

Near-continuous time Reinforcement Learning for continuous state-action spaces

Lorenzo Croissant

CROISSANT@CEREMADE.DAUPHINE.FR
CEREMADE, CNRS, Université Paris-Dauphine, Université PSL, and Criteo AI Lab, Paris, France

Marc Abeille

M.ABEILLE@CRITEO.COM
Criteo AI Lab, Paris, France

Bruno Bouchard

BOUCHARD@CEREMADE.DAUPHINE.FR
CEREMADE, CNRS, Université Paris-Dauphine, Université PSL, Paris, France

Editors: Claire Vernade and Daniel Hsu

Abstract

We consider the reinforcement learning problem of controlling an unknown dynamical system to maximise the long-term average reward along a single trajectory. Most of the literature considers system interactions that occur in discrete time and discrete state-action spaces. Although this standpoint is suitable for games, it is often inadequate for systems in which interactions occur at a high frequency, if not in continuous time, or those whose state spaces are large if not inherently continuous. Perhaps the only exception is the linear quadratic framework for which results exist both in discrete and continuous time. However, its ability to handle continuous states comes with the drawback of a rigid dynamic and reward structure. This work aims to overcome these shortcomings by modelling interaction times with a Poisson clock of frequency ε^{-1} which captures arbitrary time scales from discrete ($\varepsilon = 1$) to continuous time ($\varepsilon \downarrow 0$). In addition, we consider a generic reward function and model the state dynamics according to a jump process with an arbitrary transition kernel on \mathbb{R}^d . We show that the celebrated optimism protocol applies when the sub-tasks (learning and planning) can be performed effectively. We tackle learning by extending the eluder dimension framework and propose an approximate planning method based on a diffusive limit ($\varepsilon \downarrow 0$) approximation of the jump process. Overall, our algorithm enjoys a regret of order $\tilde{O}(\sqrt{T})$ or $\tilde{O}(\varepsilon^{1/2}T + \sqrt{T})$ with the approximate planning. As the frequency of interactions blows up, the approximation error $\varepsilon^{1/2}T$ vanishes, showing that $\tilde{O}(\sqrt{T})$ is attainable in near-continuous time.

Keywords: Online Reinforcement Learning, Stochastic Control, Continuous State-Space, Diffusion Approximation, Optimism in the Face of Uncertainty, Eluder Dimension

1. Introduction

Controlling a dynamical system to drive it to optimal long-term average behaviour is a key challenge in many applications, ranging from mechanical engineering to econometrics. Reinforcement Learning (RL) aims to do so when the system is a priori unknown by tackling jointly both the control and the statistical inference of the system. This joint objective is even more important in the online version of the problem, in which one interacts with the system along a single trajectory (no resets or episodes). In the last decades, the insights of Bandit Theory (see e.g. [Lattimore and Szepesvári \(2020\)](#)) have been leveraged to tackle the RL problem, while addressing the inherent exploration-exploitation dilemma that naturally arises in sequential decision-making (see e.g. [Szepesvári \(2010, § 4.2\)](#)).

Continuous time. While the RL literature has so far focused on discrete-time settings, many real-world systems involve interactions which occur at discrete times with a continuous-time state process, often with a very high frequency of interactions. This could be due, e.g., to the limitations of a sensor or actuator in a continuous environment or because event times rely on an exogenous noise process. A natural approach to planning in such systems is to directly model the problem in continuous time, as is common in finance (Chen and Yao, 2001; Cont and Tankov, 2004; Obizhaeva and Wang, 2013). While the resulting continuous-time stochastic control problems are well studied, they appear to conflict with the sample-based nature of statistical learning theory that fundamentally takes place in discrete time.

Near-continuous time. In order to subsume both perspectives on the problem, we consider interactions governed by a Poisson clock, setting the expected inter-arrival time of the clock to a parameter $\varepsilon \in (0, 1)$. This gives us control over a continuum of situations from discrete time ($\varepsilon = 1$) to continuous time ($\varepsilon \downarrow 0$). We call this embedding of discrete time into \mathbb{R}_+ *near-continuous time* because it allows us to consider the regime in which $\varepsilon \ll 1$ (which is of interest for modelling high-frequency systems). Nonetheless, even for $\varepsilon \gg 0$, it enriches the usual discrete-time analysis with new perspectives and mathematical tools.

Modelling dynamics. An essential prerequisite for modelling real-world systems is the ability to capture complex (non-linear) dynamics and rich reward signals defined over continuous state and action spaces. With this in mind, we focus on the model-based approach where the transition and the reward function belong to a parameterised class of functions. Achieving this level of generality poses challenges regarding all three key sub-tasks of RL, which are: planning, learning, and the exploration-exploitation trade-off.

Planning in continuous systems. For (discrete-time) dynamics on finite state-action spaces, the planning problem falls under the umbrella of Markov Decision Processes (MDPs) which have been extensively reviewed in Puterman (2005). The finite nature of MDPs is at the heart of their theoretical and computational success. Their extension to countable or even continuous state spaces is, however, non-trivial; see e.g. Bertsekas (2011, § 4.6, p.245) for a review of the challenges. Perhaps the only exception which retains those nice theoretical and computational properties is the celebrated Linear Quadratic (LQ) framework (Kalman, 1960). However, both frameworks are limited in their expressive power. In contrast, continuous-time stochastic control theory has demonstrated how to effectively solve the control problem for arbitrary regular dynamics on continuous state spaces. It enjoys a rich and mature literature (Arisawa and Lions, 1998; Arapostathis et al., 2012; Lions, 1983), both on the theoretical aspects as well as numerical solvers based on Partial Differential Equations (PDEs), another storied field (Kushner and Dupuis, 2001; Barles and Souganidis, 1991; Bonnans and Zidani, 2003). The near-continuous time framework lies between the two theories, and recent results of Abeille et al. (2022) show how to navigate between them and approximately solve the planning problem in the high-frequency interactions regime by solving its diffusive counterpart.

Learning non-linear systems. Similar to the planning problem, the natural way to move beyond finite Markov chain models and towards continuous state dynamics is through linear models. The least-squares estimator enjoys strong theoretical guarantees including adaptive confidence sets that can be efficiently maintained online (see e.g. Abbasi-yadkori

et al. (2011)). Extensions (Russo and Van Roy, 2013; Osband and Van Roy, 2014) showed how to extend this approach to richer model classes through the use of Non-Linear Least Squares (NLLS). This framework subsumes standard least squares and has been successful in many dynamics by retaining its key properties regarding confidence sets. While providing a protocol for learning with NLLS, Russo and Van Roy (2013) characterised the trade-off between the richness of the model and the hardness of its learning through two quantities of the model class: the log-covering number, and the eluder dimension which summarises the difficulty of turning the information from data into predictive power.

