



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

Reda Ouhamma, Debabrota Basu, Odalric-Ambrym Maillard

► To cite this version:

Reda Ouhamma, Debabrota Basu, Odalric-Ambrym Maillard. Bilinear Exponential Family of MDPs: Frequentist Regret Bound with Tractable Exploration & Planning. EWRL 2022 – European Workshop on Reinforcement Learning, Sep 2022, Milan, Italy. hal-03790997v1

HAL Id: hal-03790997

<https://hal.science/hal-03790997v1>

Submitted on 30 Sep 2022 (v1), last revised 26 Oct 2023 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

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

Reda Ouhamma* **Debabrota Basu** **Odalric-Ambrym Maillard**

Univ. Lille, Inria, CNRS,
Centrale Lille, UMR 9189 CRISAL,
F-59000 Lille, France

reda.ouhamma@gmail.com , debabrota.basu@inria.fr , odalric.maillard@inria.fr

Abstract

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

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 $t = 1, \dots, H$, the agent draws action a_t^k from a (possibly time-dependent) policy $\pi_t(s_t^k)$, observes the reward $r(s_t^k, a_t^k) \in [0, 1]$, and transits to a state $s_{t+1}^k \sim \mathbb{P}(\cdot | s_t^k, a_t^k)$ according to the transition function \mathbb{P} . The performance of a policy π is measured by the total expected reward V_1^π 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

*<https://redaouhamma.github.io/>

precisely as required. This tension between learning the unknown environment and reward maximization is quantified as *regret*: the typical performance measure of an episodic RL algorithm. *Regret* is defined as the difference between the *expected cumulative reward* or *value* collected by the optimal agent that knows the environment and the expected cumulative reward or value obtained by an agent that has to learn about the unknown environment. Formally, the regret over K episodes is

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

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] and 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 planning² 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 $\int \mathbb{P}(s' | s, a) V_h(s)$ 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 design 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

²By 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 MDPs with functional representations.

Algorithm	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
EXP-UCRL [CGM21]	$\sqrt{d^2 H^4 K}$ (Frequentist)	✗	✗	N.A	Bilinear Exp Family (BEF) known rewards
SMRL [LLS ⁺ 21]	$\sqrt{d^2 H^4 K}$	✗	✗	N.A	BEF, known rewards
UCRL-VTR [AJS ⁺ 20]	$\sqrt{d^2 H^4 K}$	✗	✗	N.A	Linear mixture model
\mathcal{F} -PHE-LSVI [ICN ⁺ 21]	$\text{poly}(d_E H) \sqrt{K H}$	✓	✗	✗	Eluder dimension, Tabular
PHE-LSVI (linear-RL)	$\sqrt{d^3 H^4 K}$				Anti-concentration
UC-MatrixRL [YW20]	$\sqrt{d^2 H^5 K}$	✗	✗	N.A	Linear factor MDP
OPT-RLSVI [ZBB ⁺ 20]	$\sqrt{d^4 H^5 K}$	✓	✓	✗	Linear V
BEF-RLSVI (this work)	$\sqrt{d^3 H^3 K}$	✓	✓	✓	Bilinear Exp Family

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.

3. *Analysis*: We carefully develop an analysis of BEF-RLSVI that yields $\tilde{O}(\sqrt{d^3 H^3 K})$ 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 (BEF) model coined in [CGM21], extend it to parametric rewards, and we state a novel observation about linearity of this representation.

Bilinear exponential family. We consider transitions and rewards from the BEF. Specifically,

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

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

where $\varphi \in (\mathbb{R}_+^q)^{\mathcal{S} \times \mathcal{A}}$ and $\psi \in (\mathbb{R}_+^p)^{\mathcal{S}}$ are known feature mappings, and $B \in \mathbb{R}^p$ is a known matrix. The reward and transition parameters are $\theta^p, \theta^r \in \mathbb{R}^d$. $M_{\theta^p} \stackrel{\text{def}}{=} \sum_{i=1}^d \theta_i A_i$, where $(A_i)_{1 \leq i \leq d}$ are known matrices. The log partition function: $Z_{s,a}^p(\theta^p) \stackrel{\text{def}}{=} \log \int_{\mathcal{S}} \exp(\psi(\tilde{s})^\top M_{\theta^p} \varphi(s, a)) d\tilde{s}$, and Z^r is defined similarly. Finally, we emphasize a minor difference with the original BEF model: like [LLS⁺21], we omit a base measure of the form $h(s, \tilde{s}, a)$ from the model, note that all the examples provided in [CGM21] still hold with this slight restriction.

