concentration, ref. Appendix A in Abeille and Lazaric (2017))** Let  $\bar{\xi}_{tk} \sim \mathcal{N}(0, H\nu_k(\delta)\Sigma_{tk}^{-1})$ . For any  $\delta > 0$ , with probability  $1 - \delta$

$$\|\bar{\xi}_{tk}\|_{\Sigma_{tk}} \leq c\sqrt{Hd\nu_k(\delta)\log(d/\delta)} \quad (26)$$

402 for some absolute constant  $c$ .

**Lemma 18 (Gaussian anti-concentration, ref. Appendix A in Abeille and Lazaric (2017))** Let  $\xi \sim \mathcal{N}(0, I_d)$ , for any  $u \in \mathbb{R}^d$  with  $\|u\| = 1$ , we have:

$$\mathbb{P}(u^\top \xi \geq 1) \geq \Phi(-1),$$

403 where  $\Phi$  is the normal CDF.

404 Thanks to lower bounds on the error function, we have the following bound on the probability of anti-concentration  
 405  $\Phi(-1) \geq 1/(4\sqrt{e\pi})$ .

## 406 Appendix D. Technical results

### 407 D.1 A transportation lemma

408 For any function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , we define its span as  $\mathbb{S}(f) := \max_{x \in \mathcal{X}} f(x) - \min_{x \in \mathcal{X}} f(x)$ . For a probability  
 409 distribution  $P$  supported on the set  $\mathcal{X}$ , let  $\mathbb{E}_P[f] := \mathbb{E}_P[f(X)]$  and  $\mathbb{V}_P[f] := \mathbb{V}_P[f(X)] = \mathbb{E}_P[f(X)^2] - \mathbb{E}_P[f(X)]^2$   
 410 denote the mean and variance of the random variable  $f(X)$ , respectively. We now state the following transportation  
 411 inequalities, which can be adapted from Boucheron et al. (2013) (Lemma 4.18).

**Lemma 19 (Transportation inequalities)** Assume  $f$  is such that  $\mathbb{S}(f)$  and  $\mathbb{V}_P[f]$  are finite. Then it holds

$$\begin{aligned} \forall Q \lll P, \quad \mathbb{E}_Q[f] - \mathbb{E}_P[f] &\leq \sqrt{2\mathbb{V}_P[f]\text{KL}(Q, P)} + \frac{2\mathbb{S}(f)}{3}\text{KL}(Q, P) \\ \forall Q \lll P, \quad \mathbb{E}_P[f] - \mathbb{E}_Q[f] &\leq \sqrt{2\mathbb{V}_P[f]\text{KL}(Q, P)} \end{aligned}$$

### 412 D.2 Bregman divergence

For a Legendre function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$ , the Bregman divergence between  $\theta', \theta \in \mathbb{R}^d$  associated with  $F$  is defined as  
 $B_F(\theta', \theta) := F(\theta') - F(\theta) - (\theta' - \theta)^\top \nabla F(\theta)$ . Now, for any fixed  $\theta \in \mathbb{R}^d$ , we introduce the function

$$B_{F,\theta}(\lambda) := B_F(\theta + \lambda, \theta) = F(\theta + \lambda) - F(\theta) - \lambda^\top \nabla F(\theta).$$

It then follows that  $B_{F,\theta}$  is a convex function, and we define its dual as

$$B_{F,\theta}^*(x) = \sup_{\lambda \in \mathbb{R}^d} (\lambda^\top x - B_{F,\theta}(\lambda))$$

We have for any  $\theta, \theta' \in \mathbb{R}^d$ :

$$B_F(\theta', \theta) = B_{F,\theta'}^*(\nabla F(\theta) - \nabla F(\theta')) \quad (27)$$

To see this, we observe that

$$\begin{aligned} &B_{F,\theta'}^*(\nabla F(\theta) - \nabla F(\theta')) \\ &= \sup_{\lambda \in \mathbb{R}^d} \lambda^\top (\nabla F(\theta) - \nabla F(\theta')) - [F(\theta' + \lambda) - F(\theta') - \lambda^\top \nabla F(\theta')] \\ &= \sup_{\lambda \in \mathbb{R}^d} \lambda^\top \nabla F(\theta) - F(\theta' + \lambda) + F(\theta'). \end{aligned}$$