Optimistic exploration. Optimism in the Face of Uncertainty (OFU) has proven highly successful in sequential decision-making from bandits to RL. The works of Jaksch et al. (2010); Auer and Ortner (2006); Bartlett and Tewari (2009) showed how to extend the celebrated UCB (Auer et al., 2002) algorithm from bandits to finite MDPs; later, extensions were made to continuous states in the LQ setting, see e.g. Abbasi-Yadkori and Szepesvári (2011); Abeille and Lazaric (2020); Cohen et al. (2019) and references therein. Extension from bandit to MDP and then to LQ raised new challenges that persist in our setting. First, the agent should not revise its behaviour too often to prevent dithering, which requires the design of a lazy update scheme. Second, generic continuous states-spaces models come with inherent unboundedness, and one must carefully address stability issues.

In this work, we consider the near-continuous time system interaction model and propose an optimistic algorithm for online reinforcement learning in the average reward setting¹. Our approach builds on the work of Abeille et al. (2022) to introduce continuous-time tools for studying the planning problem and on extending the work of Russo and Van Roy (2013) to our near-continuous time and unbounded state setting to perform the learning with NLLS. Underlying the extension of both these two approaches is a careful treatment of the state boundedness which we perform with Lyapunov stability arguments. Our algorithm enjoys the $\tilde{O}(\sqrt{T})$ complexity of the discrete-state case, including technical generalisations of standard learning complexity metrics. Further, leveraging another result by Abeille et al. (2022), we demonstrate an efficient approximate planning method in the regime $\varepsilon \downarrow 0$ which yields a regret scaling with $\tilde{O}(\varepsilon^{1/2}T + \sqrt{T})$. In the asymptotic ($\varepsilon \downarrow 0$), the approximation error vanishes showing that $\tilde{O}(\sqrt{T})$ is attainable even in high-frequency settings.

2. Setting

We consider an agent interacting with its environment to maximise a long-term average reward. At each interaction, it observes the current state of the system $x \in \mathbb{R}^d$, takes action $a \in \mathbb{A} \subset \mathbb{R}^{d_{\mathbb{A}}}$, and receives reward $r(x, a)$, for $r : \mathbb{R}^d \times \mathbb{A} \rightarrow \mathbb{R}$. The system then transitions to a state denoted by, say, x' according to

$$x' = x + \mu_{\theta^*}(x, a) + \Sigma\xi \quad \text{with} \quad \xi \sim \mathcal{N}(0, \mathbf{I}_d),$$

$\Sigma \in \mathbb{R}^{d \times d}$, and in which $\mu_{\theta^*} : \mathbb{R}^d \times \mathbb{A} \rightarrow \mathbb{R}^d$ is the deterministic motion of the system². Contrasting with the standard setting, we consider here the interactions to occur in a

-
1. Also known as, *average cost per stage*, *long-run average*, or *ergodic* setting.
 2. While the additive noise structure is a design choice that simplifies the analysis, the choice of parameterising the drift as $x + \mu_{\theta^*}(x, a)$ instead of $\mu_{\theta^*}(x, a)$ does not affect its generality and is made only for convenience. See Assumption 1 below.

random fashion, which we model by an independent Poisson process of intensity ε^{-1} . As such, ε parameterises the mean wait time between events and gives us direct control over the frequency of interactions.

State dynamics. Let $\Omega := \mathbb{D}$ be the space of càdlàg functions from $[0, +\infty)$ to \mathbb{R}^d , and let \mathbb{P} be a probability measure on Ω . We formalise the interaction time and the noise process as a marked \mathbb{P} -compound Poisson process $(N_t)_{t \in \mathbb{R}_+}$ of intensity $\varepsilon^{-1} \geq 1$. We denote by $(\tau_n)_{n \in \mathbb{N}}$ its arrival (interaction) times, with $\tau_0 := 0$, and by $(\xi_n)_{n \in \mathbb{N}}$ its marks, which are independent of everything else and drawn i.i.d. according to the centred standard Gaussian measure ν on \mathbb{R}^d . We encode the information available at time $t \in \mathbb{R}_+$ in the σ -algebra $\mathcal{F}_t := \sigma((\tau_n, \xi_n)_{\tau_n \leq t})$ and with the filtration \mathbb{F} defined as the completion of $(\mathcal{F}_t)_{t \in \mathbb{R}_+}$. Let \mathcal{A} be the set of \mathbb{F} -adapted \mathbb{A} -valued processes, referred to as *controls*. For any initial state $x_0 \in \mathbb{R}^d$ and $\alpha \in \mathcal{A}$, we let X^{α, θ^*} denote the pathwise-unique solution of

$$\begin{cases} X_{\tau_n}^{\alpha, \theta^*} = X_{\tau_{n-1}}^{\alpha, \theta^*} + \mu_{\theta^*}(X_{\tau_{n-1}}^{\alpha, \theta^*}, \alpha_{\tau_{n-1}}) + \Sigma \xi_n \\ X_{\tau_0}^{\alpha, \theta^*} = x_0 \end{cases} . \quad (1)$$

In (1), we model the dynamic according to a jump process and X^{α, θ^*} is then defined at any time $t \in \mathbb{R}_+$ by considering that it is piece-wise constant on each interval $[\tau_{n-1}, \tau_n)$, $n \in \mathbb{N}^*$. Although involved, this definition allows us to define the state process at any time and feature the interplay of the Poisson ($N_t \in \mathbb{N}$) and wall-time ($t \in \mathbb{R}_+$) clocks.

Reinforcement learning problem. In our model-based paradigm, ignorance about the system is condensed to a single parameter set $\Theta \subset \mathbb{R}^{d_\Theta}$ containing the unknown nominal parameter θ^* . To single out the RL challenges, we further assume that θ^* only affects the drift assuming other quantities (i.e. Σ , ε , and r) are known to the agent. For any $x_0 \in \mathbb{R}^d$, we evaluate the performance of any strategy $\alpha \in \mathcal{A}$ with the long-term average reward criterion³ defined by

$$\rho_{\theta^*}^\alpha(x_0) := \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\sum_{n=1}^{N_T} r(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) \right] . \quad (2)$$

The goal of the agent is to accumulate as much reward as possible, i.e. to compete with the best an omniscient agent can achieve: $\rho_{\theta^*}^\alpha(x_0) := \sup_{\alpha \in \mathcal{A}} \rho_{\theta^*}^\alpha(x_0)$. We evaluate the quality of a learning algorithm generating α according to its regret.

Definition 1 For any $T \in \mathbb{R}_+$, $x_0 \in \mathbb{R}^d$, and $\alpha \in \mathcal{A}$, the regret of α is

$$\mathcal{R}_T(\alpha) := T \rho_{\theta^*}^\alpha(x_0) - \sum_{n=1}^{N_T} r(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) . \quad (3)$$

Noticing that N_T is the number of events up to time T , the definitions of the optimal performance (2) and the regret (3) again feature the interplay between the wall-clock (T) and Poisson clock (N_T): the agent's realised trajectory uses the Poisson clock, which governs interactions, while the ideal performance is understood per unit of wall-clock time.

3. Notice that this criterion is strategically equivalent to a discrete-time control problem with the same transition dynamics as X^{α, θ^*} . Thus, all of the following results also apply in the discrete-time case.

