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

**Anonymous Author(s)**
Affiliation
Address
`email`

## 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 $\tilde{\mathcal{O}}(\sqrt{d^3 H^3 K})$, 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 $K$.

## 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 [SB18].

**Episodic RL.** In this paper, we consider the episodic finite-horizon MDP formulation of RL, in short *Episodic RL* [ORVR13, AOM17, DLB17]. 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. In episodic RL, the agent interacts with the environment in episodes consisting of $H$ steps. Episode $k$ starts by observing state $s_1^k$. Then, for $t = 1, \ldots H$, 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 $s_{t+1}^k \sim \mathbb{P}(. \mid s_t^k, a_t^k)$ according to 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

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

**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 [AOM17], or maintains a posterior on the parameters to perform posterior sampling [ORVR13], or perturbs the value (or state-action value) function estimates with calibrated noise [OVRW16]. Though the first two approaches induce theoretically optimal exploration, they might not yield tractable algorithms for large/continuous state-action spaces as they either involve optimization in the optimistic set or maintaining a high-dimensional posterior. Thus, *we focus on extending the third approach of* Randomized Least-Square Value Iteration (RLSVI) *framework, and inject noise only in rewards to perform tractable exploration.*

*The second challenge*, which emerges *for continuous state-action spaces, is to learn a parametric functional approximation of either the value function or the rewards and transitions* in order to perform planning and exploration. Different functional representations (or models), such as linear [JYWJ20], bilinear [DKL$^+$21], and bilinear exponential families [CGM21], are studied in literature to develop optimal algorithms for episodic RL with continuous state-action spaces. Since the linear assumption is restrictive in real-life -where non-linear structures are abundant-, generalized representations have obtained more attention recently [CGM21, LLS$^+$21, DKL$^+$21, FKQR21]. The bilinear exponential family model is of special interest as it is expressive enough to represent tabular MDPs (discrete state-action), factored MDPs [KK99], linear MDPs [JYWJ20], linearly controlled dynamical systems (such as Linear Quadratic Regulators [AYS11]) as special cases [CGM21]. Thus, in this paper, *we study the bilinear exponential family of MDPs, i.e. the episodic RL setting where the rewards and transition functions can be modelled with bilinear exponential families.*

*The third challenge is to perform tractable planning[1] given the perturbation for exploration and the model class.* Existing work [OVR14, CGM21] assumes an oracle to perform planning and yield policies that aren't explicit. The main difficulty in such planning approaches is that dynamic programming requires calculating $\int \mathbb{P}(s' \mid s, a)V_h(s)$ for all $(s, a)$ pairs. This is not trivial unless the transition is assumed to be linear and decouples $s'$ from $(s, a)$, which is not known to hold except for tabular MDPs. Much ink has been spilled about this challenge recently, *e.g.* [DKWY19] asks when misspecified linear representations are enough for a polynomial sample complexity in several settings. [SS20, LSW20, VRD19] provide positive answers for specific linear settings. In this paper, *we aim to address this issue by designing a tractable planner for the bilinear exponential family representation.*

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

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

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

1. *Formalism:* We assume that neither rewards nor transitions are known, whereas existing efforts on the bilinear exponential family of MDPs assume knowledge of rewards. This makes the addressed problem harder, practical, and more general. We also observe that though the transition model can represent non-linear dynamics, it implies a linear behavior (see Section 2) in a Reproducible Kernel Hilbert Space (RKHS). This observation contributes to the tractability of planning.

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

---

[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 [RZSD21] | $\sqrt{d^2 H^3 K}$ (Bayesian) | ✗ | ✓ | N.A | Gaussian $\mathbb{P}$ Known rewards |
| LSVI-PHE [ICN$^+$21] | $\sqrt{d^3 H^4 K}$ (Freq.) | ✓ | ✓ | ✗ | Generalized $V$ approx Tabular, anti-concentration |
| OPT-RLSVI [ZBB$^+$20] | $\sqrt{d^4 H^5 K}$ (Freq.) | ✓ | ✓ | ✗ | Linear $V$ |
| EXP-UCRL [CGM21] | $\sqrt{d^2 H^4 K}$ (Freq.) | ✗ | ✗ | N.A | Bilinear Exp family known rewards |
| BEF-RLSVI This work | $\sqrt{d^3 H^3 K}$ (Freq.) | ✓ | ✓ | ✓ | Bilinear Exp family |

3. *Analysis:* We carefully develop an analysis of BEF-RLSVI that yields $\tilde{\mathcal{O}}(\sqrt{d^3 H^3 K})$ regret which improves the existing regret bound for bilinear exponential family of MDPs with known reward by a factor of $\sqrt{H}$ (Section 3.2). Our analysis (Section 5) builds on existing analyses of RLSVI-type algorithms [OVRW16], but contrary to them, we remove the need to handcraft a clipping of the value functions [ZBB$^+$20]. We also do not need to *assume* anti-concentration bounds as we can explicitly control it by the injected noise. This was not done previously except for the linear MDPs. We illustrate this comparison in Table 1. We highlight three technical tools that we used to improve the previous analyses: 1) Using transportation inequalities instead of the simulation lemma reduces a $\sqrt{H}$ factor compared to [RZSD21], 2) Leveraging the observation that true value functions are bounded enables using an improved elliptical lemma (compared to [CGM21]), and 3) Noticing that the norm of features can only be large for a finite amount of time allows us to forgo clipping and reduce a $\sqrt{d}$ factor from the regret compared to [ZBB$^+$20].

## 2 Bilinear exponential family of MDPs

In this section, we introduce the bilinear exponential family model coined in [CGM21] and extend it to 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,

$$\mathbb{P}\left(\tilde{s} \mid s, a\right) = \exp\left(\psi(\tilde{s})^\top M_{\theta^\mathrm{P}} \varphi(s, a) - Z_{s,a}^\mathrm{P}(\theta^\mathrm{P})\right), \tag{1}$$

$$\mathbb{P}\left(r \mid s, a\right) = \exp\left(r\, B^\top M_{\theta^\mathrm{r}} \varphi(s, a) - Z_{s,a}^\mathrm{r}(\theta^\mathrm{r})\right), \tag{2}$$

where $\varphi \in (\mathbb{R}_+^q)^{\mathcal{S} \times \mathcal{A}}$ and $\psi \in (\mathbb{R}_+^p)^{\mathcal{S}}$ are known feature functions, and $B \in \mathbb{R}^p$ is a known scaling factor. The unknown reward and transition parameters are $\theta^\mathrm{P}, \theta^\mathrm{r} \in \mathbb{R}^d$. $M_{\theta^\cdot} \overset{\text{def}}{=} \sum_{i=1}^d \theta_i A_i$, where $A_i$ is a known $p \times q$ matrix for each $i$. Finally, $Z$ denotes the log partition function:

$$Z_{s,a}^\mathrm{P}(\theta^\mathrm{P}) \overset{\text{def}}{=} \log \int_{\mathcal{S}} \exp\left(\psi(\tilde{s})^\top M_{\theta^\mathrm{P}} \varphi(s, a)\right) d\tilde{s},$$

$Z^\mathrm{r}$ is defined similarly. We denote $V_{\theta^\mathrm{P}, \theta^\mathrm{r}, h}^\pi$, respectively $Q_{\theta^\mathrm{P}, \theta^\mathrm{r}, h}^\pi$, the value, respectively state-action value function for policy $\pi$ in the MDP parameterized by $(\theta^\mathrm{P}, \theta^\mathrm{r})$ at time $h$. A policy $\pi^\star$ is *optimal* if for all $s \in \mathcal{S}$, $V_{\theta,h}^{\pi^\star}(s) = \max_{\pi \in \Pi} V_{\theta,h}^\pi(s)$. A learning algorithm minimizes the (pseudo-)regret defined as:

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

**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,

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

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

$$\mathbb{P}\left(s' \mid s, a\right) = \langle \phi^{\text{p}}(s, a), \mu^{\text{p}}(s') \rangle_{\mathcal{H}}, \tag{4}$$

where $\mathcal{H}$ is the RKHS associated with the kernel, $\mu^{\text{p}}(s') = (2\pi)^{-p/2} k\left(\psi(s'), .\right) \exp\left(\|\psi\left(s'\right)\|^2/2\right)$, and $\phi^{\text{p}}(s, a) = k\left(M_{\theta^{\text{p}}}^{\top}\varphi(s, a), .\right) \exp\left(\|M_{\theta^{\text{p}}}\varphi(s, a)\|^2/2 - Z_{s,a}(\theta^{\text{p}})\right)$. Equation (4) shows that $s'$ is decoupled from $(s, a)$, we see hereafter why this is crucial to reducing the complexity of planning.

**Remark 1.** *Up to our knowledge, [RZSD21] is the only work providing an example of linear transition kernel for RL with continuous state-action spaces. They consider Gaussian transitions with an unknown mean ($f^{\star}(s, a)$) and known variance ($\sigma^2$). Actually, linear $f^{\star}$ is a special case of the bilinear exponential family model, where $\psi(s') = (s', \|s'\|^2)$ and $M_{\theta}\varphi(s, a) = (f_{\theta}(s, a)/\sigma^2, -1/\sigma^2)$.*

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

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

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

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

When $\phi^{\text{p}}, \mu^{\text{p}} \in \mathbb{R}^{\tau}$, we can obtain $Q_{\theta^{\text{p}}, \theta^{\text{r}}, h}$ by computing $\tau$ integrals per timestep, reducing the state-action space complexity to $\tau$ only. For our model, although $\phi^{\text{p}}$ and $\mu^{\text{p}}$ are infinite dimensional, we show in Section 4 (§ planning) that the planning complexity is still significantly reduced.

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

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

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

`BEF-RLSVI` is based on `RLSVI` [OVRW16] framework with the distinction that we only perturb the reward parameters and not all the parameters of the value function. `RLSVI` algorithms are reminiscent of Thompson Sampling, yet more tractable with better control over the probability to be optimistic.

---

**Algorithm 1** `BEF-RLSVI`

---

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

---

We can see that Algorithm 1 performs exploration by a Gaussian perturbation of the reward parameter (Line 4). Contrary to optimistic approaches, this method is explicit and also more efficient since it does not a involve high-dimensional optimization.

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

We can approximate Line 3 of Algorithm 2 with $\mathcal{O}(pH^3 K \log(HK))$ complexity and without harming the learning process (*cf.* § planning, Section 4). Therefore, here, planning is tractable.

## 3.2 `BEF-RLSVI`: regret upper-bound

We state the standard smoothness assumptions on the model [CGM21, JBNW17, LMT21].

**Assumption 1.** *There exist constants $\alpha^{\mathrm{p}}, \alpha^{\mathrm{r}}, \beta^{\mathrm{p}}, \beta^{\mathrm{r}} > 0$, such that the representation model satisfies:*

$$\forall (s,a) \in \mathcal{S} \times \mathcal{A}, \forall \theta, x \in \mathbb{R}^d \quad \alpha^{\mathrm{p}} \leq x^\top C_{s,a}^\theta[\psi] x \leq \beta^{\mathrm{p}}$$

$$\forall (s,a) \in \mathcal{S} \times \mathcal{A}, \forall \theta, x \in \mathbb{R}^d \quad \alpha^{\mathrm{r}} \leq \mathbb{V}\mathrm{ar}_{s,a}^\theta(r)\, x^\top B^\top B x \leq \beta^{\mathrm{r}}$$

*where* $\mathbb{C}_{s,a}^\theta[\psi(s')] \triangleq \mathbb{E}_{s' \sim \mathbb{P}_\theta | s,a}\left[\psi(s')\psi(s')^\top\right] - \mathbb{E}_{s' \sim \mathbb{P}_\theta | s,a}[\psi(s')] \mathbb{E}_{s' \sim \mathbb{P}_\theta | s,a}\left[\psi(s')^\top\right]$ *and*

$\mathbb{V}\mathrm{ar}_{s,a}^\theta(r) \triangleq \left(\mathbb{E}_{s,a}^\theta[r^2] - \mathbb{E}_{s,a}^\theta[r]^2\right)$ *is the variance of the reward under $\theta$.*

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

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

**Theorem 2** (Regret bound). *Let $\mathbb{A} \triangleq (\mathrm{tr}(A_i A_j^\top))_{i,j \in [d]}$ and $G_{s,a} \triangleq (\varphi(s,a)^\top A_i^\top A_j \varphi(s,a))_{i,j \in [d]}$. Under Assumption 1 and further considering that*

*1. $\max\{\|\theta^{\mathrm{r}}\|_{\mathbb{A}}, \|\theta^{\mathrm{p}}\|_{\mathbb{A}}\} \leq B_{\mathbb{A}}$, $\|\mathbb{A}^{-1} G_{s,a}\| \leq B_{\varphi,\mathbb{A}}$ and $\mathbb{E}_{\theta^{\mathrm{r}}}[r(s,a)] \in [0,1]$ for all $(s,a)$.*

*2. noise $\xi_k \sim \mathcal{N}(0, x_k(\bar{G}_k^{\mathrm{p}})^{-1})$ satisfies $x_k \geq \left(H\sqrt{\frac{\beta^{\mathrm{p}} \beta^{\mathrm{p}}(K,\delta)}{\alpha^{\mathrm{p}} \alpha^{\mathrm{r}}}} + \frac{\sqrt{\beta^{\mathrm{r}} \beta^{\mathrm{r}}(K,\delta) \min\{1, \frac{\alpha^{\mathrm{p}}}{\alpha^{\mathrm{r}}}\}}}{2\alpha^{\mathrm{r}}}\right)^2 \propto dH^2$,*