Now an optimal  $\lambda$  must satisfy  $\nabla F(\theta) = \nabla F(\theta' + \lambda)$ . One possible choice is  $\lambda = \theta - \theta'$ . Since, by definition,  $F$  is strictly convex, the supremum will indeed be attained at  $\lambda = \theta - \theta'$ . Plug-in this value, we obtain

$$B_{F,\theta'}^*(\nabla F(\theta) - \nabla F(\theta')) = (\theta - \theta')^\top \nabla F(\theta) - F(\theta) + F(\theta') = B_F(\theta', \theta).$$

413 Note that (27) holds for any convex function  $F$ . Only difference is that, in this case,  $B_F(\cdot, \cdot)$  will not correspond to the  
 414 Bregman divergence.

415 **D.3 Properties of the bilinear exponential family**

416 In this section, we detail some useful results related to exponential families in our model.

## 417 D.3.1 DERIVATIVES

**Lemma 20** (*Gradients*) *We provide the derivatives of the log-partitions in closed form. As usual with exponential families, these are intimately linked to moments of the random variable. We have:*

$$(\nabla_i Z_{s,a}^p)(\theta) = \mathbb{E}_{s,a}^\theta [\psi(s')]^\top A_i \varphi(s, a).$$

And

$$(\nabla_i Z_{s,a}^r)(\theta) = \mathbb{E}_{s,a}^\theta [r] B^\top A_i \varphi(s, a).$$

**Proof** We prove the lemma as follows

$$\begin{aligned} (\nabla_i Z_{s,a}^p)(\theta) &= \int_{\mathcal{S}} \psi(s')^\top A_i \varphi(s, a) \frac{\exp\left(\sum_{i=1}^d \theta_i \psi(s')^\top A_i \varphi(s, a)\right)}{\int_{\mathcal{S}} \exp\left(\sum_{i=1}^d \theta_i \psi(s')^\top A_i \varphi(s, a)\right) ds'} ds' \\ &= \mathbb{E}_{s,a}^\theta [\psi(s')]^\top A_i \varphi(s, a) \\ (\nabla_i Z_{s,a}^r)(\theta) &= \int_{\mathcal{S}} r B^\top A_i \varphi(s, a) \frac{\exp\left(r \sum_{i=1}^d \theta_i B^\top A_i \varphi(s, a)\right)}{\int_{\mathcal{S}} \exp\left(r \sum_{i=1}^d \theta_i B^\top A_i \varphi(s, a)\right) dr} dr \\ &= \mathbb{E}_{s,a}^\theta [r] B^\top A_i \varphi(s, a) \end{aligned}$$

418

419

**Lemma 21** (*Hessians*) *The entries of the Hessians of the log partition functions are given by*

$$(\nabla_{i,j}^2 Z_{s,a}^p)(\theta) = \varphi(s, a)^\top A_i^\top \mathbb{C}_{s,a}^\theta [\psi(s')] A_j \varphi(s, a),$$

420 where  $\mathbb{C}_{s,a}^\theta [\psi(s')] \stackrel{\text{def}}{=} \mathbb{E}_{s,a}^\theta [\psi(s') \psi(s')^\top] - \mathbb{E}_{s,a}^\theta [\psi(s')] \mathbb{E}_{s,a}^\theta [\psi(s')^\top]$ .

Similarly,

$$(\nabla_{i,j}^2 Z_{s,a}^r)(\theta) = \text{Var}_{s,a}^\theta(r) \times \varphi(s, a)^\top A_i^\top B B^\top A_j \varphi(s, a),$$

421 where  $\text{Var}_{s,a}^\theta(r) \stackrel{\text{def}}{=} \left(\mathbb{E}_{s,a}^\theta [r^2] - \mathbb{E}_{s,a}^\theta [r]^2\right)$  is the variance of the reward under  $\theta$ .