2.1. Working Assumptions

Of particular interest in our approach is the high-frequency regime in which $\varepsilon \downarrow 0$. In this framework, many interactions occur per unit of time, each of which is of negligible impact both in terms of dynamics and reward. This regime can be encoded by introducing, for any parameter $\theta \in \Theta$, rescaled coefficients $(\bar{\mu}_\theta, \bar{\Sigma}, \bar{r})$ connected to the original parametrisation by

$$\mu_\theta = \varepsilon \bar{\mu}_\theta, \quad \Sigma = \varepsilon^{\frac{1}{2}} \bar{\Sigma}, \quad \text{and } r = \varepsilon \bar{r}.$$

In this rescaled parametrisation, $\bar{\mu}_\theta$, $\bar{\Sigma}$, and \bar{r} are understood as independent of ε . To improve legibility, we will make alternating use of both representations (μ_θ, Σ, r) and $(\bar{\mu}_\theta, \bar{\Sigma}, \bar{r})$. While the scaling of μ_θ and r in ε arises naturally, the one of Σ is a design choice: we consider the covariance $\Sigma \Sigma^\top$ to be linear in ε . Known as the diffusive regime, this preserves stochasticity⁴ as $\varepsilon \downarrow 0$.

We now impose regularity assumptions on the drift and reward signal, uniformly over the possible parametrisations and controls $(\alpha, \theta) \in \mathcal{A} \times \Theta$. We take $\|\cdot\|$ to be the Euclidian norm on \mathbb{R}^d and $\|\cdot\|_{\text{op}}$ for the operator norm on $\mathbb{R}^{d \times d}$ associated to $\|\cdot\|$.

Assumption 1 *The map $(\bar{\mu}, \bar{r})$ is continuous, and there is $L_0 > 0$ such that for all $(\theta, a) \in \Theta \times \mathbb{A}$*

$$L_0 > \sup_{x \in \mathbb{R}^d} \frac{\|\bar{\mu}_\theta(x, a)\|}{1 + \|x\|} + \sup_{x \neq x'} \frac{\|\bar{\mu}_\theta(x, a) - \bar{\mu}_\theta(x', a)\|}{\|x - x'\|} + \sup_{x \in \mathbb{R}^d} \|\bar{r}(x, a)\| + \sup_{x \neq x'} \frac{\|\bar{r}(x, a) - \bar{r}(x', a)\|}{\|x - x'\|}.$$

Furthermore, $L_0 > \|\bar{\Sigma}\|_{\text{op}}$ and $\bar{\Sigma} \bar{\Sigma}^\top \succeq \varsigma \text{I}_d$ for some $\varsigma > 0$, where \succeq denotes the Loewner order.

Assumption 1 mainly imposes regularity on both $\bar{\mu}_\theta$ and \bar{r} through a Lipschitz condition. We also assume rewards to be bounded, which may be relaxed, but doing so is highly technical and involves trading off the growth of r with the stability of the process (see Assumption 2). Note that we do not assume boundedness of $\bar{\mu}_\theta$. Finally, we assume non-degeneracy of the noise by requiring $\bar{\Sigma}$ to be full rank.

We conclude with Assumption 2 to ensure the stability of the state process. Let $\mathbb{R}_*^d := \mathbb{R}^d \setminus \{0\}$ and $\mathbb{R}_+ := (0, +\infty)$. For $k \in \mathbb{N}$, let $\mathcal{C}^k(\mathbb{R}_*^d; \mathbb{R}_+)$ denote the set of k -times continuously differentiable functions from \mathbb{R}_*^d to \mathbb{R}_+ . Let ∇ and ∇^2 denote the gradient and Hessian operator respectively.

Assumption 2 *There is $(\ell_\mathcal{V}, L_\mathcal{V}, \mathbf{c}_\mathcal{V}, M_\mathcal{V}, M'_\mathcal{V}) \in \mathbb{R}_+^5$ and a Lyapunov function $\mathcal{V} \in \mathcal{C}^2(\mathbb{R}_*^d; \mathbb{R}_+)$ satisfying, for any $(x, x', a, \theta) \in \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{A} \times \Theta$, $x \neq x'$, and $\varepsilon \in (0, 1)$:*

- (i) $\ell_\mathcal{V} \|x - x'\| \leq \mathcal{V}(x - x') \leq L_\mathcal{V} \|x - x'\|$,
- (ii) $\sup_{x \in \mathbb{R}_*^d} \|\nabla \mathcal{V}(x)\| \leq M_\mathcal{V}$ and $\sup_{x \in \mathbb{R}_*^d} \|\nabla^2 \mathcal{V}(x)\|_{\text{op}} \leq M'_\mathcal{V}$,
- (iii) $\mathcal{V}(x + \varepsilon \bar{\mu}(x, a) - x' - \varepsilon \bar{\mu}(x', a)) \leq (1 - \varepsilon \mathbf{c}_\mathcal{V}) \mathcal{V}(x - x')$. (4)

4. Another common, but more rigid, regime is to consider $\Sigma = \varepsilon \bar{\Sigma}$, whose limit regime is deterministic and known as the fluid limit, see [Fernandez-Tapia et al. \(2016\)](#).

Assumption 2 is a Lyapunov-like condition through the function \mathcal{V} . The condition (i.) requires that \mathcal{V} behaves similarly to a norm, while (ii.) asks that \mathcal{V} be smoothly differentiable everywhere but at 0 and (iii.) imposes a contraction condition on the drifts.

Connection to linear stability. Stability theory has been extensively studied in the special case of linear dynamics. In this case, we recover Assumption 2 from the Continuous Algebraic Riccati Equation (CARE; see e.g. Lancaster and Rodman (1995), § 4.4). Considering linear dynamics $\bar{\mu}_\theta(x, a) = \bar{A}x + \bar{B}a$ (given matrices (\bar{A}, \bar{B}) of appropriate dimensions), continuous stability is guaranteed when the eigenvalues of \bar{A} have negative real-part or, equivalently, by the existence of a positive semi-definite matrix P solving the CARE $\bar{A}^\top P + P\bar{A} = -I_d$. For this P , its associated norm $\mathcal{V} = \|\cdot\|_P$ is the appropriate Lyapunov function for Assumption 2. Indeed, conditions (i.) and (ii.) follow as \mathcal{V} is a norm and, for $\varepsilon \leq 1/2\lambda_{\max}(P)$, we have

$$\begin{aligned} \mathcal{V}(x + \varepsilon\bar{\mu}(x, a) - x' - \varepsilon\bar{\mu}(x', a))^2 &= (x - x')^\top (P + \varepsilon\bar{A}^\top P + \varepsilon P\bar{A} + \varepsilon^2 P)(x - x') \\ &= (x - x')^\top (P - \varepsilon I_d + \varepsilon^2 P)(x - x') \\ &\leq (x - x')^\top (P - \varepsilon P/\lambda_{\max}(P) + \varepsilon^2 P)(x - x') \\ &\leq (1 - \varepsilon/2\lambda_{\max}(P))\mathcal{V}(x - x')^2. \end{aligned}$$

Taking the square-root and using $\sqrt{1 - \varepsilon/2\lambda_{\max}(P)} \leq 1 - \varepsilon/4\lambda_{\max}(P)$ leads to (iii.) with $c_\mathcal{V} = 1/4\lambda_{\max}(P)$.