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

$$\mathcal{R}(K) \leq \sqrt{KH}\Bigg[ \underbrace{2H\left(\sqrt{\frac{2\beta^{\mathrm{p}}}{\alpha^{\mathrm{p}}}}\beta^{\mathrm{p}}(K,\delta)\gamma_K^{\mathrm{p}} + (1+\sqrt{\gamma_K^{\mathrm{r}}})\sqrt{\log(1/\delta^2)}\right)}_{\text{Transition concentration} \approx dH} + \underbrace{\beta^{\mathrm{r}}\sqrt{\frac{\beta^{\mathrm{r}}(n,\delta)\gamma_K^{\mathrm{r}}}{2\alpha^{\mathrm{r}}}}}_{\text{Reward concentration} \approx d}$$

$$+ \underbrace{c\beta^{\mathrm{r}}\sqrt{x_K d\gamma_K^{\mathrm{r}}\log(dK/\delta)} + \frac{\beta^{\mathrm{r}}\sqrt{x_K d\gamma_K^{\mathrm{r}}\log(e/\delta^2)}}{\Phi(-1)}(1+\sqrt{\log(d/\delta)})}_{\text{Noise concentration} \approx d^{3/2}H}\Bigg]$$

$$+ \sqrt{H\gamma_K^{\mathrm{r}}}\Bigg[ \underbrace{\beta^{\mathrm{r}}C_d\left(\sqrt{\frac{\beta^{\mathrm{r}}(K,\delta)}{2\alpha^{\mathrm{r}}}} + c\sqrt{x_K d\log(dK/\delta)}\right)}_{\text{Estimation error for no clipping} \approx dH}$$

$$+ \underbrace{\frac{\beta^{\mathrm{r}} d\sqrt{x_K}}{\Phi(-1)}(1+\sqrt{\log(d/\delta)})\sqrt{C_d\left(1+\frac{\alpha^{\mathrm{r}}B_{\varphi,A}H}{\eta}\right)}}_{\text{Learning error for no clipping} \approx (dH)^{3/2}}\Bigg],$$

*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 $\gamma_K^i \triangleq d\log(1 + \frac{\beta^i}{\eta}B_{\varphi,\mathbb{A}}HK)$. Also,*

$C_d \triangleq \frac{3d}{\log(2)}\log\left(1 + \frac{\alpha^{\mathrm{r}}\|\mathbb{A}\|_2^2 B_{\varphi,\mathbb{A}}^2}{\eta \log(2)}\right)$, *$\Phi$ is the Gaussian CDF, and $c$ is a universal constant.*

Theorem 2 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 bilinear exponential family model, $K$ is the number of episodes, and $H$ is the horizon of an episode. We now clarify how this contrasts with related literature.

145 *Comparison with other bounds.* The closest work to ours is [CGM21] as it considers the same
146 model for transitions but with known rewards. They propose a UCRL-type and PSRL-type algorithm,
147 which achieve a regret of order $\widetilde{O}(\sqrt{d^2 H^4 K})$. There are two notable algorithmic differences with
148 our work. First, they do exploration using intractable-optimistic upper bounds or high-dimensional
149 posteriors, while we do it with explicit perturbation. The second difference is in planning. While
150 they assume access to a planning oracle, we do it explicitly with pseudo-polynomial complexity
151 (Section 4). Moreover, we improve the regret bound by a $\sqrt{H}$ factor thanks to an improved analysis,
152 (*cf.* Lemma 18). But similar to all RLSVI-type algorithms, we pick up an extra $\sqrt{d}$ (*cf.* [AL17]).

153 [ZBB+20] proposes a variant of RLSVI for continuous state-action spaces, where there are low-rank
154 models of transitions and rewards. They show a regret bound $R(K) = \widetilde{O}(\sqrt{d^4 H^5 K})$, which is larger
155 than that of BEF-RLSVI by $O(\sqrt{dH^2})$. In algorithm design, we improve on their work by removing
156 the need to carefully clip the value function. Analytically, our model allows us to use transportation
157 inequalities (*cf.* Lemma 13) instead of the simulation lemma, which saves us a $\sqrt{H}$ factor.

158 [RZSD21] considers Gaussian transitions, i.e. $s' = f^*(s,a) + \epsilon$ such that $\epsilon \sim \mathcal{N}\left(0, \sigma^2\right)$. This is a
159 particular case of our model. They propose to use Thompson Sampling, and have the merit of being
160 the first to have observed linearity of the value function from this transition structure. But they do not
161 connect it to the finite dimensional approximation of [RR07] unlike us (Section 4). Finally, they show
162 a Bayesian regret bound of $O(\sqrt{d^2 H^3 K})$. This notion of regret is weaker than frequentist regret,
163 hence this result is not directly comparable with Theorem 2.

164 *Tightness of regret bound.* A lower bound for episodic RL with continuous state-action spaces is
165 still missing. However, for tabular RL, [DMKV21] proves a lower bound of order $\Omega(\sqrt{H^3 SAK})$.
166 If we represent a tabular MDP in our model, we would need $d = S^2 \times A$ parameters (Section 4.3,
167 [CGM21]). In this case, our bound becomes $R(K) = O(\sqrt{(S^2 A)^3 H^3 K})$, which is clearly not tight
168 is $S$ and $A$. This is understandable due to the relative generality of our setting. We are however
169 positively surprised that **our bound is tight in terms of its dependence on $H$ and $K$.**

## 170 4 Algorithm design: building blocks of BEF-RLSVI

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

172 **Estimation of parameters.** We estimate transitions and rewards from observations similar to
173 EXP-UCRL [CGM21], *i.e.* by using a penalized maximum likelihood estimator

$$\hat{\theta}^{\mathtt{p}}(k) \in \underset{\theta \in \mathbb{R}^d}{\arg\min} \sum_{t=1}^{k} \sum_{h=1}^{H} -\log \mathbb{P}_\theta\left(s_{h+1}^t \mid s_h^t, a_h^t\right) + \eta \operatorname{pen}(\theta).$$

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

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

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

$$\hat{\theta}^{\mathtt{r}}(k) \in \underset{\theta \in \mathbb{R}^d}{\arg\min} \sum_{t=1}^{k} \sum_{h=1}^{H} -\log \mathbb{P}_\theta\left(r_t \mid s_h^t, a_h^t\right) + \eta \operatorname{pen}(\theta), \tag{7}$$

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

178 **Exploration.** A significant challenge in RL is handling exploration in continuous spaces. The majority
179 of the literature is split between intractable, upper confidence bound-style optimism or Thompson
180 sampling algorithms with high-dimensional posterior and guarantees only in terms of Bayesian
181 regret. In BEF-RLSVI, we adopt the approach of reward perturbation motivated by the RLSVI-
182 framework [ZBB+20, OVRW16]. We show that perturbing the reward estimation can guarantee

6

optimism with a constant probability, *i.e.* there exists $\nu \in (0,1]$ such that for all $k \in [K]$ and $s_1^k \in \mathcal{S}$,

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

[ZBB+20] proves that this suffices to bound the learning error. However, their method clashes with not clipping the value function, as it modifies the probability of optimism. Thus, [ZBB+20] proposes an involved clipping procedure to handle the issue of unstable values. Instead, by careful geometric analysis (*cf.* Lemma 19), we bound the occurrences of the unstable values, and in turn, upper bound the regret without clipping. Note that unlike [ICN+21], `BEF-RLSVI` does not guarantee that the estimated value function is optimistic but still is able to control the learning error (*cf.* Section 5).

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

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

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

For details about the finite-dimensional approximation of our transition kernel, refer to Appendix E. Now, we highlight the schematic of a finite-dimensional approximation of $\phi^{\mathrm{p}}$ and $\psi^{\mathrm{p}}$. We proceed in three steps. **1)** We have with high probability $\mathbb{S}(V_{\hat{\theta}^{\mathrm{p}},\tilde{\theta}^{\mathrm{r}},h}) \leq dH^{3/2}$ (Section 5). **2)** If we have a uniform $\epsilon$-approximation of $\mathbb{P}_{\theta^{\mathrm{p}}}$, we show that using it incurs at most an extra $\mathcal{O}(\epsilon dH^{5/2}K)$ regret. **3)** Finally, following [RR07], we approximate uniformly the shift invariant kernels, here the RBF in Equation (4), within $\epsilon$ error and with features of dimensions $\mathcal{O}(p\epsilon^{-2}\log\frac{1}{\epsilon^2})$, where $p$ is dimension of $\psi$. Associating these three elements and choosing $\epsilon = 1/\sqrt{(H^2K)}$, we establish our claim.

# 5 Theoretical analysis: proof outline

To convey the novelties in our analysis, we provide a proof sketch for Theorem 2. We start by decomposing the regret into an estimation loss and a learning error, as given below

$$R(K) = \sum_{k=1}^K (V^\star_{\theta^{\mathrm{p}},\theta^{\mathrm{r}},1} - V^{\pi_k}_{\theta^{\mathrm{p}},\theta^{\mathrm{r}},1})(s_{1k}) = \sum_{k=1}^K (\underbrace{V^\star_{\theta^{\mathrm{p}},\theta^{\mathrm{r}},1} - V^{\pi_k}_{\hat{\theta}^{\mathrm{p}},\tilde{\theta}^{\mathrm{r}},1}}_{learning} + \underbrace{V^{\pi_k}_{\hat{\theta}^{\mathrm{p}},\tilde{\theta}^{\mathrm{r}},1} - V^{\pi_k}_{\theta^{\mathrm{p}},\theta^{\mathrm{r}},1}}_{Estimation})(s_{1k}). \quad (9)$$

For the **estimation error**, we use smoothness arguments with concentrations of parameters up to some novelties. Regarding the **learning error**, we show that the injected noise ensures a constant probability of anti-concentration. Applying Assumption 1 and Lemma 18 leads to the upper-bound.

## 5.1 Bounding the estimation error

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

$$V^\pi_{\hat{\theta}^{\mathrm{p}},\tilde{\theta}^{\mathrm{r}}}(s_{1k}) - V^\pi_{\theta^{\mathrm{p}},\theta^{\mathrm{r}}}(s_{1k}) = \underbrace{V^\pi_{\hat{\theta}^{\mathrm{p}},\theta^{\mathrm{r}}}(s_{1k}) - V^\pi_{\theta^{\mathrm{p}},\theta^{\mathrm{r}}}(s_{1k})}_{\text{transition estimation}} + \underbrace{V^\pi_{\hat{\theta}^{\mathrm{p}},\tilde{\theta}^{\mathrm{r}}}(s_{1k}) - V^\pi_{\hat{\theta}^{\mathrm{p}},\theta^{\mathrm{r}}}(s_{1k})}_{\text{reward estimation}} \quad (10)$$

**Transition estimation** Since the reward parameter is exact, the value function's span is $\leq H$. Then, using the transportation of Lemma 13 we obtain the bound $H\sum_{h=1}^H \sqrt{2\,\mathrm{KL}_{s_{hk},a_{hk}}(\theta^{\mathrm{p}},\hat{\theta}^{\mathrm{p}})}$. We notice that since the reward parameter is exact, the bound is actually $H\min\{1,\sum_{h=1}^H \sqrt{2\,\mathrm{KL}_{s_{hk},a_{hk}}(\theta^{\mathrm{p}},\hat{\theta}^{\mathrm{p}})}\}$. Using Lemma 18 under Assumption 1, we win a $\sqrt{H}$ factor compared to the analysis of [CG19].

**Reward estimation** Previous work uses clipping to help control this error, but in this case it can hinder the optimism probability by biasing the noise. [ZBB+20] proposes an involved clipping depending on the norms $\|(A_i\varphi(s_h^k,a_h^k))_{i\in[d]}\|_{(\bar{G}_k^{\mathrm{p}})^{-1}}$, which is somewhat delicate to analyze and

deploy. We remedy the situation acting solely in the proof. First let's define what we call the set of "bad rounds": $\left\{ k \in [K], \exists h : \|(A_i \varphi(s_h^k, a_h^k))_{i \in [d]}\|_{(\bar{G}_k^{\mathrm{p}})^{-1}} \geq 1 \right\}$, these rounds are why clipping is necessary. Thanks to Lemma 19, we know that the number of such rounds is at most $\mathcal{O}(d)$. Surprisingly, it depends neither on $H$ nor on $K$. We show that the "bad rounds" incur at most $O(d^{3/2}H^2)$ regret, independent of $K$. Therefore, our algorithm can forgo clipping for free.