**Proof** We prove these formulas by differentiating under the integral sign.

$$\begin{aligned} (\nabla_{i,j}^2 Z_{s,a}^p)(\theta) &= \int_{\mathcal{S}} \psi(s')^\top A_i \varphi(s, a) \psi(s')^\top A_j \varphi(s, a) \frac{\exp\left(\sum_{i=1}^d \theta_i \psi(s')^\top A_i \varphi(s, a)\right)}{\int_{\mathcal{S}} \exp\left(\sum_{i=1}^d \theta_i \psi(s')^\top A_i \varphi(s, a)\right) ds'} ds' \\ &\quad - \int_{\mathcal{S}} \psi(s')^\top A_i \varphi(s, a) \frac{\exp\left(\sum_{i=1}^d \theta_i \psi(s')^\top A_i \varphi(s, a)\right)}{\int_{\mathcal{S}} \exp\left(\sum_{i=1}^d \theta_i \psi(s')^\top A_i \varphi(s, a)\right) ds'} ds' (\nabla_j Z_{s,a})(\theta) \\ &= \mathbb{E}_{s,a}^\theta \left[ \psi(s')^\top A_i \varphi(s, a) \psi(s')^\top A_j \varphi(s, a) \right] \\ &\quad - \mathbb{E}_{s,a}^\theta \left[ \psi(s')^\top A_i \varphi(s, a) \right] \mathbb{E}_{s,a}^\theta \left[ \psi(s')^\top A_j \varphi(s, a) \right] \\ &= \varphi(s, a)^\top A_i^\top \left( \mathbb{E}_{s,a}^\theta [\psi(s') \psi(s')^\top] - \mathbb{E}_{s,a}^\theta [\psi(s')] \mathbb{E}_{s,a}^\theta [\psi(s')^\top] \right) A_j \varphi(s, a) \\ &= \varphi(s, a)^\top A_i^\top \mathbb{C}_{s,a}^\theta [\psi(s')] A_j \varphi(s, a), \end{aligned}$$

where we introduce in the last line the  $p \times p$  covariance matrix given by

$$C_{s,a}^\theta [\psi(s')] = \mathbb{E}_{s,a}^\theta [\psi(s') \psi(s')^\top] - \mathbb{E}_{s,a}^\theta [\psi(s')] \mathbb{E}_{s,a}^\theta [\psi(s')^\top]$$

422 The proof of the form of the Hessian for the reward partition function follows the same steps as above. ■  
 423

**Lemma 22** (*KL Divergences*) For any two  $\theta, \theta'$  and for some pair  $(s, a)$ ,

$$\exists \tilde{\theta} \in [\theta, \theta']_\infty, \quad \text{KL}(P_\theta^p(\cdot | s, a), P_{\theta'}^p(\cdot | s, a)) = \frac{1}{2} (\theta - \theta')^\top (\nabla^2 Z_{s,a}^p)(\tilde{\theta}) (\theta - \theta'),$$

424 where  $[\theta, \theta']_\infty$  denotes the  $d$ -dimensional hypercube joining  $\theta$  to  $\theta'$ .

Similarly

$$\exists \tilde{\theta} \in [\theta, \theta']_\infty, \quad \text{KL}(P_\theta^r(\cdot | s, a), P_{\theta'}^r(\cdot | s, a)) = \frac{1}{2} (\theta - \theta')^\top (\nabla^2 Z_{s,a}^r)(\tilde{\theta}) (\theta - \theta').$$

**Proof** We start by writing:

$$\log \left( \frac{P_\theta^p(s' | s, a)}{P_{\theta'}^p(s' | s, a)} \right) = \sum_{i=1}^d (\theta_i - \theta'_i) \psi(s')^\top A_i \varphi(s, a) - Z_{s,a}^p(\theta) + Z_{s,a}^p(\theta'),$$

then

$$\begin{aligned} \text{KL}(P_\theta^p(\cdot | s, a), P_{\theta'}^p(\cdot | s, a)) &= \sum_{i=1}^d (\theta_i - \theta'_i) \mathbb{E}_{s,a}^\theta [\psi(s')]^\top A_i \varphi(s, a) - Z_{s,a}^p(\theta) + Z_{s,a}^p(\theta') \\ &= \frac{1}{2} (\theta - \theta')^\top (\nabla^2 Z_{s,a}^p)(\tilde{\theta}) (\theta - \theta'), \end{aligned}$$