3. Main results

Our main contribution is a demonstration of the OFU protocol in the near-continuous time continuous state-action RL problem. The ingredients of OFU are: learning from accumulated data to design confidence sets; lazy updates to trade off policy revision and learning guarantees; and planning amongst plausible parameterisations. We summarise this protocol in Algorithm 1.

Algorithm 1 OFU- \mathbb{R}^d

Input: confidence level δ , initial state x_0 , initial control ϖ_0
for $n \in \mathbb{N}^*$ **do**
 At time τ_n , receive $r(X_{\tau_{n-1}}^{\varpi, \theta^*}, \varpi_{\tau_{n-1}})$ and $X_{\tau_n}^{\varpi, \theta^*}$.
 if n satisfies (7) **then**
 $n_k \leftarrow n, k \leftarrow k + 1,$
 Compute $\hat{\theta}_{n_k}$ using (5) and $\mathcal{C}_{n_k}(\delta/3)$ with (6).
 $\tilde{\theta}_k \leftarrow \operatorname{argmax}_{\theta \in \mathcal{C}_{n_k}(\delta/3)} \rho_\theta^*$
 $\pi_k \leftarrow \pi_{\tilde{\theta}_k}^*$ using (8)
 end if
 Play $\varpi_{\tau_n} := \pi_k(X_{\tau_n}^{\varpi, \theta^*})$.
end for

Learning. Our algorithm proceeds by episodes, indexed by $k \in \mathbb{N}$ with n_k denoting the start of the k^{th} episode. At each n_k , Algorithm 1 revises its knowledge using the Non-Linear Least-Square fit and the associated confidence set $\mathcal{C}_{n_k}(\delta)$, defined (for $\beta_n(\delta)$ given

and discussed in (12), for all $n \in \mathbb{N}$) by

$$\hat{\theta}_{n_k} \in \operatorname{argmin}_{\theta \in \Theta} \sum_{n=0}^{n_k-1} \left\| X_{\tau_{n+1}}^{\varpi, \theta^*} - X_{\tau_n}^{\varpi, \theta^*} - \mu_{\theta}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) \right\|^2, \quad (5)$$

$$\mathcal{C}_{n_k}(\delta) := \left\{ \theta \in \Theta : \sqrt{\sum_{n=0}^{n_k-1} \left\| \mu_{\theta}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) - \mu_{\hat{\theta}_{n_k}}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) \right\|^2} \leq \beta_{n_k}(\delta) \right\}. \quad (6)$$

Lazy Updates. Our episodic scheme follows the same rationale as in Jaksch et al. (2010); Abbasi-Yadkori and Szepesvári (2011), and triggers updates as soon as enough information is collected. Formally, it constructs a sequence of episodes $\{S_k\}_{k \in \mathbb{N}}$ whose starting times are defined by $n_0 := 0$ and, for any $k \in \mathbb{N}$, n_{k+1} is the first time $n > n_k$ satisfying (7)

$$\sqrt{\sup_{\theta \in \mathcal{C}_{n_k}(\delta)} \sum_{i=0}^n \left\| \mu_{\theta}(X_{\tau_i}^{\varpi, \theta^*}, \varpi_{\tau_i}) - \mu_{\hat{\theta}_{n_k}}(X_{\tau_i}^{\varpi, \theta^*}, \varpi_{\tau_i}) \right\|^2} > 2\beta_n(\delta). \quad (7)$$

Planning. Algorithm 1 requires us to be able to plan using any $\theta \in \mathcal{C}_{n_k}(\delta) \subset \Theta$, and as such we will extend the definitions of $X^{\alpha, \theta}$, $\rho_{\theta}^{\alpha}(x_0)$, $\rho_{\theta}^*(x_0)$ to any $(\alpha, \theta) \in \mathcal{A} \times \Theta$ by replacing θ^* by θ in (1) and (2). An optimal Markov control for ρ_{θ}^* can be obtained by solving an integral (Hamilton-Jacobi-Bellman) equation of the form

$$\varepsilon \rho_{\theta}^*(x) = \max_{a \in \mathbb{A}} \{ \mathbb{E}[W_{\theta}^*(x + \mu_{\theta}(x, a) + \Sigma\xi)] - W_{\theta}^*(x) + r(x, a) \} \quad \forall x \in \mathbb{R}^d; \quad (8)$$

Let \mathcal{A} be the set of measurable maps from \mathbb{R}^d to \mathbb{A} . Any map $\pi_{\theta}^* \in \mathcal{A}$ such that $\pi_{\theta}^*(x)$ is an argument of the maximum in (8) for all $x \in \mathbb{R}^d$ is an optimal policy. Algorithm 1 thus obtains via (8) an optimal control for ρ_{θ}^* which we denote by $\pi_{\theta}^* := \pi_{\theta}^* \circ X^{\pi_{\theta}^*, \theta^*}$ ⁵.

3.1. Regret Bound

Stability. Working with unbounded processes and generic drift requires us to prevent state blow-up, which could degrade regret regardless of learning. In Proposition 2 we combine the Lyapunov stability of (4) with concentration arguments to show that unstable trajectories can only happen with low probability. A detailed proof is given in Appendix B.

Proposition 2 *Under Assumptions 1 and 2, there is a function $H_{\delta}(n) = \mathcal{O}(\sqrt{\log(n\delta^{-1})})$ such that for any $\delta \in (0, 1)$, $\alpha \in \mathcal{A}$, $x_0 \in \mathbb{R}^d$, and $\theta \in \Theta$ we have*

$$\mathbb{P} \left(\sup_{t \in \mathbb{R}_+} \frac{\|X_t^{\alpha, \theta}\|}{H_{\delta}(N_t)} \geq 1 \right) \leq \delta. \quad (9)$$

Learning. In order for the regret analysis to be meaningful, the learning complexity metrics of $\mathcal{F}_{\Theta} := (\mu_{\theta})_{\theta \in \Theta}$ (covering number and eluder dimension) must be adapted for

5. We make this notational confusion between the policy π_{θ}^* and the control process it generates in order to write $\rho_{\theta}^{\pi}(x)$ and $X^{\pi, \theta}$ instead of $\rho_{\theta}^{\pi \circ X^{\pi, \theta}}$ and $X^{\pi \circ X^{\pi, \theta}, \theta}$. when $\pi \in \mathcal{A}$

unbounded functions on unbounded processes. Indeed, attempting to Ξ -cover the class of linear functions from $\mathbb{R}^d \rightarrow \mathbb{R}^d$ would make the regret bound vacuous for every $\Xi > 0$. Proposition 2 allows us to restrict ourselves to the set of states effectively traversed by the state-process X^{α, θ^*} . While the quantities have the same intuitions, there is an intricate technical contribution in this extension whose details we defer to Appendix C.