**Remark 2.** *If it wasn't for the episodic nature of our setting, we could have used the forward algorithm to eliminate the span control issue. We refer to [Vov01, AW01] for a description of this algorithm, [OMP21] for a stochastic analysis, and Section 4 therein for an application to linear bandits.*

## 5.2 Bounding the learning error

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

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

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

The first and second terms are perturbation free, we handle them similarly to the estimation error, *i.e.* using concentration arguments for $\hat{\theta}^{\mathrm{p}}$ and $\hat{\theta}^{\mathrm{r}}$. For the third term, we use transportation of rewards (Lemma 17) and anti-concentration of $\xi_k$ (Lemma 12). We find that with probability at least $1 - 2\delta$

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

where $c(n, \delta) = \left( \sqrt{\beta^{\mathrm{p}} \beta^{\mathrm{p}}(n, \delta)/\alpha^{\mathrm{p}}} + \sqrt{\beta^{\mathrm{r}} \beta^{\mathrm{r}}(n, \delta) \min\{1, \alpha^{\mathrm{p}}/\alpha^{\mathrm{r}}\}/(2\alpha^{\mathrm{r}})} \right)$. Since $\xi_k \sim \mathcal{N}(0, x_k (\bar{G}_k^{\mathrm{p}})^{-1})$ and $x_k \geq H^2 c(n, \delta)^2$, we get $\mathbb{P}\left( V_{\hat{\theta}^{\mathrm{p}}, \tilde{\theta}^{\mathrm{r}}, 1}^\pi (s_1) - V_{\theta^{\mathrm{p}}, \theta^{\mathrm{r}}, 1}^\star (s_1) \geq 0 \right) \geq \Phi(-1)$, where $\Phi$ is the normal CDF. This is ensured by the anti-concentration property of Gaussian random variables, see Lemma 12.

**From stochastic optimism to error control:** Existing algorithms require the value function to be optimistic (*i.e.* negative learning error) with large probability. Contrary to them, `BEF-RLSVI` only requires the estimated value to be optimistic with a constant probability. When it is, the learning happens. Otherwise, the policy is still close to a good one thanks to the decreasing estimation error, and the learning still happens. This part of the proof is similar in spirit to that of [ZBB+20].

*Upper bound on $V_1^\star$:* 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 $\bar{O}_k \triangleq \{V_{\hat{\theta}^{\mathrm{p}}, \tilde{\theta}_k^{\mathrm{r}}, 1}(s_1^k) - V_1^\star(s_1^k) \geq 0\}$. This implies that $\quad V_1^\star(s_1^k) \leq \mathbb{E}_{\bar{\xi}_k | \bar{O}_k} [V_{\hat{\theta}^{\mathrm{p}}, \hat{\theta}^{\mathrm{r}} + \bar{\xi}_k, 1}(s_1^k)]$.

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

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

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

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

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

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

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

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

*General functional representation.* [DSL+18] provides the first convergence guarantee for general nonlinear function representations in the Maximum Entropy RL setting, where entropy of a policy is used as a regularizer to induce exploration. Thus, the analysis cannot address episodic RL, where we have to explicitly ensure exploration with optimism. [WSY20] proposes a framework that leverages the optimism with confidence bound approach for general functional representations with bounded Eluder dimensions, which is a complexity measure in RL. However, knowing the Eluder dimension is crucial for the optimistic confidence bound in their algorithm. Eluder dimension is not known for MDPs except linear and tabular MDPs. *To concretize our design, we focus on the general but explicit bilinear exponential family of MDPs than any abstract representation.*

*Bilinear exponential family of MDPs.* Exponential families are studied widely in RL theory, from bandits to MDPs [LMT21, KKM13, FCGS10, KH06], as an expressive parametric family to design theoretically-grounded model-based algorithms. [CGM21] first studies episodic RL with Bilinear Exponential Family (BEF) of transitions, which is linear in both state-action pairs and the next-state. It proposes a regularized log-likelihood method to estimate the model parameters, and two optimistic algorithms with upper confidence bounds and posterior sampling. Due to its generality to unifiedly model tabular MDPs, factored MDPs, linear MDPs, and linearly controlled dynamical systems, the BEF-family of MDPs has received increasing attention [LLS+21]. [LLS+21] estimates the model parameters based on score matching that enables them to replace regularity assumption on the log-partition function with Fisher-information and assumption on the parameters. Both [CGM21, LLS+21] achieve a worst-case regret of order $\tilde{O}(\sqrt{d^2 H^4 K})$ for known reward. On a different note, [DKL+21, FKQR21] also introduces a new structural framework for generalization in RL, called bilinear classes as it requires the Bellman error to be upper bounded by a bilinear form. Instead of using bilinear forms to capture non-linear structures, this class is not identical to BEF class of MDPs, and studying the connection is out of the scope of this paper. Specifically, *we address the shortcomings of the existing works on BEF-family of MDPs that assume known rewards, absence of RLSVI-type algorithms, and access to oracle planners.*

*Tractable planning and linearity.* Planning is a major byproduct of the chosen functional represen-tation. In general, planning can incur high computational complexity if done naïvely. Specially, [DKWY19] shows that for some settings, even with a linear $\epsilon$-approximation of the $Q$-function, a planning procedure able to produce an $\epsilon$-optimal policy has a complexity at least $2^H$. Thus, different works [SS20, LSW20, VRD19] propose to leverage different low-dimensional representations of value functions or transitions to perform efficient planning. Here, we take note from [RZSD21] that Gaussian transitions induce an explicit linear value function in an RKHS. And generalize this observation with the bilinear exponential. Moreover, using uniformly good features [RR07] to approximate transition dynamics from our model enables us to design a tractable planner. We provide a detailed discussion of this approximation in Section 4. More practically, [RZSD21, NY21] use representations given by random Fourier features [RR07] to approximate the transition dynamics and provide experiments validating the benefits of this approach for high-dimensional Atari-games.

## 7 Conclusion and future work

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

# References

[AL17] Marc Abeille and Alessandro Lazaric. Linear thompson sampling revisited. In *Artificial Intelligence and Statistics*, pages 176–184. PMLR, 2017. (Cited on 6, 28)

[AOM17] Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*, pages 263–272. PMLR, 2017. (Cited on 1, 2)

[AW01] Katy S Azoury and Manfred K Warmuth. Relative loss bounds for on-line density estimation with the exponential family of distributions. *Machine Learning*, 43(3):211–246, 2001. (Cited on 8)

[AYS11] Yasin Abbasi-Yadkori and Csaba Szepesvári. Regret bounds for the adaptive control of linear quadratic systems. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 1–26. JMLR Workshop and Conference Proceedings, 2011. (Cited on 2)

[BLM13] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013. (Cited on 17, 29)

[CG19] Sayak Ray Chowdhury and Aditya Gopalan. Online learning in kernelized markov decision processes. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3197–3205. PMLR, 2019. (Cited on 7, 25)

[CGM21] Sayak Ray Chowdhury, Aditya Gopalan, and Odalric-Ambrym Maillard. Reinforcement learning in parametric mdps with exponential families. In *International Conference on Artificial Intelligence and Statistics*, pages 1855–1863. PMLR, 2021. (Cited on 2, 3, 5, 6, 9, 18, 24)

[DKL+21] Simon Du, Sham Kakade, Jason Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong Wang. Bilinear classes: A structural framework for provable generalization in rl. In *International Conference on Machine Learning*, pages 2826–2836. PMLR, 2021. (Cited on 2, 9)

[DKWY19] Simon S Du, Sham M Kakade, Ruosong Wang, and Lin F Yang. Is a good representation sufficient for sample efficient reinforcement learning? *arXiv preprint arXiv:1910.03016*, 2019. (Cited on 2, 9)

[DLB17] Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying pac and regret: Uniform pac bounds for episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 30, 2017. (Cited on 1)

[DMKV21] Omar Darwiche Domingues, Pierre Ménard, Emilie Kaufmann, and Michal Valko. Episodic reinforcement learning in finite mdps: Minimax lower bounds revisited. In *Algorithmic Learning Theory*, pages 578–598. PMLR, 2021. (Cited on 6)

[DSL+18] Bo Dai, Albert Shaw, Lihong Li, Lin Xiao, Niao He, Zhen Liu, Jianshu Chen, and Le Song. Sbeed: Convergent reinforcement learning with nonlinear function approximation. In *International Conference on Machine Learning*, pages 1125–1134. PMLR, 2018. (Cited on 9)

[FCGS10] Sarah Filippi, Olivier Cappe, Aurélien Garivier, and Csaba Szepesvári. Parametric bandits: The generalized linear case. *Advances in Neural Information Processing Systems*, 23, 2010. (Cited on 9)

[FKQR21] Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021. (Cited on 2, 9)

[ICN+21] Haque Ishfaq, Qiwen Cui, Viet Nguyen, Alex Ayoub, Zhuoran Yang, Zhaoran Wang, Doina Precup, and Lin F Yang. Randomized exploration for reinforcement learning with general value function approximation. *arXiv preprint arXiv:2106.07841*, 2021. (Cited on 3, 7)

[JBNW17] Kwang-Sung Jun, Aniruddha Bhargava, Robert Nowak, and Rebecca Willett. Scalable generalized linear bandits: Online computation and hashing. *arXiv preprint arXiv:1706.00136*, 2017. (Cited on 5)

[JYWJ20] Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR, 2020. (Cited on 2, 34)

[KH06] Branislav Kveton and Milos Hauskrecht. Solving factored mdps with exponential-family transition models. In *ICAPS*, pages 114–120, 2006. (Cited on 9)

[KK99] Michael Kearns and Daphne Koller. Efficient reinforcement learning in factored mdps. In *IJCAI*, volume 16, pages 740–747, 1999. (Cited on 2)

[KKM13] Nathaniel Korda, Emilie Kaufmann, and Remi Munos. Thompson sampling for 1-dimensional exponential family bandits. *Advances in neural information processing systems*, 26, 2013. (Cited on 9)

[LLS+21] Gene Li, Junbo Li, Nathan Srebro, Zhaoran Wang, and Zhuoran Yang. Exponential family model-based reinforcement learning via score matching. *arXiv preprint arXiv:2112.14195*, 2021. (Cited on 2, 9)

[LMT21] Yangyi Lu, Amirhossein Meisami, and Ambuj Tewari. Low-rank generalized linear bandit problems. In *International Conference on Artificial Intelligence and Statistics*, pages 460–468. PMLR, 2021. (Cited on 5, 9)

[LS20] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020. (Cited on 33)

[LSW20] Tor Lattimore, Csaba Szepesvari, and Gellert Weisz. Learning with good feature representations in bandits and in rl with a generative model. In *International Conference on Machine Learning*, pages 5662–5670. PMLR, 2020. (Cited on 2, 9)

[NY21] Ofir Nachum and Mengjiao Yang. Provable representation learning for imitation with contrastive fourier features. *Advances in Neural Information Processing Systems*, 34, 2021. (Cited on 9)

[OMP21] Reda Ouhamma, Odalric-Ambrym Maillard, and Vianney Perchet. Stochastic online linear regression: the forward algorithm to replace ridge. *Advances in Neural Information Processing Systems*, 34:24430–24441, 2021. (Cited on 8)

[ORVR13] Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via posterior sampling. *Advances in Neural Information Processing Systems*, 26, 2013. (Cited on 1, 2)

[OVR14] Ian Osband and Benjamin Van Roy. Model-based reinforcement learning and the eluder dimension. *Advances in Neural Information Processing Systems*, 27, 2014. (Cited on 2)

[OVRW16] Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value functions. In *International Conference on Machine Learning*, pages 2377–2386. PMLR, 2016. (Cited on 2, 3, 4, 6)

[RR07] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Advances in neural information processing systems*, 20, 2007. (Cited on 6, 7, 9, 34)

[RZSD21] Tongzheng Ren, Tianjun Zhang, Csaba Szepesvári, and Bo Dai. A free lunch from the noise: Provable and practical exploration for representation learning. *arXiv preprint arXiv:2111.11485*, 2021. (Cited on 3, 4, 6, 9, 17)

[SB18] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018. (Cited on 1)

[SS20] Roshan Shariff and Csaba Szepesvári. Efficient planning in large mdps with weak linear function approximation. *arXiv preprint arXiv:2007.06184*, 2020. (Cited on 2, 9)

[Vov01] Volodya Vovk. Competitive on-line statistics. *International Statistical Review*, 69(2):213–248, 2001. (Cited on 8)

[VRD19] Benjamin Van Roy and Shi Dong. Comments on the du-kakade-wang-yang lower bounds. *arXiv preprint arXiv:1911.07910*, 2019. (Cited on 2, 9)

[WSY20] Ruosong Wang, Ruslan Salakhutdinov, and Lin F Yang. Reinforcement learning with general value function approximation: Provably efficient approach via bounded eluder dimension. *arXiv preprint arXiv:2005.10804*, 2020. (Cited on 9)

[ZBB$^+$20] Andrea Zanette, David Brandfonbrener, Emma Brunskill, Matteo Pirotta, and Alessandro Lazaric. Frequentist regret bounds for randomized least-squares value iteration. In *International Conference on Artificial Intelligence and Statistics*, pages 1954–1964. PMLR, 2020. (Cited on 3, 6, 7, 8)

# Checklist

1. For all authors...

    (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]

    (b) Did you describe the limitations of your work? [Yes]

    (c) Did you discuss any potential negative societal impacts of your work? [N/A] This is a purely theoretical contribution

    (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]

2. If you are including theoretical results...

    (a) Did you state the full set of assumptions of all theoretical results? [Yes]

    (b) Did you include complete proofs of all theoretical results? [Yes] See the appendices

3. If you ran experiments...

    (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [N/A]

    (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]

    (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]

    (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

    (a) If your work uses existing assets, did you cite the creators? [Yes] We cite creator of the bilinear exponential family model.