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

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

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

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

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

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

where \mathcal{H} is the RKHS associated with the kernel, $\mu^p(s') = (2\pi)^{-p/2} k(\psi(s'), \cdot) \exp(\|\psi(s')\|^2/2)$, and $\phi^p(s, a) = k(M_{\theta^p}^\top \varphi(s, a), \cdot) \exp(\|M_{\theta^p} \varphi(s, a)\|^2/2 - Z_{s,a}(\theta^p))$. 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 ($f^*(s, a)$) and known variance (σ^2). Actually, linear f^* 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:

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

We must approximate the integral at the R.H.S. for $(s, a) \in \mathcal{S} \times \mathcal{A}$. For a tabular MDP with $|\mathcal{S}|$ states and $|\mathcal{A}|$ actions, we need to evaluate $(Q_h^\pi)_{h \in [H]}$, i.e. to approximate $|\mathcal{S}| \times |\mathcal{A}| \times H$ integrals per episode, which can be very expensive. However, with the linear transition model of Equation (4), although ϕ^p and μ^p are infinite dimensional, we show in Section 4 (§ planning) that the planning complexity becomes polynomial in the problem parameters.

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 δ , constants α^p, η and $(x_k)_{k \in [K]} \in \mathbb{R}^+$
 - 2: **for** episode $k = 1, 2, \dots$ **do**
 - 3: Observe initial state s_1^k
 - 4: Sample noise $\xi_k \sim \mathcal{N}(0, x_k(\bar{G}_k^p)^{-1})$ such that

$$\bar{G}_k^p = \frac{\eta}{\alpha^p} \mathbb{A} + \sum_{\tau=1}^{k-1} \sum_{h=1}^H (\varphi(s_h^\tau, a_h^\tau)^\top A_i^\top A_j \varphi(s_h^\tau, a_h^\tau))_{i,j \in [d]}$$
 - 5: Perturb reward parameter: $\tilde{\theta}^r(k) = \hat{\theta}^r(k) + \xi_k$
 - 6: Compute $(Q_{\hat{\theta}^p, \tilde{\theta}^r, h}^k)_{h \in [H]}$ via Bellman-backtracking, see Algorithm 2
 - 7: **for** $h = 1, \dots, H$ **do**
 - 8: Pull action $a_h^k = \arg \max_a Q_{\hat{\theta}^p, \tilde{\theta}^r, h}(s_h^k, a)$
 - 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 $\hat{\theta}^p(k), \hat{\theta}^r(k)$, see Equation (5) and Equation (7)
 - 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 involve high-dimensional optimization.

Algorithm 2 Bellman Backtracking

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

We can approximate Line 3 of Algorithm 2 with $\mathcal{O}(pH^3K \log(HK))$ 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 $\alpha^p, \alpha^r, \beta^p, \beta^r > 0$, such that the representation model satisfies:*

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

where $C_{s,a}^\theta [\psi(s')] \triangleq \mathbb{E}_{s' \sim \mathbb{P}_\theta|s,a} [\psi(s') \psi(s')^\top] - \mathbb{E}_{s' \sim \mathbb{P}_\theta|s,a} [\psi(s')] \mathbb{E}_{s' \sim \mathbb{P}_\theta|s,a} [\psi(s')^\top]$ and $\mathbb{V}_{s,a}^\theta(r) \triangleq \left(\mathbb{E}_{s,a}^\theta[r^2] - \mathbb{E}_{s,a}^\theta[r]^2 \right)$ is the variance of the reward under θ .

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

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

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

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

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, $C_d \triangleq \frac{3d}{\log(2)} \log \left(1 + \frac{\alpha^r \|\mathbb{A}\|_2^2 B_{\varphi, \mathbb{A}}^2}{\eta \log(2)} \right)$, Φ 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.

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

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

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

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

4 Algorithm design: building blocks of BEF-RLSVI

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

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

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

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

$$\sum_{t=1}^k \sum_{h=1}^H \left(\psi(s_{h+1}^t) - \mathbb{E}_{s_h^t, a_h^t}^{\hat{\theta}_k^p} [\psi(s')] \right)^{\top} A_i \varphi(s_h^t, a_h^t) = \eta \nabla_i \text{pen}(\hat{\theta}_k^p). \quad (5)$$

Equation (5) can be solved in closed form for simple distributions, like Gaussian, but it can involve integral approximations for other distribution (cf. Appendix F). We estimate the parameter for reward, i.e. θ_r , similarly

$$\hat{\theta}^r(k) \in \arg \min_{\theta \in \mathbb{R}^d} \sum_{t=1}^k \sum_{h=1}^H -\log \mathbb{P}_{\theta}(r_t | s_h^t, a_h^t) + \eta \text{pen}(\theta), \quad (6)$$

$$\implies \sum_{t=1}^k \sum_{h=1}^H \left(r_t - \mathbb{E}_{s_h^t, a_h^t}^{\hat{\theta}_k^r} [r] \right) B^{\top} A_i \varphi(s_h^t, a_h^t) = \eta \nabla_i \text{pen}(\hat{\theta}_k^r) \quad \forall i \in [d]. \quad (7)$$