For $R > 0$, let $\mathcal{B}_2(R) \subset \mathbb{R}^d$ denotes the Euclidean ball of radius R at 0. For any set $S \subset \mathbb{R}^d$, the S -effective covering number of \mathcal{F}_Θ is the covering number of $\mathcal{F}_\Theta|_S := \{f|_{S \times \mathbb{A}} : f \in \mathcal{F}_\Theta\}$. By Proposition 2, we can work with H_δ by formally defining $\mathcal{N}_n^\varepsilon$ as the size of the smallest cover $\mathcal{C}_n^\varepsilon$ of \mathcal{F}_Θ such that

$$\sup_{\mu_1 \in \mathcal{F}_\Theta} \min_{\mu_2 \in \mathcal{C}_n^\varepsilon} \sup_{x \in \mathcal{B}_2(H_\delta(n))} \|\mu_1(x) - \mu_2(x)\| \leq \frac{\varepsilon \|\bar{\Sigma}\|_{\text{op}}^2}{n}. \quad (10)$$

The Ξ -eluder dimension (for $\Xi \in \mathbb{R}_+$) of a function class \mathcal{F}_Θ , introduced in Russo and Van Roy (2013) and denoted $\text{dim}_E(\mathcal{F}_\Theta, \Xi)$, is a notion of dimension which is perfectly tailored to converting fit errors into prediction errors. We defer to Russo and Van Roy (2013) for its technical definition. For any set $S \subset \mathbb{R}^d$, the S -effective eluder dimension is the eluder dimension of $\mathcal{F}_\Theta|_S$, which we denote by $\text{dim}_E^S(\mathcal{F}_\Theta, \Xi)$. For $n \in \mathbb{N}^*$, let $B_n := \mathcal{B}_2(\sup_{i \in [n]} \|X_{\tau_i}\|)$ and let us define the sequence of effective eluder dimensions $(d_{E,n})_{n \in \mathbb{N}^*}$ by

$$d_{E,n} := \text{dim}_E^{B_n} \left(\mathcal{F}_\Theta, \frac{2\varepsilon^{\frac{1}{2}}}{\sqrt{n}} \right)$$

for all $n \in \mathbb{N}^*$ and $u \in \mathbb{R}_+$.

Remark 3 *While eluder dimension is perfectly tailored to the needs of regret bounds of optimistic algorithms, it remains a somewhat abstract measure. For clarity, we reproduce some known bounds for unmodified eluder dimension from Osband and Van Roy (2014).*

- (i.) *If $\mathcal{F}_\Theta = \{f|f(x) = \theta\phi(x)\}$, for a kernel $\phi : \mathbb{R}^d \rightarrow B_\phi$ (in which B_ϕ is a ball of radius $k_\phi > 0$ in \mathbb{R}^{d_ϕ}) and $\theta \in \mathbb{R}^{d \times d_\phi}$, then $\text{dim}_E(\mathcal{F}_\Theta, \Xi) \leq \tilde{O}(dd_\phi \log(1 + k_\phi k_\Theta \Xi^{-1}))$ in which $k_\Theta := \sup_{\theta \in \Theta} \|\theta\|_{\text{op}}$.*
- (ii.) *Let $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be a component-wise independent function with $\sup_{i \in [d]} \|\partial_i g\| \leq L_g < +\infty$ and $\inf_{i \in [d]} \|\partial_i g\| \geq \ell_g > 0$ on \mathbb{R}^d and let $\kappa_g := L_g/\ell_g$. If $\mathcal{F}_\Theta = \{f|f(x) = g(\theta\phi(x))\}$ for (ϕ, θ) as above, then $\text{dim}_E(\mathcal{F}_\Theta, \Xi) \leq \tilde{O}(dd_\phi \kappa_g^2 \log(1 + k_\phi k_\Theta \Xi^{-1}))$*

Our extension makes these bounds more applicable to unbounded processes, for instance in (i.) by allowing kernels which do not map \mathbb{R}^d to a bounded set, so long as they map each ball B_n to a ball B_ϕ . One example are linear dynamics ($\phi = \text{Id}$) which incur a complexity of $\tilde{O}(d^2 \log(1 + k_\Theta \sup_{i \in [n]} \|X_{\tau_i}^{\alpha, \theta^}\| \Xi^{-1}))$.*

Theorem 4 *Under Assumptions 1 and 2, for any $\delta \in (0, 1)$, $x_0 \in \mathbb{R}^d$, there is a constant $C \in \mathbb{R}_+$ independent of ε such that Algorithm 1 achieves*

$$R_T(\varpi) \leq C \sqrt{d_{E, \lceil T\varepsilon^{-1} \rceil} \log(\mathcal{N}_{\lceil T\varepsilon^{-1} \rceil}^\varepsilon) T \log(T\delta^{-1})} \quad (11)$$

with probability at least $1 - \delta$, in which $d_{E, \lceil T\varepsilon^{-1} \rceil}$ is the effective $2\varepsilon/\sqrt{T}$ -eluder dimension and $\log(\mathcal{N}_{\lceil T\varepsilon^{-1} \rceil}^\varepsilon)$ is the effective $\varepsilon^2 \|\bar{\Sigma}\|_{\text{op}}^2/T$ -log-covering number.

Theorem 4 exhibits the scaling in the complexity measures expected from Russo and Van Roy (2013), both eluder dimension and log-covering numbers, as well as the usual $\sqrt{T \log(T\delta^{-1})}$ dependency from UCB.

4. Ideas of the Proof

Working on the high-probability event of Proposition 2 allows us to handle the unbounded state in learning, planning, and optimism.

4.1. Learning

Confidence Sets. The crux of our analysis is incorporating Proposition 2 into the NLLS method of Russo and Van Roy (2013). Restricting the domain of \mathcal{F}_Θ allows us to handle the richness of unbounded models and states while following Russo and Van Roy (2013) to define confidence sets. Let $\delta \in (0, 1)$, set $\beta_0 := \varepsilon^{\frac{1}{2}}$, and let

$$\beta_n(\delta) := \beta_0 \vee 2\varepsilon^{\frac{1}{2}} \|\bar{\Sigma}\|_{\text{op}} \left(\sqrt{1 + 2 \left(\sqrt{2 \log \left(\frac{4\pi^2 n^3}{3\delta} \right) + \sqrt{2\kappa_n(\delta)}} \right) + \sqrt{\kappa_n(\delta)}} \right) \quad (12)$$

in which

$$\kappa_n(\delta) := \log \left(\frac{2\pi^2 n^2 \varepsilon \mathcal{N}_n^\varepsilon}{3\delta} (\|\bar{\Sigma}\|_{\text{op}}^2 + 8nL_0^2(1 + H_\delta(n))) \right).$$

Using this choice $(\beta_n)_{n \in \mathbb{N}}$ and replacing n_k by n in (6) formally defines the confidence sets $(\mathcal{C}_n(\delta))_{n \in \mathbb{N}}$. For any $\alpha \in \mathcal{A}$, the probability that the state process X_t^{α, θ^*} outgrows $H_\delta(N_t)$ is small and, thus, this confidence set will hold with high probability as shown by Proposition 5.