    (b) Did you mention the license of the assets? [N/A]

    (c) Did you include any new assets either in the supplemental material or as a URL? [No]

    (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]

    (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]

5. If you used crowdsourcing or conducted research with human subjects...

    (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]

    (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]

    (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

# Appendix

## Table of Contents

## A  Notations

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

Table 2: Notations

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

## 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:

$$\mathcal{R}(K) = \sum_{k=1}^{K}(V_{\theta^p,\theta^r,1}^{\star} - V_{\theta^p,\theta^r,1}^{\pi_k})(s_1^k) = \sum_{k=1}^{K}\Big(\underbrace{V_{\theta^p,\theta^r,1}^{\star} - V_{\hat{\theta}^p,\tilde{\theta}^r,1}^{\pi_k}}_{learning} + \underbrace{V_{\hat{\theta}^p,\tilde{\theta}^r,1}^{\pi_k} - V_{\theta^p,\theta^r,1}^{\pi_k}}_{Estimation}\Big)(s_1^k)$$

### B.1    Estimation error

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

$$V_{\hat{\theta}^p,\tilde{\theta}^r}^{\pi}(s_1^k) - V_{\theta^p,\theta^r}^{\pi}(s_1^k) = \underbrace{V_{\hat{\theta}^p,\theta^r}^{\pi}(s_1^k) - V_{\theta^p,\theta^r}^{\pi}(s_1^k)}_{\text{transition estimation}} + \underbrace{V_{\hat{\theta}^p,\tilde{\theta}^r}^{\pi}(s_1^k) - V_{\hat{\theta}^p,\theta^r}^{\pi}(s_1^k)}_{\text{reward estimation}},$$

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

**Lemma 3.** *The estimation error satisfies, with probability at least $1 - 5\delta$*

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

$$+ \left[\sqrt{KHd\log\left(1 + \alpha^r\eta^{-1}B_{\varphi,\mathbb{A}}n\right)} + C_d\sqrt{Hd\log(1 + \alpha\eta^{-1}B_{\varphi,\mathbb{A}}H)}\right] \times \left(\sqrt{\frac{\beta^r(n,\delta)}{2\alpha^r}}\right.$$

$$\left. + c\sqrt{(\max_k x_k)d\log(dK/\delta)}\right)\beta^r + \sqrt{2KHd\log\left(1 + \alpha^r\eta^{-1}B_{\varphi,\mathbb{A}}n\right)\log(1/\delta)}$$

*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 $\gamma_K^i \triangleq d\log(1 + \frac{\beta^i}{\eta}B_{\varphi,\mathbb{A}}HK)$. Also,*

*$C_d \triangleq \frac{3d}{\log(2)}\log\left(1 + \frac{\alpha^r\|\mathbb{A}\|_2^2 B_{\varphi,\mathbb{A}}^2}{\eta\log(2)}\right)$, and $c$ is a universal constant.*

*Proof.* It follows directly by combining Lemma 4 and Lemma 5 using a union bound. □

#### B.1.1    Transition estimation

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

**Lemma 4.** *We have, with probability at least $1 - 2\delta$*

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

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

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

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

$$V_{\hat{\theta}^{\mathrm{p}},\theta^r,1}^{\pi}(s_1^k) - V_{\theta^{\mathrm{p}},\theta^r,1}^{\pi}(s_1^k)$$

$$= \mathbb{E}_{s_1^k,a_1^k}^{\hat{\theta}^{\mathrm{p}}}\left[V_{\hat{\theta}^{\mathrm{p}},\theta^r,2}^{\pi}\right] - \mathbb{E}_{s_1^k,a_1^k}^{\theta^{\mathrm{p}}}\left[V_{\hat{\theta}^{\mathrm{p}},\theta^r,2}^{\pi}\right] + \mathbb{E}_{s_1^k,a_1^k}^{\theta^{\mathrm{p}}}[V_{\hat{\theta}^{\mathrm{p}},\theta^r,2}^{\pi} - V_{\theta^{\mathrm{p}},\theta^r,2}^{\pi}]$$

$$= \mathbb{E}_{s_1^k,a_1^k}^{\hat{\theta}^{\mathrm{p}}}\left[V_{\hat{\theta}^{\mathrm{p}},\theta^r,2}^{\pi}\right] - \mathbb{E}_{s_1^k,a_1^k}^{\theta^{\mathrm{p}}}\left[V_{\hat{\theta}^{\mathrm{p}},\theta^r,2}^{\pi}\right] + V_{\hat{\theta}^{\mathrm{p}},\theta^r,2}^{\pi}(s_{2k}) - V_{\theta^{\mathrm{p}},\theta^r,2}^{\pi}(s_{2k}) + \zeta_1^k$$

$$= \sum_{h=1}^{H} \mathbb{E}_{s_{hk},a_{hk}}^{\hat{\theta}^{\mathrm{p}}}\left[V_{\hat{\theta}^{\mathrm{p}},\theta^r,h+1}^{\pi}\right] - \mathbb{E}_{s_{hk},a_{hk}}^{\theta^{\mathrm{p}}}\left[V_{\hat{\theta}^{\mathrm{p}},\theta^r,h+1}^{\pi}\right] + \zeta_{hk}$$

where $\zeta_{hk} = \mathbb{E}_{s_{hk},a_{hk}}^{\theta^{\mathrm{p}}}[V_{\hat{\theta}^{\mathrm{p}},\theta^r,h+1}^{\pi} - V_{\theta^{\mathrm{p}},\theta^r,h+1}^{\pi}] - \left(V_{\hat{\theta}^{\mathrm{p}},\theta^r,h+1}^{\pi}(s_{h+1k}) - V_{\theta^{\mathrm{p}},\theta^r,h+1}^{\pi}(s_{h+1k})\right)$ is a martingale sequence, and the last equality comes by induction. Here we consider the true reward parameter which verifies $|\mathbb{E}_{\theta^r}[r(s,a)]| \leq 1$ by assumption, therefore $|\zeta_{hk}| \leq 2H$. Using the Azuma-Hoeffding inequality [BLM13], 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 13, indeed

$$\mathbb{E}_{s_{hk},a_{hk}}^{\hat{\theta}^{\mathrm{p}}}\left[V_{\hat{\theta}^{\mathrm{p}},\theta^r,h+1}^{\pi}\right] - \mathbb{E}_{s_{hk},a_{hk}}^{\theta^{\mathrm{p}}}\left[V_{\hat{\theta}^{\mathrm{p}},\theta^r,h+1}^{\pi}\right] \leq H\sqrt{2\,\mathrm{KL}_{s_{hk},a_{hk}}(\theta^{\mathrm{p}},\hat{\theta}^{\mathrm{p}})}$$

$$\leq H\min\left\{1, \sqrt{2\,\mathrm{KL}_{s_{hk},a_{hk}}(\theta^{\mathrm{p}},\hat{\theta}^{\mathrm{p}})}\right\},$$

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

**Remark 3.** *Traditionally, the expected value difference bound follows from the simulation lemma [RZSD21]. The simulation lemma incurs an extra $\sqrt{H}$ factor compared to our bound.*

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

$$\sum_{k=1}^{K} V_{\hat{\theta}^{\mathrm{p}},\theta^r}(s_1^k) - V_{\theta^{\mathrm{p}},\theta^r}^{\pi}(s_1^k)$$

$$\leq H\sum_{k=1}^{K}\min\left\{1, \sum_{h=1}^{H}\sqrt{2\,\mathrm{KL}_{s_{hk},a_{hk}}(\theta^{\mathrm{p}},\hat{\theta}^{\mathrm{p}})}\right\} + 2H\sqrt{2KH\log(1/\delta)} \quad (11)$$

*2) Bounding the sum of KL divergences.* we explicit the bound of inequality (11) using Assumption 1 along with properties of the exponential family (*cf.* Section D.3). We have for all $(s,a)$,

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

This implies that

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

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

From Corollary 8, with probability at least $1 - \delta$ and for all $k \in \mathbb{N}$

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

Also, using Lemma 18, we have

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

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

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

**Remark 4.** *Notice that the minimum with 1 is crucial, indeed, without it the bound deteriorates by a factor $H$ as was the case in [CGM21].*

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

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

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

### B.1.2 Reward estimation

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

**Lemma 5.** *With probability at least $1 - 3\delta$, the following result holds true.*

$$\sum_{k=1}^{K}V_{\hat{\theta}^{\mathrm{p}}, \tilde{\theta}^{\mathrm{r}}, 1}^{\pi}(s_1^k) - V_{\hat{\theta}^{\mathrm{p}}, \theta^{\mathrm{r}}, 1}^{\pi}(s_1^k) \leq \left(\sqrt{\frac{\beta^{\mathrm{r}}(K, \delta)}{2\alpha^{\mathrm{r}}}} + c\sqrt{(\max_{k \leq K}x_k)d\log(dK/\delta)}\right)\beta^{\mathrm{r}}$$

$$\times \left(\sqrt{C_d\left(1 + \frac{\alpha^{\mathrm{r}}B_{\varphi, A}H}{\eta}\right)} + \sqrt{K\log(e/\delta^2)}\right)\sqrt{Hd\log\left(1 + \alpha^{\mathrm{r}}\eta^{-1}B_{\varphi, \mathbb{A}}HK\right)},$$

*where $\beta^{\mathrm{p}}(K, \delta) \triangleq \frac{\eta}{2}B_{\mathbb{A}}^2 + \gamma_K^{\mathrm{p}} + \log(1/\delta)$, and $\gamma_K^{\mathrm{p}} \triangleq d\log(1 + \frac{\beta^{\mathrm{p}}}{\eta}B_{\varphi, \mathbb{A}}HK)$. Also, $C_d \triangleq \frac{3d}{\log(2)}\log\left(1 + \frac{\alpha^{\mathrm{r}}\|\mathbb{A}\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta\log(2)}\right)$, and $c$ is a universal constant.*

*Proof.* The reward estimation error in Equation (10) can be written explicitly. Indeed, using Lemma 17

$$V_{\hat{\theta}^{\mathrm{p}}, \tilde{\theta}^{\mathrm{r}}, 1}^{\pi}(s_1^k) - V_{\hat{\theta}^{\mathrm{p}}, \theta^{\mathrm{r}}, 1}^{\pi}(s_1^k) = \mathbb{E}_{(\tilde{s}_h)_{1 \leq h \leq H}\sim\pi|\hat{\theta}^{\mathrm{p}}, s_1^k}\left[\sum_{h=1}^{H}\frac{\mathbb{V}\mathrm{ar}_{\tilde{s}_h, \pi(\tilde{s}_h)}(r)}{2}B^{\top}M_{\tilde{\theta}^{\mathrm{r}} - \theta^{\mathrm{r}}}\varphi(\tilde{s}_h, \pi(\tilde{s}_h))\right]$$

$$\leq \mathbb{E}\left[\sum_{h=1}^{H}\frac{\mathbb{V}\mathrm{ar}_{\tilde{s}_h, \pi(\tilde{s}_h)}(r)}{2}\|\tilde{\theta}^{\mathrm{r}} - \theta^{\mathrm{r}}\|_{\bar{G}_k^{\mathrm{r}}}\|(B^{\top}A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^{\mathrm{r}})^{-1}}\right]$$

$$\leq \|\tilde{\theta}^{\mathrm{r}} - \theta^{\mathrm{r}}\|_{\bar{G}_k^{\mathrm{r}}}\mathbb{E}\left[\sum_{h=1}^{H}\frac{\mathbb{V}\mathrm{ar}_{\tilde{s}_h, \pi(\tilde{s}_h)}(r)}{2}\|(B^{\top}A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^{\mathrm{r}})^{-1}}\right]$$

$$\leq \|\tilde{\theta}^{\mathrm{r}} - \theta^{\mathrm{r}}\|_{\bar{G}_k^{\mathrm{r}}}\frac{\beta^{\mathrm{r}}}{2}\mathbb{E}\Big[\underbrace{\sum_{h=1}^{H}\|(A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^{\mathrm{r}})^{-1}}}_{\stackrel{\mathrm{def}}{=}\widetilde{\mathrm{traj}}_k}\Big],$$

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

**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:

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

540 *1)* We know that $\|(A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1\le i\le d}(A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1\le i\le d}^\top\|_2^2 \le \|\mathbb{A}\|_2^2 B_{\varphi,\mathbb{A}}^2$. Consequently,
541 according to Lemma 19

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

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

544 This further yields

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

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

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

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

$$\|(A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1\le i\le d}\|_{(\bar{G}_k^{\mathtt{r}})^{-1}} = \sqrt{\mathrm{tr}((A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1\le i\le d}^\top (\bar{G}_k^{\mathtt{r}})^{-1}(A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1\le i\le d})}$$

$$= \sqrt{\mathrm{tr}\left(\left(1 + \frac{\alpha^{\mathtt{r}}B_{\varphi,A}H}{\eta}\right)(\bar{G}_{k+H}^{\mathtt{r}})^{-1}G_{s,a}\right)}$$

$$\le \sqrt{\left(1 + \frac{\alpha^{\mathtt{r}}B_{\varphi,A}H}{\eta}\right)}\|(A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1\le i\le d}\|_{(\bar{G}_{k+H}^{\mathtt{r}})^{-1}}$$

547 Since $\|(A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1\le i\le d}\|_{(\bar{G}_{k+H}^{\mathtt{r}})^{-1}} \le 1$, we have $\|(A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1\le i\le d}\|_{(\bar{G}_{k+H}^{\mathtt{r}})^{-1}} \le$
548 $\min\left\{1, \|(A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1\le i\le d}\|_{(\bar{G}_k^{\mathtt{r}})^{-1}}\right\}$. Consequently