425 where in the last line, we used, by a Taylor expansion, that  $Z_{s,a}(\theta') = Z_{s,a}(\theta) + (\nabla Z_{s,a}(\theta))^\top (\theta' - \theta) + \frac{1}{2}(\theta - \theta')^\top (\nabla^2 Z_{s,a}(\tilde{\theta})) (\theta - \theta')$  for some  $\tilde{\theta} \in [\theta, \theta']_\infty$ .  
 426

427 The proof of the form of the KL divergence for the reward follows the same steps as above. ■  
 428

### 429 D.3.2 A TRANSPORTATION LEMMA FOR REWARDS

**Lemma 23** We provide a closed-form formula for the difference of expected rewards under two distinct parameters:

$$\exists \theta_3 \in [\theta_1, \theta_2], \quad \mathbb{E}_{s,a}^{\theta_1} [r] = \mathbb{E}_{s,a}^{\theta_2} [r] + \frac{\text{Var}_{s,a}^{\theta_3}(r)}{2} B^\top M_{\theta_1 - \theta_2} \varphi(s, a)$$

**Proof** Let's recall the gradient of the reward log partition function:

$$(\nabla_i Z_{s,a}^r)(\theta^r) = \mathbb{E}_{s,a}^{\theta^r} [r] B^\top A_i \varphi(s, a)$$

then for all  $\theta^{r'}$  we have:

$$\mathbb{E}_{s,a}^{\theta^{r'}} [r] = \frac{1}{B^\top M_{\theta^{r'}} \varphi(s, a)} \nabla_i Z_{s,a}^r(\theta^r)^\top \theta^{r'}$$

Let  $\theta_1, \theta_2 \in \mathbb{R}^d$ , using Taylor-Cauchy's formula there exists  $\theta_3 \in [\theta_1, \theta_2]$  such that:

$$\mathbb{E}_{s,a}^{\theta_1} [r] = \mathbb{E}_{s,a}^{\theta_2} [r] + \frac{1}{2B^\top M_{\theta^{r'}} \varphi(s, a)} (\theta_1 - \theta_2)^\top \nabla^2 Z_{s,a}^r(\theta_3)^\top \theta^{r'}$$

We know that  $(\nabla_{i,j}^2 Z_{s,a}^r)(\theta) = \text{Var}_{s,a}^\theta(r) \times \varphi(s, a)^\top A_i^\top B B^\top A_j \varphi(s, a)$ , choosing  $\theta^{r'} = \theta_1 - \theta_2$  we find:

$$\mathbb{E}_{s,a}^{\theta_1} [r] = \mathbb{E}_{s,a}^{\theta_2} [r] + \frac{\text{Var}_{s,a}^{\theta_3}(r)}{2} B^\top M_{\theta_1 - \theta_2} \varphi(s, a).$$

430  
 431

432 **D.4 Elliptical potentials and elliptical lemma**

## 433 D.4.1 ELLIPTICAL LEMMA

434 Here we show a lemma that is popular for regret control in linear MDPs and linear Bandits.

 435 First, consider the notations:  $G_{s,a} := (\varphi(s, a)^\top A_i^\top A_j \varphi(s, a))_{1 \leq i, j \leq d}$ ,  $\bar{G}_n^e \equiv \bar{G}_{(k-1)H}^e := G_n + (\alpha^e)^{-1} \eta A$ , and

 436  $G_n \equiv G_{(k-1)H} := \sum_{\tau=1}^{k-1} \sum_{h=1}^H G_{s_\tau^e, a_h^e}$ . Where  $e$  represents either  $r$  or  $p$ , we omit the superscript  $e$  w.l.o.g in the
 437 rest of this section.

**Lemma 24** (Elliptical lemma and variant for bounded potentials) *Let  $c \in \mathbb{R}^+$ , we can bound the sum of feature norms as follows*

$$\sum_{t=1}^T \min\left\{c, \sum_{h=1}^H \left\| \bar{G}_n^{-1/2} G_{s,a} \bar{G}_n^{-1/2} \right\| \right\} \leq \frac{c}{\log(1+c)} d \log(1 + \alpha \eta^{-1} B_{\varphi, \mathbb{A}} n).$$

 438 where  $B_{\varphi, \mathbb{A}} := \sup_{s,a} \|\mathbb{A}^{-1} G_{s,a}\|$ .