Proposition 5 (Adapted from Osband and Van Roy (2014, Prop. 5)) *Under Assumptions 1 and 2, for any $x_0 \in \mathbb{R}^d$, and $\delta > 0$,*

$$\mathbb{P} \left(\left\{ \theta^* \in \bigcap_{n=1}^{\infty} \mathcal{C}_n(\delta) \right\} \cap \left\{ \sup_{n \in \mathbb{N}^*} \frac{\|X_{\tau_n}^{\varpi, \theta^*}\|}{H_\delta(n)} \leq 1 \right\} \right) \geq 1 - \delta, \quad (13)$$

Prediction error. A well-posed confidence set is not sufficient for low-regret approaches in the OFU paradigm. This high confidence (low fit error) of the NLLS estimator must be translated as a low online prediction error. In Proposition 6 we obtain first- and second-order prediction error bounds from the effective eluder dimension. In Proposition 6 the order notation \tilde{O} hides terms that are poly-logarithmic in N_t and d_{E, N_t} whose the full details are given in Appendix C.2.

Proposition 6 *Under Assumptions 1 and 2, for any $\delta \in (0, 1)$, $\alpha \in \mathcal{A}$, $x_0 \in \mathbb{R}^d$, and $t \in \mathbb{R}_+$, we have with probability at least $1 - \delta$*

$$\sum_{n=1}^{N_t} \left\| \mu_{\hat{\theta}_n} (X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) - \mu_{\theta^*} (X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) \right\| \leq \tilde{O} \left(\sqrt{\varepsilon d_{E, N_t} \log(\cdot \mathcal{N}_{N_t}^\varepsilon)} N_t + \varepsilon d_{E, N_t} \right), \quad (14)$$

and

$$\sum_{n=1}^{N_t} \left\| \mu_{\hat{\theta}_n} (X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) - \mu_{\theta^*} (X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) \right\|^2 \leq \tilde{O} (d_{E, N_t} \log(\cdot \mathcal{N}_{N_t}^\varepsilon)). \quad (15)$$

Lazy updates. We leverage the second order bound (15) of Proposition 6 to define our lazy-update scheme (7). We show in Appendix E that this scheme does not degrade the speed at which Algorithm 1 learns by more than a constant factor, while also ensuring that the policy is only updated logarithmically in the number of interactions up to any horizon.

4.2. Planning

For a given $\theta \in \Theta$, the well-posedness of the control problem $\rho_\theta^*(x_0)$ and its resolution are non-trivial.

Proposition 7 (Adapted from (Abeille et al., 2022, Thm. 2.3 & Rem. 2.4.))

Under Assumptions 1 and 2, there is $L_W \in \mathbb{R}_+$, independent of ε , such that for any $\theta \in \Theta$

- (i.) *The map $x \mapsto \rho_\theta^*(x)$ is constant, taking only one value which we denote by $\rho_\theta^* \in \mathbb{R}$;*
- (ii.) *There is an L_W -Lipschitz function W_θ^* such that*

$$\varepsilon \rho_\theta^* = \max_{a \in \mathbb{A}} \{ \mathbb{E}[W_\theta^*(x + \mu_\theta(x, a) + \Sigma \xi)] - W_\theta^*(x) + r(x, a) \} \quad \forall x \in \mathbb{R}^d; \quad (16)$$

- (iii.) *There is $\pi_\theta^* \in \mathcal{A}$, such that for all $x \in \mathbb{R}^d$, $\pi_\theta^*(x)$ maximises the right hand side in (16), and $\pi_\theta^* \circ X^{\pi_\theta^*, \theta}$ is an optimal Markov control, i.e. $\rho_\theta^{\pi_\theta^*}(\cdot) \equiv \rho_\theta^*$.*

Proposition 7.(i.) shows that the control problem ρ_θ^* is independent of the initial conditions and meaningfully ergodic, which follows from the stability analysis of the process using (4). Points (ii.) and (iii.) show that there is an optimal policy, which can be computed by solving the HJB equation (16). Unfortunately (16) is an integral equation with low regularity, owing to the non-local jumps of the system, which complicates its analysis and the construction of numerical solvers.

4.3. Regret Decomposition

To sketch the proof of Theorem 4, we work on the high-probability event of Proposition 5, and omit martingale measurability issues this could cause. We will also ignore the randomness of jump times and consider $T \lesssim \varepsilon N_T$, with \lesssim denoting inequality up to a constant. Appendix E is dedicated to a complete proof. Recall that ϖ denotes the control generated by Algorithm 1.

Proof Let $k : \mathbb{N} \rightarrow \mathbb{N}$ map an event n to the episode $k(n)$ to which it belongs and let $\theta_n := \tilde{\theta}_{k(n)}$. We begin the regret decomposition by applying the HJB equation (16) to the rewards collected along the trajectory $r(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n})$ in the definition of the regret. Conditioning as appropriate, this yields

$$\mathcal{R}_T(\varpi) = T\rho_{\theta^*}^* - \varepsilon \sum_{n=1}^{N_T} \rho_{\theta_n}^{\pi_{\theta_n}^*}(0) \quad (R_1)$$

$$+ \sum_{n=1}^{N_T} \mathbb{E}[W_{\theta_n}^*(\tilde{X}_{\tau_{n+1}}^{\varpi, \theta_n}) | \mathcal{F}_{\tau_n}] - W_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*}) \quad (R_2)$$

in which $\tilde{X}_{\tau_{n+1}}^{\varpi, \theta} := X_{\tau_n}^{\varpi, \theta^*} + \mu_{\theta}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) + \Sigma \xi_{n+1}$, for $(n, \theta) \in \mathbb{N} \times \Theta$, is a counterfactual one-step transition assuming parameter $\theta \in \Theta$.

On the event of Proposition 5, θ^* is in $\cap_{n \in \mathbb{N}} \mathcal{C}_n(\delta)$ and the optimism of Algorithm 1 ensures that (R1) is negative. For (R2), the identity

$$\tilde{X}_{\tau_{n+1}}^{\varpi, \theta} = \tilde{X}_{\tau_{n+1}}^{\varpi, \theta^*} - \mu_{\theta^*}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) + \mu_{\theta}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n})$$

combined with the L_W -Lipschitzness of W_{θ}^* from Proposition 8, yields

$$R_2 \leq L_W \sum_{n=1}^{N_T} \left\| \mu_{\theta_n}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) - \mu_{\theta^*}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) \right\| \quad (R_4)$$

$$+ \sum_{n=1}^{N_T} \mathbb{E}[W_{\theta_n}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) - W_{\theta_{n+1}}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}] \quad (R_5)$$

$$+ \sum_{n=1}^{N_T} \mathbb{E}[W_{\theta_{n+1}}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}] - W_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*}), \quad (R_6)$$

by adding and subtracting $\mathbb{E}[W_{\theta_{n+1}}^*(\tilde{X}_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}] = \mathbb{E}[W_{\theta_{n+1}}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}]$. (R6) is a martingale term, which we can bound using concentration theory. Our lazy update-scheme ensures that $\theta_n \neq \theta_{n+1}$ only $\mathcal{O}(\log(N_T))$ times by time T , keeping (R5) small.