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

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

$$\sum_{k\in\mathcal{T}}\sum_{h\in[H]} V_{\hat{\theta}^{\mathtt{p}},\tilde{\theta}^{\mathtt{r}},1}^\pi(s_1^k) - V_{\hat{\theta}^{\mathtt{p}},\theta^{\mathtt{r}},1}^\pi(s_1^k) \le \|\tilde{\theta}^{\mathtt{r}} - \theta^{\mathtt{r}}\|_{\bar{G}_k^{\mathtt{r}}}\frac{\beta^{\mathtt{r}}}{2}$$

$$\sqrt{\frac{3d}{\log(2)}\log\left(1 + \frac{\alpha^{\mathtt{r}}\|\mathbb{A}\|_2^2 B_{\varphi,\mathbb{A}}^2}{\eta\log(2)}\right)\left(1 + \frac{\alpha^{\mathtt{r}}B_{\varphi,A}H}{\eta}\right)Hd\log(1 + \alpha^{\mathtt{r}}\eta^{-1}B_{\varphi,\mathbb{A}}H)} \quad (14)$$

550 **Remark 5.** *The bad rounds analysis is one of our most important contributions as it enables us to*
551 *forgo clipping without consequences. Consequently, this is a novel method to control the reward*
552 *estimation error that improves on existing work for whom clipping was essential.*

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

$$\zeta_k' \overset{\text{def}}{=} \mathrm{traj}_k - \mathbb{E}_{(\tilde{s}_h)_{1\le h\le H}\sim\pi|\hat{\theta}^{\mathtt{p}},s_1^k}\left[\widetilde{\mathrm{traj}}_k\right].$$

554 where $\widetilde{\mathrm{traj}}_k$ is the same quantity as $\mathrm{traj}$ but with a random realization of state transitions.
555 Since all feature norms are smaller than one, $(\zeta_k')_k$ is a martingale sequence with $|\zeta_k'| \le$
556 $\sqrt{Hd\log(1 + \alpha^{\mathtt{r}}\eta^{-1}B_{\varphi,\mathbb{A}}HK)}$. We deduce that with probability at least $1 - \delta$:

$$\sum_{k=1}^K \zeta_k' \le \sqrt{2KHd\log(1 + \alpha^{\mathtt{r}}\eta^{-1}B_{\varphi,\mathbb{A}}HK)\log(1/\delta)}$$

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

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

558 The last inequality follows from controlling the concentration of the reward parameter. First we ob-
559 serve that (Corollary 10) with probability at least $1 - \delta$, uniformly over $k \in \mathbb{N}$, $\left\| \theta^r - \hat{\theta}^r(k) \right\|^2_{\bar{G}^r_k} \leq$
560 $\frac{2}{\alpha^r} \beta^r(k, \delta)$. Second, we also have that for all $k \geq 1$, with probability at least $1 - \delta$, $\|\xi_k\|_{G^r_k} \leq$
561 $c\sqrt{x_k d \log(d/\delta)}$, we then use a union bound. Combining with Equation (14) we find

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

562 This concludes the proof. $\qquad\square$

563 **Remark 6.** *If we use Lemma 17 without the martingale difference sequence, it will lead to a linear*
564 *regret. Indeed, the span of the sum of norms over an episode is of order $\sqrt{H}$. Using the martingale*
565 *technique instead allows us to retrieve a telescopic sum controlled using the elliptical lemma, this is*
566 *essential to obtaining a sub-linear regret bound.*

## B.2  Learning error

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

570 **Lemma 6.** *If the variance parameter of the injected noise $(\xi_k)_k$ satisfies*

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

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

$$\sum_{k=1}^K V^\star_1(s_1^k) - V^\pi_{\hat{\theta}^p, \hat{\theta}^r + \bar{\xi}_k, 1}(s_1^k) \leq \frac{d\beta^r \sqrt{x_k} \left( 1 + \sqrt{\log(d/\delta)} \right)}{\Phi(-1)} \sqrt{H \log\left(1 + \alpha^r \eta^{-1} B_{\varphi, \mathbb{A}} HK\right)}$$
$$\times \left( \sqrt{C_d \left( 1 + \frac{\alpha^r B_{\varphi, A} H}{\eta} \right)} + \sqrt{K \log(e/\delta^2)} \right),$$

572 *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 $\gamma_K^i \triangleq d \log(1 + \frac{\beta^i}{\eta} B_{\varphi, \mathbb{A}} HK)$. Also*
573 $C_d \overset{\text{def}}{=} \frac{3d}{\log(2)} \log\left(1 + \frac{\alpha^r \|\mathbb{A}\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)}\right)$, *and $\Phi$ is the normal CDF.*

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

### B.2.1  Stochastic optimism

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

$$(V_{\hat{\theta}^p, \tilde{\theta}^r, 1} - V^\star_{\hat{\theta}^p, \theta^r, 1})(s_1) \geq (Q^\star_{\hat{\theta}^p, \tilde{\theta}^r, 1} - Q^\star_1)(s_1, \pi^\star(s_1))$$
$$\geq \underbrace{V^{\pi^\star}_{\hat{\theta}^p, \theta^r}(s_1) - V^{\pi^\star}_{\hat{\theta}^p, \theta^r}(s_1)}_{\text{first term}} + \underbrace{V^{\pi^\star}_{\hat{\theta}^p, \hat{\theta}^r}(s_1) - V^{\pi^\star}_{\hat{\theta}^p, \theta^r}(s_1)}_{\text{second term}} + \underbrace{V^{\pi^\star}_{\hat{\theta}^p, \tilde{\theta}^r}(s_1) - V^{\pi^\star}_{\hat{\theta}^p, \hat{\theta}^r}(s_1)}_{\text{third term}}$$

**First term.** By assumption, the expected reward under the true parameter satisfies $\mathbb{E}_{\theta^{\mathtt{r}}}[r(s,a)] \in [0,1]$, then $\mathbb{S}\left(\sum_{t=1}^{H} \mathbb{E}_{\theta^{\mathtt{r}}}[r(s_t, \pi(s_t))]\right) \le H$. Consequently, the first term can be controlled using Lemma 13

$$V_{\theta^{\mathtt{p}}, \theta^{\mathtt{r}}}^{\pi^\star}(s_1) - V_{\hat{\theta}^{\mathtt{p}}, \theta^{\mathtt{r}}}^{\pi^\star}(s_1) \le H\sqrt{\mathrm{KL}(P_{\hat{\theta}^{\mathtt{p}}}(s_2, \ldots, s_H), P_{\theta^{\mathtt{p}}}(s_2, \ldots, s_H))}$$

$$\le H\sqrt{\mathbb{E}_{(\tilde{s}_t)_{t\in[H]}\sim\hat{\theta}^{\mathtt{p}}|s_1^k}\left[\sum_{t=1}^{H} \psi(\tilde{s}_{t+1})^\top M_{\hat{\theta}^{\mathtt{p}}-\theta^{\mathtt{p}}}\varphi(\tilde{s}_t, \pi^\star(\tilde{s}_t)) + Z_{\theta^{\mathtt{p}}}^{\mathtt{p}}(\tilde{s}_t, \pi^\star(\tilde{s}_t)) - Z_{\hat{\theta}^{\mathtt{p}}}^{\mathtt{p}}(\tilde{s}_t, \pi^\star(\tilde{s}_t))\right]}$$

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

$$\mathbb{E}_{(\tilde{s}_t)_{t\in[H]}\sim\hat{\theta}^{\mathtt{p}}|s_1^k}\left[\psi(\tilde{s}_{t+1})^\top M_{\hat{\theta}^{\mathtt{p}}-\theta^{\mathtt{p}}}\varphi(\tilde{s}_t, \pi^\star(\tilde{s}_t)) + Z_{\theta^{\mathtt{p}}}^{\mathtt{p}}(\tilde{s}_t, \pi^\star(\tilde{s}_t)) - Z_{\hat{\theta}^{\mathtt{p}}}^{\mathtt{p}}(\tilde{s}_t, \pi^\star(\tilde{s}_t))\right]$$

$$= \frac{1}{2}(\hat{\theta}^{\mathtt{p}} - \theta^{\mathtt{p}})^\top \mathbb{E}_{(\tilde{s}_t)_{t\in[H]}\sim\hat{\theta}^{\mathtt{p}}|s_1^k}\left[\nabla_{s_h, \pi^\star(s_h)}^2 Z^{\mathtt{p}}(\theta_h)\right](\hat{\theta}^{\mathtt{p}} - \theta^{\mathtt{p}})$$

$$\le \frac{\beta^{\mathtt{p}}}{2}\mathbb{E}_{(\tilde{s}_t)_{t\in[H]}\sim\hat{\theta}^{\mathtt{p}}|s_1^k}\left[\|\hat{\theta}^{\mathtt{p}} - \theta^{\mathtt{p}}\|_{G_{\tilde{s}_h, \pi^\star(\tilde{s}_h)}}^2\right].$$

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

$$V_{\theta^{\mathtt{p}}, \theta^{\mathtt{r}}}^{\pi^\star}(s_1) - V_{\hat{\theta}^{\mathtt{p}}, \theta^{\mathtt{r}}}^{\pi^\star}(s_1) \le H\sqrt{\frac{\beta^{\mathtt{p}}}{2}\sum_{h=1}^{H}\mathbb{E}_{(\tilde{s}_t)_{t\in[H]}\sim\hat{\theta}^{\mathtt{p}}|s_1^k}\left[\|\hat{\theta}^{\mathtt{p}} - \theta^{\mathtt{p}}\|_{G_{\tilde{s}_h, \pi^\star(\tilde{s}_h)}}^2\right]}$$

$$\le H\sqrt{\frac{\beta^{\mathtt{p}}}{2}}\left\|\hat{\theta}^{\mathtt{p}} - \theta^{\mathtt{p}}\right\|_{\sum_{h=1}^{H}\mathbb{E}_{(\tilde{s}_t)_{t\in[H]}\sim\hat{\theta}^{\mathtt{p}}|s_1^k}\left[G_{\tilde{s}_h, \pi^\star(\tilde{s}_h)}\right]}$$

$$\le H\sqrt{\frac{\beta^{\mathtt{p}}}{2}}\left\|\hat{\theta}^{\mathtt{p}} - \theta^{\mathtt{p}}\right\|_{u_k u_k^\top}$$

$$\le H\sqrt{\frac{\beta^{\mathtt{p}}}{2}}\left\|(\bar{G}_k^{\mathtt{p}})^{-1/2}u_k u_k^\top(\bar{G}_k^{\mathtt{p}})^{-1/2}\right\|\|\hat{\theta}^{\mathtt{p}} - \theta^{\mathtt{p}}\|_{\bar{G}_k^{\mathtt{p}}}$$

$$\le H\sqrt{\frac{\beta^{\mathtt{p}}}{2}}\|u_k\|_{(\bar{G}_k^{\mathtt{p}})^{-1}}\|\hat{\theta}^{\mathtt{p}} - \theta^{\mathtt{p}}\|_{\bar{G}_k^{\mathtt{p}}}$$

The third line follows because $\forall x \in \mathbb{R}^d, \quad \|x\|_{\sum_{i=1} a_i a_i^\top} \le \|x\|_{(\sum_{i=1} a_i)(\sum_{i=1} a_i)^\top}$, and the last one follows because $\mathrm{tr}(AB) \le \mathrm{tr}(A)\,\mathrm{tr}(B)$ for any two real positive semi-definite matrices $A$ and $B$. We deduce, with probability at least $1 - \delta$:

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

**Second term.** We have

$$V_{\hat{\theta}^{\mathtt{p}}, \hat{\theta}^{\mathtt{r}}}^{\pi^\star}(s_1) - V_{\hat{\theta}^{\mathtt{p}}, \theta^{\mathtt{r}}}^{\pi^\star}(s_1) = \mathbb{E}_{(\tilde{s}_t)_{t\in[H]}\sim\hat{\theta}^{\mathtt{p}}|s_1^k}\left[\sum_{t=1}^{H}\frac{\mathbb{V}\mathrm{ar}^{\theta_t^{\mathtt{r}}}(r)}{2}B^\top M_{\hat{\theta}^{\mathtt{r}}-\theta^{\mathtt{r}}}\varphi(\tilde{s}_t, \pi^\star(\tilde{s}_t))\right]$$

$$= (\hat{\theta}^{\mathtt{r}} - \theta^{\mathtt{r}})^\top \mathbb{E}_{(\tilde{s}_t)_{t\in[H]}\sim\hat{\theta}^{\mathtt{p}}|s_1^k}\left[\sum_{t=1}^{H}\frac{\mathbb{V}\mathrm{ar}^{\theta_t^{\mathtt{r}}}(r)}{2}(A_i\varphi(\tilde{s}_t, \pi^\star(\tilde{s}_t)))_{i\in[d]}\right]B$$

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

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 $\beta^{\mathtt{r}}$. We conclude that with probability at least $1 - \delta$,

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

21

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

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

**Third term.** We have

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

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

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

Choosing $x_k \geq \left( H \sqrt{\frac{\beta^{\mathrm{p}} \beta^{\mathrm{p}}(k, \delta)}{\alpha^{\mathrm{p}} \alpha^{\mathrm{r}}}} + \frac{\sqrt{\beta^{\mathrm{r}} \beta^{\mathrm{r}}(k, \delta) \min\{1, \frac{\alpha^{\mathrm{p}}}{\alpha^{\mathrm{r}}}\}}}{2\alpha^{\mathrm{r}}} \right)$ and using Lemma 12, we find that the perturbed value function is optimistic with probability at least $\Phi(-1)$.