Further, we have

$$\sum_{t=1}^T \sum_{h=1}^H \left\| \bar{G}_n^{-1/2} G_{s,a} \bar{G}_n^{-1/2} \right\| \leq 2d \log(1 + \alpha \eta^{-1} B_{\varphi, \mathbb{A}} n) + \frac{3dH}{\log(2)} \log\left(1 + \frac{\alpha \|A\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)}\right)$$

**Proof** First we have

$$\begin{aligned} \left\| \bar{G}_n^{-1/2} G_{s,a} \bar{G}_n^{-1/2} \right\| &= \sqrt{\text{tr}(\bar{G}_n^{-1/2} G_{s,a} \bar{G}_n^{-1/2} \bar{G}_n^{-1/2} G_{s,a} \bar{G}_n^{-1/2})} \\ &\leq \text{tr}(\bar{G}_n^{-1/2} G_{s,a} \bar{G}_n^{-1/2}) = \text{tr}(\bar{G}_n^{-1} G_{s,a}) = \text{tr}(\mathbf{a}_h^\top \bar{G}_n^{-1} \mathbf{a}_h) \end{aligned}$$

 439 the last line is because  $G_{s,a} = \mathbf{a}_h \mathbf{a}_h^\top$ , where  $\mathbf{a}_h = (A_i \varphi(s_h, a_h))_{i \in [d]}$ .

**First result.** Consider  $h \in [H]$ , denote  $(\lambda_{h,i})_{i \in [d]}$  the eigenvalues of  $\mathbf{a}_h^\top \bar{G}_n^{-1} \mathbf{a}_h$ .  $\bar{G}_n$  is positive definite hence  $\lambda_{h,i} > 0, \forall h, i$ , then

$$\begin{aligned} \min\left\{c, \sum_{h=1}^H \text{tr}(\mathbf{a}_h^\top \bar{G}_n^{-1} \mathbf{a}_h)\right\} &= \min\left\{c, \sum_{h=1}^H \sum_{i=1}^d \lambda_{h,i}\right\} \\ &\leq \frac{c}{\log(1+c)} \sum_{h=1}^H \sum_{i=1}^d \log(1 + \lambda_{h,i}) \quad (\log \text{ is concave}) \\ &\leq \frac{c}{\log(1+c)} \sum_{h=1}^H \log\left(\prod_{i=1}^d 1 + \lambda_{h,i}\right) = \frac{c}{\log(1+c)} \sum_{h=1}^H \log \det(I + \mathbf{a}_h^\top \bar{G}_n^{-1} \mathbf{a}_h) \\ &\leq \frac{c}{\log(1+c)} \log\left(\frac{\det(\bar{G}_n + \sum_{h=1}^H G_{s_h, a_h})}{\det(\bar{G}_n)}\right) \end{aligned}$$

where the last line follows from the matrix determinant lemma:

$$\det(\bar{G}_n + \mathbf{a}_h \mathbf{a}_h^\top) = \det(I + \mathbf{a}_h^\top \bar{G}_n^{-1} \mathbf{a}_h) \det(\bar{G}_n)$$

Therefore:

$$\sum_{t=1}^T \min\left\{c, \sum_{h=1}^H \left\| \bar{G}_n^{-1} G_{s_h^t, a_h^t} \right\| \right\} \leq \frac{c}{\log(1+c)} \sum_{t=1}^T \log \frac{\det(\bar{G}_{n+H})}{\det(\bar{G}_n)},$$

We can now control the R.H.S. of the above equation, as

$$\begin{aligned}
 \sum_{t=1}^T \log \frac{\det(\bar{G}_{n+H})}{\det(\bar{G}_n)} &= \sum_{t=1}^T \log \frac{\det(\bar{G}_{tH})}{\det(\bar{G}_{(t-1)H})} = \log \frac{\det(\bar{G}_{TH})}{\det(\bar{G}_0)} \\
 &= \log \frac{\det(\bar{G}_N)}{\det((\alpha^p)^{-1} \eta \mathbb{A})} = \log \det(I + \alpha \eta^{-1} \mathbb{A}^{-1} G_N) \\
 &\leq d \log \left( 1 + \frac{\alpha^p \eta^{-1}}{d} \text{tr}(\mathbb{A}^{-1} G_N) \right) \quad (\text{Trace-determinant (or AM-GM) inequality}) \\
 &\leq d \log(1 + \alpha^p \eta^{-1} B_{\varphi, \mathbb{A}} n)
 \end{aligned}$$

440 This concludes the proof of the first result.

441 **Second result.** First, we have  $\sup_{s,a} \|G_{s,a}\|_2 \leq \|A\|_2 B_{\varphi, \mathbb{A}}$ .