It remains to show that the lazy update-scheme, does not degrade the learning of (R4), which is controlled by improvements to Proposition 6 in Appendix C which yield

$$\sum_{n=1}^{N_T} \sup_{(\theta_1, \theta_2) \in \mathcal{C}_{k(n)}(\delta)^2} \left\| \mu_{\theta_1}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) - \mu_{\theta_2}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) \right\| \lesssim \tilde{\mathcal{O}} \left(\sqrt{d_{E, T\varepsilon^{-1}} \log(\mathcal{N}_{T\varepsilon^{-1}}^{\varepsilon}) T} \right).$$

■

5. Approximate Planning for $\varepsilon \downarrow 0$

Embedding the control problem in continuous time allowed us to extend learning complexity measures and construct an optimistic algorithm for unbounded continuous states in

Sections 3 and 4. The natural next step is to develop efficient sub-routines for each sub-task in a modular manner. Past work such as UCRL2 (Jaksch et al., 2010) gives some indications of how to do this. In this section, we present another way to improve computational efficiency which relies on the power of the near-continuous-time formulation to apply results from Abeille et al. (2022) to the planning problem.

As $\varepsilon \downarrow 0$, ρ_ε^* will enter a diffusive limit regime whose limit control problem is

$$\bar{\rho}_\theta^*(x_0) := \liminf_{T \rightarrow +\infty} \frac{1}{T} \mathbb{E} \left[\int_0^T \bar{r}(\bar{X}_t^{\bar{\alpha}, \theta}, \bar{\alpha}_t) dt \right] \quad \text{in which} \quad \begin{cases} d\bar{X}_t^{\bar{\alpha}, \theta} = \bar{\mu}_\theta(\bar{X}_t^{\bar{\alpha}, \theta}, \bar{\alpha}_t) dt + \bar{\Sigma} dW_t \\ \bar{X}_0^{\bar{\alpha}, \theta} = x_0 \end{cases} \quad (17)$$

with W denoting a \mathbb{P} -Brownian motion, $\bar{\mathbb{F}}$ its filtration, and $\bar{\mathcal{A}}$ the set \mathbb{A} -valued $\bar{\mathbb{F}}$ -predictable processes. This control problem has been extensively studied, see e.g. (Arisawa and Lions, 1998; Arapostathis et al., 2012), and Proposition 8 shows it to be well-posed.

Proposition 8 (Adapted from Abeille et al. (2022, Thm. 3.4.))

Under Assumptions 1 and 2, for any $\theta \in \Theta$,

(i.) The map $x \mapsto \bar{\rho}_\theta^*(x)$ is constant, taking only one value which we denote by $\bar{\rho}_\theta^* \in \mathbb{R}$.

(ii.) There is an L_W -Lipschitz function $\bar{W}_\theta^* \in \mathcal{C}^2(\mathbb{R}^d; \mathbb{R})$ such that

$$\bar{\rho}_\theta^* = \max_{a \in \mathbb{A}} \left\{ \bar{\mu}_\theta(x, a)^\top \nabla \bar{W}_\theta^*(x) + \bar{r}(x, a) \right\} + \frac{1}{2} \text{Tr}[\bar{\Sigma} \bar{\Sigma}^\top \nabla^2 \bar{W}_\theta^*(x)], \quad \forall x \in \mathbb{R}^d. \quad (18)$$

(iii.) There is $\bar{\pi}_\theta^* \in \mathcal{A}$ such that, for all $x \in \mathbb{R}^d$, $\bar{\pi}_\theta^*(x)$ maximises the right hand side in (18), and $\bar{\pi}_\theta^* \circ \bar{X}^{\bar{\pi}_\theta^*, \theta}$ is an optimal Markov control, i.e. $\bar{\rho}_\theta^*(\cdot) \equiv \bar{\rho}_\theta^*$.

The limit HJB equation (18) is purely differential, and, thus, local: the solution at x depends only on its derivatives at x . This is fundamentally simpler than the non-local behaviour of (16), in which there are cross-dependencies between points due to the expectation. Moreover, this diffusive PDE belongs to a well-studied family, both from the points of view of theory (Gilbarg and Trudinger, 1983; Ladyzhenskaya and Ural'tseva, 1968) and of numerics (Knabner and Angermann, 2003; Kushner, 1977).

These facts strongly motivate the use of these tools to construct approximate planning methods for (16) in the near-continuous time regime as $\varepsilon \downarrow 0$. It is important to note that this approximation is not a numerical approximation but an approximation of state process $X^{\alpha, \theta}$ by another state process $\bar{X}^{\alpha, \theta}$. This is only possible because of strong tools from the theory of viscosity solutions of PDEs available by embedding into continuous time.

Proposition 9 (Adapted from Abeille et al. (2022, Thm. 3.6.))

Under Assumptions 1 and 2, for any $\gamma \in (0, 1)$, there is a constant $C_\gamma > 0$, independent of ε , such that, for any $\theta \in \Theta$,

$$|\bar{\rho}_\theta^* - \rho_\theta^*| \leq C_\gamma \varepsilon^{\frac{\gamma}{2}} \quad \text{and} \quad \rho_\theta^* - \rho_\theta^{\bar{\pi}_\theta^*}(0) \leq C_\gamma \varepsilon^{\frac{\gamma}{2}}. \quad (19)$$

Moreover, there is a function $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ such that,

$$\varepsilon \rho_\theta^{\bar{\pi}_\theta^*}(0) = \mathbb{E}[\bar{W}_\theta^*(x + \mu_\theta(x, a) + \Sigma \xi)] - \bar{W}_\theta^*(x) + r(x, \bar{\pi}_\theta^*(x)) + e_\theta(x), \quad \forall x \in \mathbb{R}^d \quad (20)$$

and there is $C'_\gamma > 0$, independent of ε , such that $|e_\theta(x)| \leq C'_\gamma \varepsilon^{1 + \frac{\gamma}{2}} (1 + \|x\|^3)$ for all $x \in \mathbb{R}^d$.

Proposition 9, combined with (18) provides a certifiable approximation for solving the control problem (2) with off-the-shelf diffusive HJB solvers, at a cost independent of ε . An example of this methodology is seen in (Abeille et al., 2022, § 4), in which (Abeille et al., 2022, Fig. 1, p. 30) shows the reduction in computational effort. Using this method for approximate planning on top of Algorithm 1 yields Algorithm 2.

Algorithm 2 OFU-Diffusion

Input: confidence level δ , initial state x_0 , initial control ϖ'_0
for $n \in \mathbb{N}^*$ **do**
 At time τ_n , receive $r(X_{\tau_{n-1}}^{\varpi', \theta^*}, \varpi'_{\tau_{n-1}})$ and $X_{\tau_n}^{\varpi', \theta^*}$.
if n satisfies (7) **then**
 $n_k \leftarrow n, k \leftarrow k + 1,$
 Compute $\hat{\theta}_{n_k}$ using (5) and $\mathcal{C}_{n_k}(\delta/3)$ with (6).
 $\tilde{\theta}_k \leftarrow \operatorname{argmax}_{\theta \in \mathcal{C}_{n_k}(\delta/3)} \bar{\rho}_{\theta}^*$
 $\pi'_k \leftarrow \bar{\pi}_{\tilde{\theta}_k}^*$ using (18)
end if
 Play $\varpi'_{\tau_n} := \pi'_k(X_{\tau_n}^{\varpi', \theta^*})$.
end for

Proposition 9 also provides in (20) an HJB-like representation of the approximation, which provides a key with which to analyse the regret incurred when using this approximation. As seen in the sketch of proof of Theorem 10.