### B.2.2 Controlling the learning error

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

**Upper bound on $V_1^\star$.** 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 $k$:
$$
\bar{O}_k = \{ V_{\hat{\theta}^{\mathrm{p}}, \hat{\theta}^{\mathrm{r}} + \bar{\xi}_k, 1}(s_1^k) - V_1^\star(s_1^k) \geq 0 \} \tag{15}
$$
we know that $\mathbb{P}(\bar{O}_k) \geq \Phi(-1)$. This event provides the upper bound:
$$
V_1^\star(s_1^k) \leq \mathbb{E}_{\bar{\xi}_k | \bar{O}_k}[V_{\hat{\theta}^{\mathrm{p}}, \hat{\theta}^{\mathrm{r}} + \bar{\xi}_k, 1}(s_1^k)] \tag{16}
$$

**Lower bound on $V_{\hat{\theta}^{\mathrm{p}}, \tilde{\theta}^{\mathrm{r}}}$.** We define this bound with an optimization problem under concentration of the noise. Consider $\underline{V}_1(s_1^k)$ is the solution of
$$
\min_{\xi_k} V_{\hat{\theta}^{\mathrm{p}}, \hat{\theta}^{\mathrm{r}} + \xi_k, 1}(s_{1^k}) \tag{17}
$$
$$
\|\xi_k\|_{\bar{G}_k^{\mathrm{p}}} \leq \sqrt{x_k d \log(d/\delta)}, \quad \forall t \in [H]
$$

Under the concentration of our injected noise, we obtain
$$
\underline{V}_1(s_1^k) \leq V_{\hat{\theta}^{\mathrm{p}}, \tilde{\theta}^{\mathrm{r}}}(s_1^k) \tag{18}
$$

**Combining the error bounds.** Combining the upper bound of Equation (16) with the lower bound of Equation (18), we get, with probability at least $1 - \delta$:
$$
V_1^\star(s_1^k) - V_{\hat{\theta}^{\mathrm{p}}, \hat{\theta}^{\mathrm{r}} + \bar{\xi}_k, 1}(s_1^k) \leq \mathbb{E}_{\bar{\xi}_k | \bar{O}_k}[V_{\hat{\theta}^{\mathrm{p}}, \hat{\theta}^{\mathrm{r}} + \bar{\xi}_k, 1}(s_1^k) - \underline{V}_1(s_1^k)]
$$

611 Also, using the tower rule,

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

612 Therefore,

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

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

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

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

$$V_{\hat{\theta}^{\mathrm{p}},\tilde{\theta}^{\mathrm{r}},1}^{\pi}(s_1^k) - V_{\hat{\theta}^{\mathrm{p}},\hat{\theta}^{\mathrm{r}}+\underline{\xi}_k,1}^{\pi}(s_1^k)$$

$$= \mathbb{E}_{(\tilde{s}_h)_{1\leq h\leq H}\sim\pi|\hat{\theta}^{\mathrm{p}},s_1^k}\left[\sum_{h=1}^{H}\frac{\mathbb{V}\mathrm{ar}_{\tilde{s}_h,\pi(\tilde{s}_h)}(r)}{2}B^{\top}M_{\tilde{\theta}^{\mathrm{r}}-\hat{\theta}^{\mathrm{r}}-\underline{\xi}_k}\varphi(\tilde{s}_h,\pi(\tilde{s}_h))\right]$$

$$\leq \mathbb{E}\left[\sum_{h=1}^{H}\frac{\mathbb{V}\mathrm{ar}_{\tilde{s}_h,\pi(\tilde{s}_h)}(r)}{2}\|\tilde{\theta}^{\mathrm{r}}-\hat{\theta}^{\mathrm{r}}-\underline{\xi}_k\|_{\bar{G}_k^{\mathrm{p}}}\|(B^{\top}A_i\varphi(\tilde{s}_h,\pi(\tilde{s}_h)))_{1\leq i\leq d}\|_{(\bar{G}_k^{\mathrm{p}})^{-1}}\right]$$

$$\leq \|\tilde{\theta}^{\mathrm{r}}-\hat{\theta}^{\mathrm{r}}-\underline{\xi}_k\|_{\bar{G}_k^{\mathrm{p}}}\mathbb{E}\left[\sum_{h=1}^{H}\frac{\mathbb{V}\mathrm{ar}_{\tilde{s}_h,\pi(\tilde{s}_h)}(r)}{2}\|(B^{\top}A_i\varphi(\tilde{s}_h,\pi(\tilde{s}_h)))_{1\leq i\leq d}\|_{(\bar{G}_k^{\mathrm{p}})^{-1}}\right]$$

$$\leq \|\tilde{\xi}_k-\underline{\xi}_k\|_{\bar{G}_k^{\mathrm{p}}}\frac{\beta^{\mathrm{r}}}{2}\mathbb{E}\left[\sum_{h=1}^{H}\|(A_i\varphi(\tilde{s}_h,\pi(\tilde{s}_h)))_{1\leq i\leq d}\|_{(\bar{G}_k^{\mathrm{p}})^{-1}}\right]$$

616 Then,

$$\mathbb{E}_{\tilde{\xi}_k}\left[V_{\hat{\theta}^{\mathrm{p}},\tilde{\theta}^{\mathrm{r}},1}^{\pi}(s_1^k) - V_{\hat{\theta}^{\mathrm{p}},\hat{\theta}^{\mathrm{r}}+\underline{\xi}_k,1}^{\pi}(s_1^k)\right]$$

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

617 Also,

$$\left|\mathbb{E}_{\xi_k|\bar{O}_k^{\mathrm{c}}}[V_{\hat{\theta}^{\mathrm{p}},\hat{\theta}^{\mathrm{r}}+\xi_k,1}(s_1^k) - \underline{V}_1(s_1^k)]\right|$$

$$\leq \frac{\beta^{\mathrm{r}}}{2}\mathbb{E}_{\tilde{\xi}_k|\bar{O}_k^{\mathrm{c}}}[\|\tilde{\xi}_k-\underline{\xi}_k\|_{\bar{G}_k^{\mathrm{p}}}]\mathbb{E}_{(\tilde{s}_h)\sim\pi|\hat{\theta}^{\mathrm{p}}}\left[\sum_{h=1}^{H}\|(A_i\varphi(\tilde{s}_h,\pi(\tilde{s}_h)))_{1\leq i\leq d}\|_{(\bar{G}_k^{\mathrm{p}})^{-1}}\right]$$

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

618 We have a bound on the expected value of the sum of feature norms in the proof of Lemma 5. Also,

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

619 The second line follows from Cauchy-Schwarz and by definition of $\underline{\xi}_k$. The last line is due to the
620 fact that $x_k(\bar{G}_k^{\mathrm{p}})^{-1} \sim \mathcal{N}(0, x_k I_d)$, which implies $\|\tilde{\xi}_k\|_{\bar{G}_k^{\mathrm{p}}}^2 \sim \mathcal{N}(0, dx_k)$. We conclude the proof by
621 taking the sum of feature norms from the proof of Lemma 5.

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

$$\sum_{k=1}^{K} V_1^\star(s_1^k) - V_{\hat{\theta}^{\mathrm{p}}, \hat{\theta}^{\mathrm{r}} + \bar{\xi}_k, 1}(s_1^k) \le \frac{\beta^{\mathrm{r}}}{\Phi(-1)}(\sqrt{x_k d} + \sqrt{x_k d \log(d/\delta)})$$

$$\left[ \sqrt{\frac{3d}{\log(2)} \log\left( 1 + \frac{\alpha^{\mathrm{r}} \|\mathbb{A}\|_2^2 B_{\varphi,\mathbb{A}}^2}{\eta \log(2)} \right) \left( 1 + \frac{\alpha^{\mathrm{r}} B_{\varphi,A} H}{\eta} \right) H d \log(1 + \alpha^{\mathrm{r}} \eta^{-1} B_{\varphi,\mathbb{A}} H)} \right.$$

$$\left. + \sqrt{K H d \log\left(1 + \alpha^{\mathrm{r}} \eta^{-1} B_{\varphi,\mathbb{A}} H K\right) \log(e/\delta^2)} \right]$$

# C   Concentrations

## C.1   Concentration of the transition parameter

625    We recall the important concentration of the maximum likelihood estimator for general bilinear
626    exponential families (*cf.* Theorem 1 of [CGM21]).

627    **Theorem 7.** *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) $(s_t, a_t)$*
628    *is $\mathcal{F}_{t-1}$ measurable, and (iii) given $(s_t, a_t)$, $s_{t+1} \sim P_{\theta^{\mathrm{p}}}^{\mathrm{p}}(\cdot \mid s_t, a_t)$ according to the exponential*
629    *family defined by Equation (1). Let $\hat{\theta}^{\mathrm{p}}(k)$ be the penalized MLE defined by Equation (6), and let*
630    *$Z_{s,a}^{\mathrm{p}}(\theta)$ be strictly convex in $\theta$ for all $(s, a)$. Then, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$,*
631    *the following holds uniformly over all $n \in \mathbb{N}$:*

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

632    *where $C_{\mathbb{A},k}^{\mathrm{p}} = \left( \int_{\mathbb{R}^d} \exp\left( -\frac{\eta}{2} \|\theta'\|_{\mathbb{A}}^2 \right) d\theta' \right) / \left( \int_{\mathbb{R}^d} \exp\left( -\sum_{t=1}^{k} \mathrm{KL}_{s_t, a_t}(\theta_k, \theta') - \frac{\eta}{2} \|\theta' - \theta_k\|_{\mathbb{A}}^2 \right) d\theta' \right)$.*

633    *Define $G_{s,a} \stackrel{\text{def}}{=} \left( \varphi(s, a)^\top A_i^\top A_j \varphi(s, a) \right)_{i,j \in [d]}$, we have*

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

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

635    A proof of this result can be found in the work [CGM21]. We provide an almost similar proof for the
636    concentration of rewards in the next section.

637    **Corollary 8.** *The previous theorem implies a simple euclidean confidence region. Indeed, with*
638    *probability at least $1 - \delta$, for all $k \in \mathbb{N}$*

$$\left\| \theta^{\mathrm{p}} - \hat{\theta}^{\mathrm{p}}(k) \right\|_{\bar{G}_n^{\mathrm{p}}}^2 \le \frac{2}{\alpha^{\mathrm{p}}} \beta^{\mathrm{p}}(k, \delta),$$

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

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

$$\frac{1}{2} \left\| \theta^{\mathrm{p}} - \hat{\theta}^{\mathrm{p}}(k) \right\|_{\bar{G}_k}^2 = \frac{(\alpha^{\mathrm{p}})^{-1} \eta}{2} \left\| \theta^{\mathrm{p}} - \hat{\theta}^{\mathrm{p}}(k) \right\|_{\mathbb{A}}^2 + \sum_{\tau=1}^{k-1} \sum_{h=1}^{H} \frac{1}{2} \left\| \theta^{\mathrm{p}} - \hat{\theta}^{\mathrm{p}}(k) \right\|_{G_{s_h^\tau, a_h^\tau}}^2$$

$$\le (\alpha^{\mathrm{p}})^{-1} \left( \frac{\eta}{2} \left\| \theta^{\mathrm{p}} - \hat{\theta}^{\mathrm{p}}(k) \right\|_{\mathbb{A}}^2 + \sum_{\tau=1}^{k-1} \sum_{h=1}^{H} \mathrm{KL}_{s_h^\tau, a_h^\tau}(\theta_k, \theta) \right).$$

641    $\square$

## C.2 Concentration of the reward parameter (contribution)

**Theorem 9.** *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 $(s_t, a_t)$, $r(s_t, a_t) \sim P_{\theta^r}^r(\cdot \mid s_t, a_t)$ according to the exponential family defined by (2). Let $\hat{\theta}^r(k)$ be the penalized MLE defined by Equation (8), and let $Z_{s,a}^r(\theta)$ 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 $k \in \mathbb{N}$:*

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

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

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

*where $\beta^r := \|B\|_2^2 \sup_{\theta, s, a} \mathbb{V}\mathrm{ar}_{s,a}^\theta(r)$.*

*Proof.* We proceed similar to the proof of Theorem 1 in [CG19].