Fix an episode  $k \in [K]$ ,  $n = (k-1)H$ , using Lemma 26, we know that the number of times  $h \in [H]$  such that  $\|\bar{G}_n^{-1} G_{s_h, a_h}\| \geq 1$  is smaller than  $\frac{3d}{\log(2)} \log \left( 1 + \frac{\alpha \|A\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)} \right)$ . Let us call  $\mathcal{T}_k := \{h \in [H] \mid \|\bar{G}_{(k-1)H}^{-1} G_{s_h, a_h}\| \leq 1\}$ , then

$$\sum_{t=1}^T \sum_{h=1}^H \left\| \bar{G}_n^{-1} G_{s_h^t, a_h^t} \right\| \leq \frac{3d}{\log(2)} \log \left( 1 + \frac{\alpha \|A\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)} \right) + \sum_{h \in \mathcal{T}_k} \min\{1, \left\| \bar{G}_n^{-1} G_{s_h^t, a_h^t} \right\|\}$$

442 the sum of the right hand side is similar to the first result. Although the sum is not contiguous, the previous bound holds  
 443 since if  $h_1 < h_2$ ,  $\det(\bar{G}_{n+h_1}) \leq \det(\bar{G}_{n+h_2})$ , this concludes the proof.  $\blacksquare$

444

**Remark 25** We can also write from the lemma in terms of  $\|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^i)^{-1}}$  by skipping the norm upper bound at the beginning of the proof:

$$\sum_{t=1}^T \min\left\{c, \sum_{h=1}^H \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^i)^{-1}}\right\} \leq \frac{c}{\log(1+c)} d \log(1 + \alpha \eta^{-1} B_{\varphi, \mathbb{A}} n).$$

and

$$\begin{aligned}
 \sum_{t=1}^T \sum_{h=1}^H \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^i)^{-1}} &\leq 2d \log(1 + \alpha \eta^{-1} B_{\varphi, \mathbb{A}} n) \\
 &+ \frac{3dH}{\log(2)} \log \left( 1 + \frac{\alpha \|A\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)} \right)
 \end{aligned}$$

445 D.4.2 ELLIPTICAL POTENTIALS: FINITE NUMBER OF LARGE FEATURE NORMS (CONTRIBUTION)

**Lemma 26** (Worst case elliptical potentials, adaptation of Exercise 19.3 Lattimore and Szepesvári (2020) for matrices) Let  $V_0 = \lambda I$  and  $a_1, \dots, a_n \in \mathbb{R}^{d \times p}$  be a sequence of matrices with  $\|a_t\|_2 \leq L$  for all  $t \in [n]$ . Let  $V_t = V_0 + \sum_{s=1}^t a_s a_s^\top$ , then

$$\left| \{t \in \mathbb{N}^*, \|a_t\|_{V_{t-1}^{-1}} \geq 1\} \right| \leq \frac{3d}{\log(2)} \log \left( 1 + \frac{L^2}{\lambda \log(2)} \right)$$

**Proof** Let  $\mathcal{T}$  be the set of rounds  $t$  when  $\|a_t\|_{V_{t-1}^{-1}} \geq 1$  and  $G_t = V_0 + \sum_{s=1}^t \mathbb{I}_{\mathcal{T}}(s) a_s a_s^\top$ . Then

$$\begin{aligned}
 \left(\frac{d\lambda + |\mathcal{T}|L^2}{d}\right)^d &\geq \left(\frac{\text{trace}(G_n)}{d}\right)^d \\
 &\geq \det(G_n) && \text{(Trace-determinant inequality)} \\
 &= \det(V_0) \prod_{t \in \mathcal{T}} \left(1 + \|a_t\|_{G_{t-1}^{-1}}^2\right) \\
 &\geq \det(V_0) \prod_{t \in \mathcal{T}} \left(1 + \|a_t\|_{V_{t-1}^{-1}}^2\right) \\
 &\geq \lambda^{d|\mathcal{T}|}
 \end{aligned}$$

where the third line follows from the matrix determinant lemma:

$$\det(\bar{G}_n + \mathbf{a}_h \mathbf{a}_h^\top) = \det(I + \mathbf{a}_h^\top \bar{G}_n^{-1} \mathbf{a}_h) \det(\bar{G}_n).$$