Exploration. A significant challenge in RL is handling exploration in continuous spaces. The majority of the literature is split between intractable, upper confidence bound-style optimism or Thompson sampling algorithms with high-dimensional posterior and guarantees only in terms of Bayesian regret. In BEF-RLSVI, we adopt the approach of reward perturbation motivated by the RLSVI-framework [ZBB⁺20, OVRW16]. We show that perturbing the reward estimation can guarantee optimism with a constant probability, *i.e.* there exists $\nu \in (0, 1]$ such that for all $k \in [K]$ and $s_1^k \in \mathcal{S}$,

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

[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. Using BEF-RLSVI, we have at step h :

$$Q_{\hat{\theta}^p, \hat{\theta}^r, h}^\pi(s, a) = \mathbb{E}_{\tilde{r}}[r(s, a)] + \left\langle \phi^p(s, a), \int_{\mathcal{S}} \mu^p(\tilde{s}) V_{\hat{\theta}^p, \hat{\theta}^r, h+1}^\pi(\tilde{s}) d\tilde{s} \right\rangle. \quad (8)$$

Then, we select the best action greedily using dynamic programming to compute $Q_h(s, a)$. Although our model yields infinite dimensional ϕ^p and ψ^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 ϕ^p and ψ^p . We proceed in three steps. **1)** We have with high probability $\mathbb{S}(V_{\hat{\theta}^p, \hat{\theta}^r, h}) \leq dH^{3/2}$ (Section 5). **2)** If we have a uniform ϵ -approximation of \mathbb{P}_{θ^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 ϵ error and with features of dimensions $\mathcal{O}(p\epsilon^{-2} \log \frac{1}{\epsilon^2})$, where p is dimension of ψ . Associating these three elements and choosing $\epsilon = 1/\sqrt{(H^2K)}$, we establish our claim.

Remark 2. The observation of linearity (*cf.* Equation. (8) and Line 3) does not reduce BEF MDPs to linear MDPs because the former holds in an RKHS. Also, linearity is not in the representation parameter. Therefore, linear RL algorithms do not readily solve the BEF MDPs.

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

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

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.

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

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 $H \sum_{h=1}^H \sqrt{2 \text{KL}_{s_{hk}, a_{hk}}(\theta^p, \hat{\theta}^p)}$. We notice that since the reward parameter is exact, the bound is actually $H \min\{1, \sum_{h=1}^H \sqrt{2 \text{KL}_{s_{hk}, a_{hk}}(\theta^p, \hat{\theta}^p)}\}$. 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 $\|(A_i \varphi(s_h^k, a_h^k))_{i \in [d]}\|_{(\bar{G}_k^p)^{-1}}$, 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": $\left\{k \in [K], \exists h : \|(A_i \varphi(s_h^k, a_h^k))_{i \in [d]}\|_{(\bar{G}_k^p)^{-1}} \geq 1\right\}$, 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 $\mathcal{O}(d^{3/2} H^2)$ regret, independent of K . Therefore, our algorithm can forgo clipping for free.

Remark 3. *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,

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

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

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

where $c(n, \delta) = \left(\sqrt{\beta^p \beta^p(n, \delta) / \alpha^p} + \sqrt{\beta^r \beta^r(n, \delta) \min\{1, \alpha^p / \alpha^r\} / (2\alpha^r)} \right)$. Since $\xi_k \sim \mathcal{N}(0, x_k (\bar{G}_k^p)^{-1})$ and $x_k \geq H^2 c(n, \delta)^2$, we get $\mathbb{P} \left(V_{\hat{\theta}^p, \hat{\theta}^r, 1}^{\pi^*}(s_1) - V_{\theta^p, \theta^r, 1}^*(s_1) \geq 0 \right) \geq \Phi(-1)$, where Φ 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^* : Draw $(\tilde{\xi}_k)_{k \in [K]}$ i.i.d copies of $(\xi_k)_{k \in [K]}$ and define the event where optimism holds as $\bar{O}_k \triangleq \{V_{\hat{\theta}^p, \hat{\theta}^r, 1}(s_1^k) - V_1^*(s_1^k) \geq 0\}$. This implies that $V_1^*(s_1^k) \leq \mathbb{E}_{\tilde{\xi}_k | \bar{O}_k} [V_{\hat{\theta}^p, \hat{\theta}^r, 1}(s_1^k)]$.