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

$$\mathrm{KL}_{s,a}\left(\theta^r, \theta^{r'}\right) := \mathrm{KL}\left(P_{\theta^r}(\cdot \mid s, a), P_{\theta^{r'}}(\cdot \mid s, a)\right) = B_{Z_{s,a}}\left(\theta^{r'}, \theta^r\right).$$

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

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

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

$$\mathbb{E}_{s,a}^{\theta^r}\left[\exp\left(\sum_{i=1}^{d} \lambda_i\left(r(s_t, a_t) - \mathbb{E}_{s_t, a_t}^{\theta^r}[r]\right) B^\top A_i \varphi(s_t, a_t)\right)\right]$$

$$= \exp\left(-\lambda^\top \nabla Z_{s,a}^r(\theta^r)\right) \int_{\mathcal{S}} \exp\left(\sum_{i=1}^{d}\left(\theta_i^r + \lambda_i\right) B^\top A_i \varphi(s, a) - Z_{s,a}^r(\theta^r)\right) dr$$

$$= \exp\left(Z_{s,a}^r(\theta^r + \lambda) - Z_{s,a}^r(\theta^r) - \lambda^\top \nabla Z_{s,a}^r(\theta^r)\right) = \exp\left(B_{Z_{s,a}^r}(\theta^r)\right)$$

This implies $\mathbb{E}\left[\exp\left(\lambda^\top S_n\right) \mid \mathcal{F}_{n-1}\right] = \exp\left(\lambda^\top S_{n-1} + B_{Z_{n_n, a_n}, \theta^r}(\lambda)\right)$ thus $\mathbb{E}\left[M_n^\lambda \mid \mathcal{F}_{n-1}\right] = M_{n-1}^\lambda$. 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 $\mathbb{E}\left[M_n^\lambda\right] = 1$. For any prior density $q(\theta)$ for $\theta$, we now define a mixture of martingales

$$M_n = \int_{\mathbb{R}^d} M_n^\lambda q\left(\theta^r + \lambda\right) d\lambda \tag{19}$$

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

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

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

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

$$B_{Z^{\mathrm{r}}_n}(\theta',\theta) = \sum_{t=1}^n \left(Z^{\mathrm{r}}_{s_t,a_t}(\theta') - Z^{\mathrm{r}}_{s_t,a_t}(\theta) - (\theta' - \theta)^\top \nabla Z^{\mathrm{r}}_{S_t,a_t}(\theta)\right) = \sum_{t=1}^n B_{Z^{\mathrm{r}}_{s_t,\alpha_t}}(\theta',\theta)$$

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

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

This yields

$$S_k = \sum_{t=1}^k \left(\nabla Z^{\mathrm{r}}_{s_t,a_t}\left(\hat{\theta}^{\mathrm{r}}(k)\right) - \nabla Z^{\mathrm{r}}_{s_t,a_t}(\theta^{\mathrm{r}})\right) + \eta\mathbb{A}\hat{\theta}^{\mathrm{r}}(k) = \nabla Z^{\mathrm{r}}_k\left(\hat{\theta}^{\mathrm{r}}(k)\right) - \nabla Z^{\mathrm{r}}_k(\theta^{\mathrm{r}}) + \eta\mathbb{A}\hat{\theta}^{\mathrm{r}}(k) \tag{21}$$

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

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

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

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

$$g_k(\lambda) = g_n(\lambda) + \sup_{\lambda\in\mathbb{R}^\star} g_k(\lambda) - g_k(\lambda^\star)$$

$$= \frac{\eta}{2}\left\|\hat{\theta}^{\mathrm{r}}(k) - \theta^{\mathrm{r}}\right\|^2_{\mathbb{A}} + \eta(\lambda - \lambda^\star)^\top \mathbb{A}(\theta^{\mathrm{r}} + \lambda^\star) + \frac{\eta}{2}\|\theta^{\mathrm{r}} + \lambda^\star\|^2_A - \frac{\eta}{2}\|\theta^{\mathrm{r}} + \lambda\|^2_{\mathbb{A}}$$

$$= B_{Z^{\mathrm{r}}_0}\left(\theta^{\mathrm{r}}, \hat{\theta}^{\mathrm{r}}(k)\right) + (\lambda - \lambda^\star)^\top \nabla Z^{\mathrm{r}}_0(\theta^{\mathrm{r}} + \lambda^\star) + Z^{\mathrm{r}}_0(\theta^{\mathrm{r}} + \lambda^\star) - Z^{\mathrm{r}}_0(\theta^{\mathrm{r}} + \lambda) \tag{23}$$

where we have introduced the Legendre function $Z^{\mathrm{r}}_0(\theta) = \frac{\eta}{2}\|\theta\|^2_{\mathbb{A}}$. We now have from (27) that

$$\sup_{\lambda\in\mathbb{R}^d}\left(\lambda^\top x_n - B_{Z^{\mathrm{r}}_n,\theta^{\mathrm{r}}}(\lambda)\right)$$

$$= B^\star_{Z^{\mathrm{r}}_n,\theta^{\mathrm{r}}}(x_n) = B^\star_{Z^{\mathrm{r}}_n,\theta^{\mathrm{r}}}\left(\nabla Z^{\mathrm{r}}_n\left(\hat{\theta}^{\mathrm{r}}(n)\right) - \nabla Z^{\mathrm{r}}_n(\theta^{\mathrm{r}})\right) = B_{Z^{\mathrm{r}}n}\left(\theta^{\mathrm{r}}, \hat{\theta}^{\mathrm{r}}(n)\right).$$

Further, any optimal $\lambda$ must satisfy

$$\nabla Z^{\mathrm{r}}_n(\theta^{\mathrm{r}} + \lambda) - \nabla Z^{\mathrm{r}}_n(\theta^{\mathrm{r}}) = x_n \Longrightarrow \nabla Z^{\mathrm{r}}_n(\theta^{\mathrm{r}} + \lambda) = \nabla Z^{\mathrm{r}}_n\left(\hat{\theta}^{\mathrm{r}}(n)\right).$$

One possible solution is $\lambda = \lambda^\star$. Now, since $Z^{\mathrm{r}}_n$ is strictly convex, the supremum is indeed attained at $\lambda = \lambda^\star$. We then have

$$\lambda^\top x_n - B_{Z^{\mathrm{r}}_n,\theta^{\mathrm{r}}}(\lambda)$$

$$= \lambda^\top x_n - B_{Z^{\mathrm{r}}_n,\theta^{\mathrm{r}}}(\lambda) + B_{Z^{\mathrm{r}}_n}\left(\theta^{\mathrm{r}}, \hat{\theta}^{\mathrm{r}}(n)\right) - \left(\lambda^\star x_n - B_{Z^{\mathrm{r}}_n,\theta^{\mathrm{r}}}(\lambda^\star)\right)$$

$$= B_{Z^{\mathrm{r}}_n}\left(\theta^{\mathrm{r}}, \hat{\theta}^{\mathrm{r}}(n)\right) + (\lambda - \lambda^\star)^\top \nabla Z^{\mathrm{r}}_n(\theta^{\mathrm{r}} + \lambda^\star) + B_{Z^{\mathrm{r}}_n,\theta^*}(\lambda^\star) - B_{Z^{\mathrm{r}}_n,\theta^*}(\lambda)$$

$$\quad - (\lambda - \lambda^\star)^\top \nabla Z^{\mathrm{r}}_n(\theta^{\mathrm{r}})$$

$$= B_{Z^{\mathrm{r}}_n}\left(\theta^{\mathrm{r}}, \hat{\theta}^{\mathrm{r}}(n)\right) + (\lambda - \lambda^\star)^\top \nabla Z^{\mathrm{r}}_n(\theta^{\mathrm{r}} + \lambda^\star) + Z^{\mathrm{r}}_n(\theta^{\mathrm{r}} + \lambda^\star) - Z^{\mathrm{r}}_n(\theta^{\mathrm{r}} + \lambda) \tag{24}$$

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

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

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

Therefore, we have from (5) that

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

An application of Markov's inequality now yields

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

**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, $M_\infty = \lim_{n\to\infty} M_n$ 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

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

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 for any random stopping time $N < \infty$. To complete the proof, we now employ a random stopping time construction as in Abbasi-Yadkori et al. (2011)

We define a random stopping time $N$ by

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

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

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

which concludes the proof of the first part.

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

$$\mathrm{KL}_{s,a}\left(\hat{\theta}^{\mathtt{r}}(n), \theta'\right) = \frac{1}{2}\sum_{i,j=1}^{d}\left(\theta' - \hat{\theta}^{\mathtt{r}}(n)\right)_i \mathbb{V}\mathrm{ar}^{\theta}_{s,a}(r) \times \varphi(s,a)^\top A_i^\top B B^\top A_j \varphi(s,a)\left(\theta' - \hat{\theta}^{\mathtt{r}}(n)\right)_j \tag{25}$$

Now (25) implies that

$$\sum_{t=1}^{n} \mathrm{KL}_{s_t,a_t}\left(\hat{\theta}^{\mathtt{r}}(n), \theta'\right) \leq \frac{\beta}{2}\sum_{t=1}^{n}\sum_{i,j=1}^{d}\left(\theta' - \hat{\theta}^{\mathtt{r}}(n)\right)_i \varphi(s_t,a_t)^\top A_i^\top A_j \varphi(s_t,a_t)\left(\theta' - \hat{\theta}^{\mathtt{r}}(n)\right)_j$$

$$= \frac{\beta^{\mathtt{r}}}{2}\left\|\theta' - \hat{\theta}^{\mathtt{r}}(n)\right\|_{\sum_{t=1}^{n} G_{s_t,a_t}}^2,$$

where $\beta^{\mathtt{r}} := \lambda_{\max}\left(BB^\top\right) \times \sup_{\theta,s,a} \mathbb{V}\mathrm{ar}^\theta_{s,a}(r)$ 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

$$C_{\mathrm{A},n} \leq \frac{\int_{\mathbb{R}^d} \exp\left(-\frac{\eta}{2}\|\theta'\|_{\mathrm{A}}^2\right) d\theta'}{\int_{\mathbb{R}^d} \exp\left(-\frac{1}{2}\left\|\theta' - \hat{\theta}^{\mathtt{r}}(n)\right\|_{\left(\beta^{\mathtt{r}}\sum_{t=1}^{n} G_{s_t,a_t} + \eta\mathrm{A}\right)}^2\right) d\theta'}$$

$$= \frac{(2\pi)^{d/2}}{\det(\eta\mathbb{A})^{1/2}} \times \frac{\det\left(\beta^{\mathtt{r}}\sum_{t=1}^{n} G_{s_t,a_t} + \eta\mathbb{A}\right)^{1/2}}{(2\pi)^{d/2}} = \det\left(I + \beta^{\mathtt{r}}\eta^{-1}\mathbb{A}^{-1}\sum_{t=1}^{n} G_{s_t,a_t}\right),$$

which completes the proof of the second part.

$\square$

**Corollary 10.** *Here also, the theorem implies a euclidean control. With probability at least $1 - \delta$ uniformly over $k \in \mathbb{N}$*

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

*where $\beta^r(k, \delta) \overset{\text{def}}{=} \beta^r_{(k-1)H}(\delta) = \frac{2}{2}B_A^2 + \log\left(2C_{A,k}^r/\delta\right)$.*

## C.3 Gaussian concentration and anti-concentration

**Lemma 11** (Gaussian concentration, ref. Appendix A in [AL17]). *Let $\bar{\xi}_{tk} \sim \mathcal{N}(0, H\nu_k(\delta)\Sigma_{tk}^{-1})$. For any $\delta > 0$, with probability $1 - \delta$*

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

*for some absolute constant $c$.*

**Lemma 12** (Gaussian anti-concentration, ref. Appendix A in [AL17]). *Let $\xi \sim \mathcal{N}(0, I_d)$, for any $u \in \mathbb{R}^d$ with $\|u\| = 1$, we have:*

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

*where $\Phi$ is the normal CDF.*

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

# D  Technical results

## D.1  A transportation lemma

For any function $f : \mathcal{X} \to \mathbb{R}$, 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 distribution $P$ supported on the set $\mathcal{X}$, let $\mathbb{E}_P[f] := \mathbb{E}_P[f(X)]$ and $\mathbb{V}_P[f] :=$
$\mathbb{V}_P[f(X)] = \mathbb{E}_P\left[f(X)^2\right] - \mathbb{E}_P[f(X)]^2$ denote the mean and variance of the random variable
$f(X)$, respectively. We now state the following transportation inequalities, which can be adapted
from [BLM13] (Lemma 4.18).

**Lemma 13.** *(Transportation inequalities) Assume $f$ is such that $S(f)$ and $\mathbb{V}_P[f]$ are finite. Then it*
*holds*

$$\forall Q \ll P, \quad \mathbb{E}_Q[f] - \mathbb{E}_P[f] \leq \sqrt{2\mathbb{V}_P[f]\mathrm{KL}(Q,P)} + \frac{2S(f)}{3}\mathrm{KL}(Q,P)$$

$$\forall Q \ll P, \quad \mathbb{E}_P[f] - \mathbb{E}_Q[f] \leq \sqrt{2\mathbb{V}_P[f]\mathrm{KL}(Q,P)}$$

## D.2  Bregman divergence

For a Legendre function $F : \mathbb{R}^d \to \mathbb{R}$, 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

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

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

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

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

$$B_F(\theta', \theta) = B_{F,\theta'}^\star\left(\nabla F(\theta) - \nabla F(\theta')\right) \tag{27}$$

To see this, we observe that

$$B_{F,\theta'}^\star\left(\nabla F(\theta) - \nabla F(\theta')\right)$$
$$= \sup_{\lambda \in \mathbb{R}^d} \lambda^\top \left(\nabla F(\theta) - \nabla F(\theta')\right) - \left[F(\theta' + \lambda) - F(\theta') - \lambda^\top \nabla F(\theta')\right]$$
$$= \sup_{\lambda \in \mathbb{R}^d} \lambda^\top \nabla F(\theta) - F(\theta' + \lambda) + F(\theta').$$

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 $\lambda = \theta - \theta'$. Plugin-in this
value, we obtain

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

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 Bregman divergence.