Rearranging and taking the logarithm shows that

$$|\mathcal{T}| \leq \frac{d}{\log(2)} \log\left(1 + \frac{|\mathcal{T}|L^2}{d\lambda}\right)$$

Abbreviate  $x = d/\log(2)$  and  $y = L^2/d\lambda$ , which are both positive. Then

$$x \log(1 + y(3x \log(1 + xy))) \leq x \log(1 + 3x^2 y^2) \leq x \log(1 + xy)^3 = 3x \log(1 + xy).$$

Since  $z - x \log(1 + yz)$  is decreasing for  $z \geq 3x \log(1 + xy)$  it follows that

$$|\mathcal{T}| \leq 3x \log(1 + xy) = \frac{3d}{\log(2)} \log\left(1 + \frac{L^2}{\lambda \log(2)}\right).$$

446  
447

■

## 448 Appendix E. Tractable planning with random Fourier transform

**A Primer on random Fourier transforms.** We start by defining the Random Fourier Transform and its most relevant property. Let us consider the transition model of Equation (1), we have

$$\mathbb{P}(s' \mid s, a, \theta) = \exp(\psi(s') M_\theta \varphi(s, a) - Z_\theta(s, a)) = \mathbb{E}_{p(w, b)} [f(\psi(s'), w, b) f(M_\theta \varphi(s, a), w, b)],$$

449 where  $f(x, w, b) = \sqrt{2} \cos(w^\top x + b)$  are the random Fourier bases.  $p(w, b) = \mathcal{N}(0, \sigma^{-2} I) \times \mathcal{U}([0, 2\pi])$ , such that  
450  $\mathcal{N}$  is the Gaussian distribution,  $\mathcal{U}$  is the Uniform distribution, and  $p(w, b)$  is a coupling among them.

451 Notice that this provides an alternative approach to decompose the transition kernel and obtain linearity of the value  
452 function. Moreover, since  $\forall x, w \in \mathbb{R}^d, b \in \mathbb{R}, |f(x, w, b)| \leq \sqrt{2}$ , we can use Hoeffding's inequality to prove that a  
453 Monte-Carlo approximation of  $\mathbb{P}(s' \mid s, a, \theta)$  using  $N$  sample pairs of  $(w, b)$  guarantees an error smaller than  $\epsilon$  with  
454 probability at least  $1 - 2 \exp(-N\epsilon^2/4)$ . [Rahimi and Recht \(2007\)](#) proves a stronger result: it provides an algorithm  
455 approximating the Gaussian kernel for which the following uniform convergence bound holds.

**Lemma 27** *Let  $\mathcal{M}$  be a compact subset of  $\mathcal{R}^p$  with diameter  $\text{diam}(\mathcal{M})$ . Then, using the explicit mapping  $\mathbf{z}$  defined in Algorithm 1 in [Rahimi and Recht \(2007\)](#) with  $N$  samples, we have*

$$\Pr \left[ \sup_{x, y \in \mathcal{M}} |\mathbf{z}(x)' \mathbf{z}(y) - k(y, x)| \geq \epsilon \right] \leq 2^8 \left( \frac{\sigma_p \text{diam}(\mathcal{M})}{\epsilon} \right)^2 \exp\left(-\frac{N\epsilon^2}{4(p+2)}\right)$$

456 where  $\sigma_p^2 \equiv E_p[\omega' \omega]$  is the second moment of the Fourier transform of  $k$ .