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

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

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

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

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

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 work: functional representation in RL with regret and tractability

Our work extends the endeavor of using functional representations for regret minimization in continuous state-action MDPs. Now, we posit our contributions in existing literature.

Kernel value function representation. [AJS⁺20] studies MDPs with a linear mixtures model then extends to an RKHS setting, this generalizes our work and that of [YW20]. However, the paper proposes an Eluder-dimension analysis, for RKHS settings this leads to the result of [YW20], *i.e.* a regret $H \log(T)^d$ higher than for BEF-RLSVI. Recently, [HKLL21] shows that for RKHS, Eluder dimension and the information gain are strictly equivalent, which brings in the extra factor.

General functional representation. The Eluder dimension is a complexity measure often used to analyze RL with general function space, [HKLL21] asserts that "common examples of where it is known to be small are function spaces (vector spaces)". [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. In the episodic setting, [WSY20] leverage the UCB approach for tabular MDPs and function spaces with bounded Eluder dimension, this strategy achieves a and achieve a $\tilde{O}(\sqrt{d^4 H^2 T})$ regret for linear MDPs. [ICN⁺21] considers the same setting, proposes an RLSVI based algorithm, and achieves a $\tilde{O}(\sqrt{d^3 H^4 K})$ for linear MDPs. However, the latter assumes an oracle perturbing the estimation to achieve anti-concentration while maintaining a bounded covering number, which is a counter-intuitive mix of boundedness and anti-concentration. Indeed, [ZBB⁺20] studied the linear MDP case, and while it managed to design an ingenious clipping verifying previous assumptions, the method is extremely intricate and the proof is involved and unlikely to extend for general value function spaces. *To concertize our design, we focus on the general but explicit BEF of MDPs than any abstract representation. We also remove the requirement to clip with a novel analysis.*

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, 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 representation. In general, planning can incur high computational complexity if done naively. Specially, [DKWY19] shows that for some settings, even with a linear ϵ -approximation of the Q -function, a

planning procedure able to produce an ϵ -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. *Thus, we propose the first algorithm with tractable planning for BEF-family.*

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 $\mathcal{O}(pH^3K \log(HK))$) 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.

Acknowledgments and Disclosure of Funding

The authors acknowledge the funding of the French National Research Agency, the French Ministry of Higher Education and Research, Inria, the MEL and the I-Site ULNE regarding project R-PILOTE-19-004-APPRENF. R. Ouhamma also acknowledges support from Ecole polytechnique.

References

- [AJS⁺20] Alex Ayoub, Zeyu Jia, Csaba Szepesvari, Mengdi Wang, and Lin Yang. Model-based reinforcement learning with value-targeted regression. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 463–474. PMLR, 13–18 Jul 2020. (Cited on 3, 9)
- [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 8, 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)
- [CPH05] Miguel A Carreira-Perpinan and Geoffrey Hinton. On contrastive divergence learning. In *International workshop on artificial intelligence and statistics*, pages 33–40. PMLR, 2005. (Cited on 35)
- [DDG⁺19] Bo Dai, Hanjun Dai, Arthur Gretton, Le Song, Dale Schuurmans, and Niao He. Kernel exponential family estimation via doubly dual embedding. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2321–2330. PMLR, 2019. (Cited on 35)
- [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 Cappé, 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)
- [HKLL21] Kaixuan Huang, Sham M Kakade, Jason D Lee, and Qi Lei. A short note on the relationship of information gain and eluder dimension. *arXiv preprint arXiv:2107.02377*, 2021. (Cited on 9)
- [ICN⁺21] Haque Ishfaq, Qiwen Cui, Viet Nguyen, Alex Ayoub, Zhuoran Yang, Zhaoran Wang, Doina Precup, and Lin Yang. Randomized exploration in reinforcement learning with general value function approximation. In *International Conference on Machine Learning*, pages 4607–4616. PMLR, 2021. (Cited on 3, 7, 9)
- [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)
- [Jør83] Bent Jørgensen. Maximum likelihood estimation and large-sample inference for generalized linear and nonlinear regression models. *Biometrika*, 70(1):19–28, 1983. (Cited on 35)

- [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, 3, 9, 35)
- [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, 10)
- [Nea01] Radford M Neal. Annealed importance sampling. *Statistics and computing*, 11(2):125–139, 2001. (Cited on 35)
- [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 10)
- [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, 7)
- [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, 10, 34)
- [RY77] Gerald S Rogers and Dennis L Young. Explicit maximum likelihood estimators for certain patterned covariance matrices. *Communications in Statistics-Theory and Methods*, 6(2):121–133, 1977. (Cited on 35)
- [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, 10, 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, 10)
- [SSW21] Abhin Shah, Devavrat Shah, and Gregory Wornell. A computationally efficient method for learning exponential family distributions. *Advances in Neural Information Processing Systems*, 34:15841–15854, 2021. (Cited on 35)
- [VGB12] Shankar Vembu, Thomas Gartner, and Mario Boley. Probabilistic structured predictors. *arXiv preprint arXiv:1205.2610*, 2012. (Cited on 35)
- [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, 10)
- [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)
- [YW20] Lin Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning*, pages 10746–10756. PMLR, 2020. (Cited on 3, 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, 9)