## D.3  Properties of the bilinear exponential family

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

### D.3.1  Derivatives

**Lemma 14.** *(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:*

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

*And*

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

*Proof.* We prove the lemma as follows

$$\left(\nabla_i Z_{s,a}^{\texttt{p}}\right)(\theta) = \int_{\mathcal{S}} \psi\left(s'\right)^{\top} A_i \varphi(s,a) \frac{\exp\left(\sum_{i=1}^{d} \theta_i \psi\left(s'\right)^{\top} A_i \varphi(s,a)\right)}{\int_{\mathcal{S}} \exp\left(\sum_{i=1}^{d} \theta_t \psi\left(s'\right)^{\top} A_i \varphi(s,a)\right) ds'} ds'$$

$$= \mathbb{E}_{s,a}^{\theta}\left[\psi\left(s'\right)\right]^{\top} A_i \varphi(s,a)$$

$$\left(\nabla_i Z_{s,a}^{\texttt{r}}\right)(\theta) = \int_{\mathcal{S}} r B^{\top} A_i \varphi(s,a) \frac{\exp\left(r \sum_{i=1}^{d} \theta_i B^{\top} A_i \varphi(s,a)\right)}{\int_{\mathcal{S}} \exp\left(r \sum_{i=1}^{d} \theta_i B^{\top} A_i \varphi(s,a)\right) dr} dr$$

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

$\square$

**Lemma 15.** *(Hessians) The entries of the Hessians of the log partition functions are given by*

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

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

*Similarly,*

$$\left(\nabla_{i,j}^2 Z_{s,a}^{\texttt{r}}\right)(\theta) = \mathbb{V}\text{ar}_{s,a}^{\theta}(r) \times \varphi(s,a)^{\top} A_i^{\top} BB^{\top} A_j \varphi(s,a),$$

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

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

$$\left(\nabla_{i,j}^2 Z_{s,a}^{\texttt{p}}\right)(\theta) = \int_{\mathcal{S}} \psi\left(s'\right)^{\top} A_i \varphi(s,a) \psi\left(s'\right)^{\top} A_j \varphi(s,a) \frac{\exp\left(\sum_{i=1}^{d} \theta_i \psi\left(s'\right)^{\top} A_i \varphi(s,a)\right)}{\int_{\mathcal{S}} \exp\left(\sum_{i=1}^{d} \theta_i \psi\left(s'\right)^{\top} A_i \varphi(s,a)\right) ds'} ds'$$

$$- \int_{\mathcal{S}} \psi\left(s'\right)^{\top} A_i \varphi(s,a) \frac{\exp\left(\sum_{i=1}^{d} \theta_i \psi\left(s'\right)^{\top} A_i \varphi(s,a)\right)}{\int_{\mathcal{S}} \exp\left(\sum_{i=1}^{d} \theta_i \psi\left(s'\right)^{\top} A_i \varphi(s,a)\right) ds'} ds' \left(\nabla_j Z_{s,a}\right)(\theta)$$

$$= \mathbb{E}_{s,a}^{\theta}\left[\psi\left(s'\right)^{\top} A_i \varphi(s,a) \psi\left(s'\right)^{\top} A_j \varphi(s,a)\right]$$

$$- \mathbb{E}_{s,a}^{\theta}\left[\psi\left(s'\right)^{\top} A_i \varphi(s,a)\right] \mathbb{E}_{s,a}^{\theta}\left[\psi\left(s'\right)^{\top} A_j \varphi(s,a)\right]$$

$$= \varphi(s,a)^{\top} A_i^{\top} \left(\mathbb{E}_{s,a}^{\theta}\left[\psi\left(s'\right)\psi\left(s'\right)^{\top}\right] - \mathbb{E}_{s,a}^{\theta}\left[\psi\left(s'\right)\right]\mathbb{E}_{s,a}^{\theta}\left[\psi\left(s'\right)^{\top}\right]\right) A_j \varphi(s,a)$$

$$= \varphi(s,a)^{\top} A_i^{\top} \mathbb{C}_{s,a}^{\theta}\left[\psi\left(s'\right)\right] A_j \varphi(s,a),$$

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

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

The proof of the form of the Hessian for the reward partition function follows the same steps as above. $\square$

**Lemma 16.** *(KL Divergences) For any two* $\theta, \theta'$ *and for some pair* $(s,a)$,

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

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

*Similarly*

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

*Proof.* We start by writing:

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

then

$$\mathrm{KL}\left(P_\theta^{\mathrm{p}}(\cdot \mid s,a), P_{\theta'}^{\mathrm{p}}(\cdot \mid s,a)\right) = \sum_{i=1}^d (\theta_i - \theta_i')\,\mathbb{E}_{s,a}^\theta\left[\psi\left(s'\right)\right]^\top A_i\varphi(s,a) - Z_{s,a}^{\mathrm{p}}(\theta) + Z_{s,a}^{\mathrm{p}}(\theta')$$

$$= \frac{1}{2}\left(\theta - \theta'\right)^\top \left(\nabla^2 Z_{s,a}^{\mathrm{p}}\right)(\tilde\theta)\left(\theta - \theta'\right),$$

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

The proof of the form of the KL divergence for the reward follows the same steps as above. $\qquad\square$

### D.3.2 A transportation lemma for rewards

**Lemma 17.** *We provide a closed-form formula for the difference of expected rewards under two distinct parameters:*

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

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

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

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

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

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:

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

We know that $\left(\nabla_{i,j}^2 Z_{s,a}^{\mathrm{r}}\right)(\theta) = \mathbb{V}\mathrm{ar}_{s,a}^\theta(r) \times \varphi(s,a)^\top A_i^\top BB^\top A_j\varphi(s,a)$, choosing $\theta^{\mathrm{r}'} = \theta_1 - \theta_2$ we find:

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

$\qquad\square$

### D.4 Elliptical potentials and elliptical lemma

### D.4.1 Elliptical lemma

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

First, consider the notations: $G_{s,a} := (\varphi(s,a)^\top A_i^\top A_j\varphi(s,a))_{1\le i,j\le d}$, $\quad \bar{G}_n^{\mathrm{e}} \equiv \bar{G}_{(k-1)H}^{\mathrm{e}} := G_n + (\alpha^{\mathrm{e}})^{-1}\eta A$, and $G_n \equiv G_{(k-1)H} := \sum_{\tau=1}^{k-1}\sum_{h=1}^H G_{s_s^\tau, a_h^\tau}$. Where $\mathrm{e}$ represents either $\mathrm{r}$ or $\mathrm{p}$, we omit the superscript $\mathrm{e}$ w.l.o.g in the rest of this section.

**Lemma 18.** *(Elliptical lemma and variant for bounded potentials) Let $c \in \mathbb{R}^+$, we can bound the sum of feature norms as follows*

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

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

*Further, we have*

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

*Proof.* First we have

$$\|\bar{G}_n^{-1/2} G_{s,a} \bar{G}_n^{-1/2}\| = \sqrt{\text{tr}(\bar{G}_n^{-1/2} G_{s,a} \bar{G}_n^{-1/2} \bar{G}_n^{-1/2} G_{s,a} \bar{G}_n^{-1/2})}$$
$$\leq \text{tr}(\bar{G}_n^{-1/2} G_{s,a} \bar{G}_n^{-1/2}) = \text{tr}(\bar{G}_n^{-1} G_{s,a}) = \text{tr}(\boldsymbol{a}_h^\top \bar{G}_n^{-1} \boldsymbol{a}_h)$$

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

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

$$\min\{c, \sum_{h=1}^H \text{tr}(\boldsymbol{a}_h^\top \bar{G}_n^{-1} \boldsymbol{a}_h)\} = \min\{c, \sum_{h=1}^H \sum_{i=1}^d \lambda_{h,i}\}$$

$$\leq \frac{c}{\log(1+c)} \sum_{h=1}^H \sum_{i=1}^d \log(1 + \lambda_{h,i}) \qquad \text{(log is concave)}$$

$$\leq \frac{c}{\log(1+c)} \sum_{h=1}^H \log(\prod_{i=1}^d 1 + \lambda_{h,i}) = \frac{c}{\log(1+c)} \sum_{h=1}^H \log\det(I + \boldsymbol{a}_h^\top \bar{G}_n^{-1} \boldsymbol{a}_h)$$

$$\leq \frac{c}{\log(1+c)} \log\left(\frac{\det(\bar{G}_n + \sum_{h=1}^H G_{s_h,a_h})}{\det(\bar{G}_n)}\right)$$

where the last line follows from the matrix determinant lemma:

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

Therefore:

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

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

$$\sum_{t=1}^T \log\frac{\det\left(\bar{G}_{n+H}\right)}{\det\left(\bar{G}_n\right)} = \sum_{t=1}^T \log\frac{\det\left(\bar{G}_{tH}\right)}{\det\left(\bar{G}_{(t-1)H}\right)} = \log\frac{\det\left(\bar{G}_{TH}\right)}{\det\left(\bar{G}_0\right)}$$

$$= \log\frac{\det\left(\bar{G}_N\right)}{\det\left((\alpha^{\text{p}})^{-1} \eta \mathbb{A}\right)} = \log\det\left(I + \alpha\eta^{-1} \,\text{A}^{-1} G_N\right)$$

$$\leq d\log\left(1 + \frac{\alpha^{\text{p}}\eta^{-1}}{d} \text{tr}\left(\mathbb{A}^{-1} G_n\right)\right) \qquad \text{(Trace-determinant (or AM-GM) inequality)}$$

$$\leq d\log\left(1 + \alpha^{\text{p}}\eta^{-1} B_{\varphi,\mathbb{A}} n\right)$$

This concludes the proof of the first result.

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

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

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

**Remark 7.** *We can also write from the lemma in terms of* $\|(A_i\varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^r)^{-1}}$ *by skipping the norm upper bound at the beginning of the proof:*

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

*and*

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

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

### D.4.2 Elliptical potentials: finite number of large feature norms (contribution)

**Lemma 19.** *(Worst case elliptical potentials, adaptation of Exercise 19.3 [LS20] for matrices) Let $V_0 = \lambda I$ and $a_1,\ldots,a_n \in \mathbb{R}^{d\times p}$ 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

$$\left(\frac{d\lambda+|\mathcal{T}|L^2}{d}\right)^d \geq \left(\frac{\text{trace}\,(G_n)}{d}\right)^d$$

$$\geq \det(G_n) \qquad\qquad \text{(Trace-determinant inequality)}$$

$$= \det(V_0)\prod_{t\in T}\left(1+\|a_t\|_{G_{t-1}^{-1}}^2\right)$$

$$\geq \det(V_0)\prod_{t\in T}\left(1+\|a_t\|_{V_{t-1}^{-1}}^2\right)$$

$$\geq \lambda^d 2^{|\mathcal{T}|}$$

where the third line follows from the matrix determinant lemma:

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

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

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

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

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

$\square$

## 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

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

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 $\mathcal{N}$ is the Gaussian distribution, $\mathcal{U}$ is the Uniform distribution, and $p(w,b)$ is a coupling among them.

Notice that this provides an alternative approach to decompose the transition kernel and obtain linearity of the value 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 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 probability at least $1 - 2\exp(-N\epsilon^2/4)$. [RR07] proves a stronger result: it provides an algorithm approximating the Gaussian kernel for which the following uniform convergence bound holds.

**Lemma 20.** *Let $\mathcal{M}$ be a compact subset of $\mathcal{R}^p$ with diameter* $\mathrm{diam}(\mathcal{M})$. *Then, using the explicit mapping $\mathbf{z}$ defined in Algorithm 1 in [RR07] with $N$ samples, we have*

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

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

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

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

Let's call $\hat{V}_h$ 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:

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

Here, we use the facts that $\mathbb{S}\left(V_{\hat{\theta}, \tilde{\theta}^\times, h}\right) \leq dH^{3/2}$ (*cf.* Section B.2) and the error in approximating $P$ is bounded by $\epsilon$, 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

$$\hat{V}_h^\pi - V_h^\pi \leq \int_{s'} (\hat{P} - P)(s' \mid s, a)\hat{V}_{h+1}^\pi(s') \, \mathrm{d}s' + \int_{s'} P(s' \mid s, a)(\hat{V}_{h+1}^\pi - V_{h+1}^\pi)(s') \, \mathrm{d}s'$$

$$= \int_{s'} (\hat{P} - P)(s' \mid s, a)(V_{h+1}^\pi + \hat{V}_{h+1}^\pi - V_{h+1}^\pi) \, \mathrm{d}s' + \int_{s'} P(s' \mid s, a)(\hat{V}_{h+1}^\pi - V_{h+1}^\pi)(s') \, \mathrm{d}s'$$

$$\leq \epsilon\left(dH^{3/2} + \sum_{j=1}^{h+1} \epsilon^j \alpha_{h+1,j}\right) + \sum_{j=1}^{h+1} \epsilon^j \alpha_{h+1,j}$$

$$\leq \epsilon(dH^{3/2} + \alpha_{h+1,1}) + \sum_{j=2}^{h+1} \epsilon^j(\alpha_{h+1,j-1} + \alpha_{h+1,j}) + \epsilon^{h+2}\alpha_{h+1,h+1}$$

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

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

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

**Remark 8.** *The reader might be tempted to combine the finite approximation using RFT with algorithms from the linear reinforcement learning literature [JYWJ20]. However, note that the dimensionality of the linear space induced by RFT is polynomial in $H$ and $K$. Consequently, applying algorithms designed with the assumption of linear value function would incur a linear regret.*