457 Further, it implies that if  $N = \Omega\left(\frac{p}{\epsilon^2} \log \frac{\sigma_p \text{diam}(\mathcal{M})}{\epsilon}\right)$ , then  $\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{M}} |\mathbf{z}(\mathbf{x})' \mathbf{z}(\mathbf{y}) - k(\mathbf{y}, \mathbf{x})| \leq \epsilon$  with constant  
 458 probability.

459 **Application to planning in BEF-RLSVI.** Since our regret analysis is done under the high probability event of bounded  
 460 estimation parameters, we know that the spaces of  $\psi(s')$  and  $M_\theta \varphi(s, a)$  are bounded and the diameter depends on the  
 461 dimensions. We abstain from explicating the exact diameter as it only influences the number of samples logarithmically.  
 462 Using  $N \approx p/\epsilon^2$  samples, we can construct a uniform  $\epsilon$ -approximation of  $\mathbb{P}(s' | s, a, \theta)$ .

Let's call  $\hat{V}_h^\pi$  the estimated value function using Algorithm 3 with the above approximation of transition. Here, we  
 elucidate the span of this estimation of value function. First we have:

$$\hat{V}_H^\pi - V_H^\pi = \int_{s'} (\hat{P} - P)(s' | s, a) r(s', \pi(s')) ds' \leq \epsilon dH^{3/2}$$

463 Here, we use the facts that  $\mathbb{S}(V_{\hat{\theta}, \hat{\theta}^x, h}) \leq dH^{3/2}$  (cf. Section B.2) and the error in approximating  $P$  is bounded by  $\epsilon$ ,  
 464 i.e.  $\sup_{s', s, a} |(\hat{P} - P)(s' | s, a)| \leq \epsilon$ .

Assume that at step  $h + 1$ , we have  $\hat{V}_{h+1}^\pi - V_{h+1}^\pi \leq \sum_{j=1}^{h+1} \epsilon^j \alpha_{h+1, j}$ . Then, we obtain

$$\begin{aligned} \hat{V}_h^\pi - V_h^\pi &\leq \int_{s'} (\hat{P} - P)(s' | s, a) \hat{V}_{h+1}^\pi(s') ds' + \int_{s'} P(s' | s, a) (\hat{V}_{h+1}^\pi - V_{h+1}^\pi)(s') ds' \\ &= \int_{s'} (\hat{P} - P)(s' | s, a) (V_{h+1}^\pi + \hat{V}_{h+1}^\pi - V_{h+1}^\pi) ds' + \int_{s'} P(s' | s, a) (\hat{V}_{h+1}^\pi - V_{h+1}^\pi)(s') ds' \\ &\leq \epsilon(dH^{3/2} + \sum_{j=1}^{h+1} \epsilon^j \alpha_{h+1, j}) + \sum_{j=1}^{h+1} \epsilon^j \alpha_{h+1, j} \\ &\leq \epsilon(dH^{3/2} + \alpha_{h+1, 1}) + \sum_{j=2}^{h+1} \epsilon^j (\alpha_{h+1, j-1} + \alpha_{h+1, j}) + \epsilon^{h+2} \alpha_{h+1, h+1} \end{aligned}$$

Using the fact that  $\alpha_{1, 1} = dH^{3/2}$  and with a proper induction, we find that:

$$\hat{V}_1^\pi - V_1^\pi \leq \epsilon dH^{5/2} \frac{1 - \epsilon^{H-h}}{1 - \epsilon} \underset{H \rightarrow \infty}{\leq} \epsilon dH^{5/2}$$

465 This concludes the proof of the arguments provided in § Planning of Section 4. This means that the extra regret due  
 466 to planning with the approximation by RFT features is of order  $\mathcal{O}(\epsilon dH^{5/2} K)$ . By choosing an  $\epsilon$  of order  $1/(H\sqrt{K})$ ,  
 467 we deduce that approximating the probability kernel with  $\mathcal{O}(pH^2 K)$  samples induces a tractable planning procedure  
 468 without harming the regret.

469 **Remark 28** *The reader might be tempted to combine the finite approximation using RFT with algorithms from the*  
 470 *linear reinforcement learning literature Jin et al. (2020). However, note that the dimensionality of the linear space*  
 471 *induced by RFT is polynomial in  $H$  and  $K$ . Consequently, applying algorithms designed with the assumption of linear*  
 472 *value function would incur a linear regret.*