Appendix

Table of Contents

A	Notations	15
B	Regret analysis	16
B.1	Estimation error	16
B.1.1	Transition estimation	16
B.1.2	Reward estimation	18
B.2	Learning error	20
B.2.1	Stochastic optimism	20
B.2.2	Controlling the learning error	22
	Upper bound on V_1^*	22
	Lower bound on $V_{\hat{\theta}_P, \hat{\theta}^*}$	22
	Combining the error bounds.	23
C	Concentrations	24
C.1	Concentration of the transition parameter	24
C.2	Concentration of the reward parameter (contribution)	25
	Step 1: Martingale construction.	25
	Step 2: Method of mixtures.	26
	Step 3: A stopped martingale and its control.	27
C.3	Gaussian concentration and anti-concentration	28
D	Technical results	29
D.1	A transportation lemma	29
D.2	Bregman divergence	29
D.3	Properties of the bilinear exponential family	29
D.3.1	Derivatives	29
D.3.2	A transportation lemma for rewards	31
D.4	Elliptical potentials and elliptical lemma	31
D.4.1	Elliptical lemma	31
D.4.2	Elliptical potentials: finite number of large feature norms (contribution)	33
E	Tractable planning with random Fourier transform	33
F	Tractable Maximum Likelihood estimation	35

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	$\stackrel{\text{def}}{=}$	number of steps in a given episode
K	$\stackrel{\text{def}}{=}$	number of episodes
T	$\stackrel{\text{def}}{=}$	KH , total number of steps
s_h^k	$\stackrel{\text{def}}{=}$	state at time h of episode k , denoted s_h when k is clear from context
a_h^k	$\stackrel{\text{def}}{=}$	action at time h of episode k , denoted a_h when k is clear from context
$r(s, a)$	$\stackrel{\text{def}}{=}$	realization of the reward in state s under action a
θ^p	$\stackrel{\text{def}}{=}$	parameter of the transition distribution, $\in \mathbb{R}^d$
θ^r	$\stackrel{\text{def}}{=}$	parameter of the reward distribution, $\in \mathbb{R}^d$
θ	$\stackrel{\text{def}}{=}$	$\in \mathbb{R}^d$ denotes either θ^r or θ^p , unless stated otherwise
$\hat{\theta}$	$\stackrel{\text{def}}{=}$	θ estimator with Maximum Likelihood unless stated otherwise
$\tilde{\theta}$	$\stackrel{\text{def}}{=}$	$\hat{\theta} + \xi$ where ξ is a chosen noise. Perturbed estimation of θ .
$[\theta_1, \theta_2]$	$\stackrel{\text{def}}{=}$	the d -dimensional ℓ_∞ hypercube joining θ_1 and θ_2
\mathbb{P}_{θ^p}	$\stackrel{\text{def}}{=}$	transition under the exponential family model with parameter θ^p
ψ	$\stackrel{\text{def}}{=}$	feature function, $\in (\mathbb{R}_+^p)^S$
φ	$\stackrel{\text{def}}{=}$	feature function, $\in (\mathbb{R}_+^q)^{S \times \mathcal{A}}$
B	$\stackrel{\text{def}}{=}$	p -dimensional vector
M_θ	$\stackrel{\text{def}}{=}$	$\sum_{i=1}^d \theta_i A_i$, where A_i are $p \times q$ matrices.
Z^r	$\stackrel{\text{def}}{=}$	the rewards' log partition function
Z^p	$\stackrel{\text{def}}{=}$	the transitions' log partition function
\mathcal{H}	$\stackrel{\text{def}}{=}$	Hilbert space where we decompose transitions
μ^p	$\stackrel{\text{def}}{=}$	feature function after decomposition, $\in (\mathbb{R}_+)^{S \times \mathcal{H}}$
ϕ^p	$\stackrel{\text{def}}{=}$	feature function after decomposition, $\in (\mathbb{R}_+)^{S \times \mathcal{A} \times \mathcal{H}}$
$G_{s,a}$	$\stackrel{\text{def}}{=}$	$(\varphi(s, a)^\top A_i^\top A_j \varphi(s, a))_{i,j \in [d]}$
\bar{G}_k^r	$\stackrel{\text{def}}{=}$	$\bar{G}_{(k-1)h}^r = \frac{\eta}{\alpha^r} \mathbb{A} + \sum_{\tau=1}^{k-1} \sum_{h=1}^H G_{s_h^\tau, a_h^\tau}$
\bar{G}_k^p	$\stackrel{\text{def}}{=}$	$\bar{G}_{(k-1)h}^p = \frac{\eta}{\alpha^p} \mathbb{A} + \sum_{\tau=1}^{k-1} \sum_{h=1}^H G_{s_h^\tau, a_h^\tau}$
$\mathbb{C}_{s,a}^\theta [\psi(s')]$	$\stackrel{\text{def}}{=}$	$\mathbb{E}_{s,a}^\theta [\psi(s') \psi(s')^\top] - \mathbb{E}_{s,a}^\theta [\psi(s')] \mathbb{E}_{s,a}^\theta [\psi(s')^\top]$
β^p	$\stackrel{\text{def}}{=}$	$\sup_{\theta, s, a} \lambda_{\max}(\mathbb{C}_{s,a}^\theta [\psi(s')])$ linked to the maximum eigenvalue of $\nabla^2 Z^p$
α^p	$\stackrel{\text{def}}{=}$	$\inf_{\theta, s, a} \lambda_{\max}(\mathbb{C}_{s,a}^\theta [\psi(s')])$ linked to the minimum eigenvalue of $\nabla^2 Z^p$
β^r	$\stackrel{\text{def}}{=}$	$\lambda_{\max}(BB^\top) \sup_{\theta, s, a} \mathbb{V}\text{ar}_{s,a}^\theta(r)$, linked to the maximum eigenvalue of $\nabla^2 Z^r$
α^r	$\stackrel{\text{def}}{=}$	$\lambda_{\min}(BB^\top) \inf_{\theta, s, a} \mathbb{V}\text{ar}_{s,a}^\theta(r)$, linked to the minimum eigenvalue of $\nabla^2 Z^r$

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:

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

B.1 Estimation error

To show that the estimation error $(\sum_{k=1}^K V_{\hat{\theta}^p, \tilde{\theta}^r, 1}^{\pi_k} - V_{\theta^p, \theta^r, 1}^{\pi_k})$ 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 $1 - 5\delta$*

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

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. \square

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 $1 - 2\delta$*

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

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

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.

I) 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:

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

where $\zeta_{hk} = \mathbb{E}_{s_{hk}, a_{hk}}^{\theta^p} [V_{\hat{\theta}^p, \theta^r, h+1}^\pi - V_{\theta^p, \theta^r, h+1}^\pi] - (V_{\hat{\theta}^p, \theta^r, h+1}^\pi(s_{h+1k}) - V_{\theta^p, \theta^r, h+1}^\pi(s_{h+1k}))$ 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

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

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

Remark 4. 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 $1 - \delta$:

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

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) ,

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

This implies that

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

where $\bar{G}_k^p \equiv \bar{G}_{(k-1)H}^p := G_k + (\alpha^p)^{-1} \eta \mathbb{A}$ 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 $k \in \mathbb{N}$

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

Also, using Lemma 18, we have

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

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

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

Remark 5. 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 $K \in \mathbb{N}$

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

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

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

where $\beta^p(K, \delta) \triangleq \frac{\eta}{2} B_{\mathbb{A}}^2 + \gamma_K^p + \log(1/\delta)$, and $\gamma_K^p \triangleq d \log(1 + \frac{\beta^p}{\eta} B_{\varphi, \mathbb{A}} H K)$. 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. The reward estimation error in Equation (10) can be written explicitly. Indeed, using Lemma 17

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

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

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

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

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

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

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

This further yields

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

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

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

Therefore, for all pairs (s, a) ,

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

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

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

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

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

Remark 6. The bad rounds analysis is one of our most important contributions as it enables us to forgo clipping without consequences. Consequently, this is a novel method to control the reward estimation error that improves on existing work for whom clipping was essential.

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

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

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

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

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

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

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

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

This concludes the proof. \square

Remark 7. If we use Lemma 17 without the martingale difference sequence, it will lead to a linear regret. Indeed, the span of the sum of norms over an episode is of order \sqrt{H} . Using the martingale technique instead allows us to retrieve a telescopic sum controlled using the elliptical lemma, this is essential to obtaining a sub-linear regret bound.

B.2 Learning error

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

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

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

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

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

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 Φ is the normal CDF.

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

B.2.1 Stochastic optimism

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

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

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

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

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

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

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

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

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

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

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

Second term. We have

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

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

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

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

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

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

Third term. We have

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

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

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

Choosing $x_k \geq \left(H \sqrt{\frac{\beta^p \beta^p(k, \delta)}{\alpha^p \alpha^r}} + \frac{\sqrt{\beta^r \beta^r(k, \delta) \min\{1, \frac{\alpha^p}{\alpha^r}\}}}{2\alpha^r} \right)$ 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^* . Let us draw $(\bar{\xi}_k)_{k \in [K]}$ i.i.d copies of $(\xi_k)_{k \in [K]}$. Define the optimism event at episode k :

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

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

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

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

$$\min_{\xi_k} V_{\hat{\theta}^p, \hat{\theta}^r + \xi_k, 1}(s_1^k) \quad (17)$$

$$\|\xi_k\|_{\bar{G}_k^p} \leq \sqrt{x_k d \log(d/\delta)}, \quad \forall t \in [H]$$

Under the concentration of our injected noise, we obtain

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

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

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

Also, using the tower rule,

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

Therefore,

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

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

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

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

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

Then,

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

Also,

$$\begin{aligned} & \left| \mathbb{E}_{\xi_k | \bar{O}_k^c} [V_{\hat{\theta}^p, \hat{\theta}^r + \xi_k, 1}(s_1^k) - \underline{V}_1(s_1^k)] \right| \\ & \leq \frac{\beta^r}{2} \mathbb{E}_{\xi_k | \bar{O}_k^c} [\|\tilde{\xi}_k - \underline{\xi}_k\|_{\bar{G}_k^p}] \mathbb{E}_{(\tilde{s}_h) \sim \pi | \hat{\theta}^p} \left[\sum_{h=1}^H \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^p)^{-1}} \right] \\ & \leq \frac{\beta^r}{2} \mathbb{E}_{\tilde{\xi}_k} [\|\tilde{\xi}_k - \underline{\xi}_k\|_{\bar{G}_k^p}] \mathbb{E}_{(\tilde{s}_h) \sim \pi | \hat{\theta}^p} \left[\sum_{h=1}^H \|(A_i \varphi(\tilde{s}_h, \pi(\tilde{s}_h)))_{1 \leq i \leq d}\|_{(\bar{G}_k^p)^{-1}} \right]. \end{aligned}$$

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

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

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

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

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

C Concentrations

C.1 Concentration of the transition parameter

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

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) is \mathcal{F}_{t-1} measurable, and (iii) given (s_t, a_t) , $s_{t+1} \sim P_{\theta^p}^p(\cdot \mid s_t, a_t)$ according to the exponential family defined by Equation (1). Let $\hat{\theta}^p(k)$ be the penalized MLE defined by Equation (5), and let $Z_{s,a}^p(\theta)$ be strictly convex in θ for all (s, a) . Then, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$, the following holds uniformly over all $n \in \mathbb{N}$:

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

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

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

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

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

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

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

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

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

Proof. The result follows from the following simple calculations:

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

□

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 (7), and let $Z_{s,a}^r(\theta)$ be strictly convex in θ for all (s, a) . Then, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$, the following holds uniformly over all $k \in \mathbb{N}$:

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

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

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

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

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

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.

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

Now, for any $\lambda \in \mathbb{R}^d$, we introduce the function $B_{Z_{n,\alpha}, \theta^r}(\lambda) = B_{Z_{n,\alpha}}(\theta^r + \lambda, \lambda)$ 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 $(S_n)_i = \sum_{t=1}^n (r(s_t, a_t) - \mathbb{E}_{s_t, a_t}^{\theta^r}[r]) B^\top A_i \varphi(s_t, a_t)$. Note that $M_n^\lambda > 0$ and it is \mathcal{F}_{n-} measurable. Furthermore, we have for all (s, a) ,

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

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

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

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

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

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

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

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

Furthermore, we have $\sum_{t=1}^n B_{Z_{s_t, a_t}^r, \theta^r}(\lambda) = B_{Z_n^r, \theta^r}(\lambda)$. From the penalized likelihood formula (7), recall that

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

This yields

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

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

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

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

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

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

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

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

Further, any optimal λ must satisfy

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

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

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

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

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

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

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

Therefore, we have from (5) that

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

An application of Markov's inequality now yields

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

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

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

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

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

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

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

which concludes the proof of the first part.

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

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

Now (25) implies that

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

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

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

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 $k \in \mathbb{N}$*

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

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

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 $1 - \delta$*

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

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 Φ 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 $f : \mathcal{X} \rightarrow \mathbb{R}$, we define its span as $\mathbb{S}(f) := \max_{x \in \mathcal{X}} f(x) - \min_{x \in \mathcal{X}} f(x)$. For a probability distribution P supported on the set \mathcal{X} , let $\mathbb{E}_P[f] := \mathbb{E}_P[f(X)]$ and $\mathbb{V}_P[f] := \mathbb{V}_P[f(X)] = \mathbb{E}_P[f(X)^2] - \mathbb{E}_P[f(X)]^2$ 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*

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

D.2 Bregman divergence

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

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

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

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

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

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

To see this, we observe that

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

Now an optimal λ 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

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

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:*

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

And

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

Proof. We prove the lemma as follows

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

□

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

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

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

Similarly,

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

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

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

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

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

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

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

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

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

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

Similarly

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

Proof. We start by writing:

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

then

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

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

The proof of the form of the KL divergence for the reward follows the same steps as above. \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:*

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

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

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

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

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

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

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

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

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

\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: $G_{s,a} := (\varphi(s, a)^\top A_i^\top A_j \varphi(s, a))_{1 \leq i, j \leq d}$, $\bar{G}_n^e \equiv \bar{G}_{(k-1)H}^e := G_n + (\alpha^e)^{-1} \eta A$, 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 e represents either r or p , we omit the superscript 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*

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

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

Further, we have

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

Proof. First we have

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

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

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

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

where the last line follows from the matrix determinant lemma:

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

Therefore:

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

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

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

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 $\|\bar{G}_n^{-1} G_{s_h, a_h}\| \geq 1$ is smaller than $\frac{3d}{\log(2)} \log\left(1 + \frac{\alpha(\|A\|_2 B_{\varphi, \mathbb{A}})^2}{\eta \log(2)}\right)$. Let us call

$\mathcal{T}_k := \{h \in [H] \mid \|\bar{G}_{(k-1)H}^{-1} G_{s_h, a_h}\| \leq 1\}$, then

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

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

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

and

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

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 $a_1, \dots, a_n \in \mathbb{R}^{d \times p}$ be a sequence of matrices with $\|a_t\|_2 \leq L$ for all $t \in [n]$. Let $V_t = V_0 + \sum_{s=1}^t a_s a_s^\top$, then

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

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

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

where the third line follows from the matrix determinant lemma:

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

Rearranging and taking the logarithm shows that

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

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

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

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

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

□

E Tractable planning with random Fourier transform

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

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

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 $\mathbb{P}(s' | s, a, \theta)$ using N sample pairs of (w, b) guarantees an error smaller than ϵ 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 $\text{diam}(\mathcal{M})$. Then, using the explicit mapping \mathbf{z} defined in Algorithm 1 in [RR07] with N samples, we have*

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

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

Further, it implies that if $N = \Omega \left(\frac{p}{\epsilon^2} \log \frac{\sigma_p \text{diam}(\mathcal{M})}{\epsilon} \right)$, 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 ϵ -approximation of $\mathbb{P}(s' | s, a, \theta)$.

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:

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

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

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

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

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

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

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 ϵ of order $1/(H\sqrt{K})$, we deduce that approximating the probability kernel with $\mathcal{O}(pH^2 K)$ samples induces a tractable planning procedure without harming the regret.

Remark 9. *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 incurs a linear regret.*

F Tractable Maximum Likelihood estimation

The maximum likelihood estimation is explicit for simple distributions like the Gaussian [RY77] and for Linearly controlled dynamical systems. But it requires integral approximations for generic transitions. However, we believe that this estimation problem is far simpler than the planning problem since the latter traditionally involves approximating an integral for all state-action pairs.

Different approximation techniques have been used in literature to handle the penalized ML estimation. For instance, *Integral Approximation* techniques are well studied for this problem. Indeed, [Nea01] proposes to handle the ML estimation using simulated annealing, a method that starts from a tractable distribution and updates it to resemble the distribution at hand. [VGB12] proposes *MCMC techniques* for approximating the partition function. [CPH05] shows that optimizing a different objective, called the *contrastive divergence* leads to a good approximation of the ML. Another line of work is related to *Score matching*, a technique that avoids approximating the partition function and is well studied in literature, see [Jør83]. More recently, [LLS⁺21] proposed an adaptation of this technique to the exact setting we consider. The latter shows that under certain conditions, that we are unable to verify, the estimation can be solved in $\mathcal{O}(d^3)$ time. Furthermore, in the case of *Bounded distribution support and natural parameter*, [SSW21] shows that for a minimally represented k -parameter Exponential family, an α -approximation of the ML can be derived in $\mathcal{O}(\text{poly}(k/\alpha))$ time. The latter assumes a specific definition of compactness of the representation as well as knowledge of the support and shows how to re-parameterize the density to a specific class of exponential families that are easier to study. Finally, [DDG⁺19] studies *exponential families such that the natural parameter belongs to an RKHS*, it proposes a method that improves over score matching in time and in memory complexity.