Theorem 10 *Under Assumptions 1 and 2, for any $\delta \in (0, 1)$, $x_0 \in \mathbb{R}^d$, and $\gamma \in (0, 1)$, there is a pair $(C_\gamma, C) \in \mathbb{R}_+^2$ of constants independent of ε such that Algorithm 2 achieves*

$$R_T(\varpi') \leq 2C_\gamma \varepsilon^{\frac{\gamma}{2}} T + C \sqrt{d_{E, \lceil T\varepsilon^{-1} \rceil} \log \left(\mathcal{N}_{\lceil T\varepsilon^{-1} \rceil}^\varepsilon \right) T \log(T\delta^{-1})} \quad (21)$$

with probability at least $1 - \delta$.

Compared to Theorem 4, Theorem 10 has an additional linear term from the approximate planning method which scales with $C_\gamma \varepsilon^{\gamma/2}$. The dependency of the constant in γ is inherited from the analysis of Abeille et al. (2022) and $C_\gamma < +\infty$ holds for $\gamma < 1$. Quantifying the behaviour of C_γ as $\gamma \uparrow 1$ is technically intricate. Nevertheless, our bound indicates that the long run approximation error vanishes as $\varepsilon \downarrow 0$ almost as fast as $\sqrt{\varepsilon}$, making it practical for systems with very high jump intensity.

Proof The proof follows the same lines as the proof of Theorem 4 and we only sketch the appropriate modifications. Instead of (16), apply the HJB-like equation (20) of Proposi-

tion 9.(iii.) which yields

$$\mathcal{R}_T(\varpi') = T\rho_{\theta^*}^* - \varepsilon \sum_{n=1}^{N_T} \rho_{\theta_n}^{\bar{\pi}_{\theta_n}^*}(0) \tag{R1}$$

$$+ \sum_{n=1}^{N_T} \mathbb{E}[\bar{W}_{\theta_n}^*(\tilde{X}_{\tau_{n+1}}^{\varpi', \theta_n}) | \mathcal{F}_{\tau_n}] - \bar{W}_{\theta_n}^*(X_{\tau_n}^{\varpi', \theta^*}) \tag{R2}$$

$$+ \sum_{n=1}^{N_T} e_{\theta_n}(X_{\tau_n}^{\varpi', \theta^*}) \tag{R3}$$

in which $\tilde{X}_{\tau_{n+1}}^{\varpi', \theta^*}$ is defined analogously to $\tilde{X}_{\tau_{n+1}}^{\varpi, \theta^*}$. Noticing that \bar{W}_{θ}^* is L_W -Lipschitz for any $\theta \in \Theta$, just as W_{θ}^* is, (R2) can be treated with the same arguments as in the proof of Theorem 4. On the event of Proposition 5, θ^* is in $\cap_{n \in \mathbb{N}} \mathcal{C}_n(\delta)$ and the optimism of Algorithm 2 ensures that $\bar{\rho}_{\theta^*}^* \leq \bar{\rho}_{\theta_n}^* = \bar{\rho}_{\theta_n}^{\bar{\pi}_{\theta_n}^*}$ for all $n \in \mathbb{N}$. Combining this with Proposition 9, show that (R1) decomposes into

$$R_1 \lesssim \varepsilon \left(\sum_{n=1}^{N_T} (\rho_{\theta^*}^* - \bar{\rho}_{\theta^*}^*) + \sum_{n=1}^{N_T} (\bar{\rho}_{\theta_n}^* - \rho_{\theta_n}^{\bar{\pi}_{\theta_n}^*}) \right) \leq 4N_T C_{\gamma} \varepsilon^{1+\frac{\gamma}{2}}.$$

Meanwhile, Proposition 9 implies that $R_3 \leq \varepsilon^{1+\frac{\gamma}{2}} N_T (1+H_{\delta})(N_T)^3$. Thus, $R_1+R_3 \lesssim C_{\gamma} \varepsilon^{\frac{\gamma}{2}} T$. ■

6. Conclusion

In this work, we proposed a general framework for the Reinforcement Learning problem of controlling an unknown dynamical system, on a continuous state-action space, to maximise the long-term average reward along a single trajectory. In particular, we focused on the understudied high-frequency systems driven by many small movements. Modelling such systems as controlled jump processes, we provided an optimistic algorithm which leverages Non-Linear Least Squares for learning and the diffusive limit regime for approximate planning. This proof of concept calls for several further refinements to be implementable in practice.

Optimism. The optimistic step of Algorithm 1 chooses $\tilde{\theta}_n$ in an inefficient manner. Like in UCRL2 (Jaksch et al., 2010), optimistic exploration can be performed at the same time as planning by solving an expanded HJB equation, i.e. (18) with the maximum now taken over $(a, \theta) \in \mathbb{A} \times \Theta$. Since our assumptions are uniform in $\theta \in \Theta$, this is possible up to a modified regret decomposition, as in Jaksch et al. (2010).

Lazy updates. The way we quantify learning progress to design the lazy update scheme (7) remains fundamentally discrete. Computationally cheaper lazy update schemes might be obtained through simpler heuristics. For instance, the scaling of the drift with ε suggests it could be possible to update periodically, directly in terms of the wall-clock time T .

Case-by-case. As a proof of concept, we endeavoured to study the RL problem in high generality. However, practical applications must use all available model information

to refine the method ad hoc. This is true for the learning method (replacing NLLS with an estimator specialised to the model at hand and bound the eluder dimension and log-covering numbers) and for numerical schemes on the PDE (18) which are built on a case-by-case basis for $d > 1$, see [Kushner and Dupuis \(2001\)](#).

References

- Yasin Abbasi-Yadkori and Csaba Szepesvári. Regret bounds for the adaptive control of linear quadratic systems. In Sham M. Kakade and Ulrike von Luxburg, editors, *Proceedings of the 24th Annual Conference on Learning Theory*, volume 19 of *Proceedings of Machine Learning Research*, pages 1–26, Budapest, Hungary, 09–11 Jun 2011. PMLR. URL <https://proceedings.mlr.press/v19/abbasi-yadkori11a.html>.
- Yasin Abbasi-yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011. URL https://proceedings.neurips.cc/paper_files/paper/2011/file/e1d5be1c7f2f456670de3d53c7b54f4a-Paper.pdf.
- Marc Abeille and Alessandro Lazaric. Efficient optimistic exploration in linear-quadratic regulators via Lagrangian relaxation. 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 23–31. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/abeille20a.html>.
- Marc Abeille, Bruno Bouchard, and Lorenzo Croissant. Diffusive limit approximation of pure jump optimal ergodic control problems, September 2022. URL <http://arxiv.org/abs/2209.15284>. arXiv:2209.15284 [math].
- Ari Arapostathis, Vivek S. Borkar, and Mrinal K. Ghosh. *Ergodic control of diffusion processes*. Number 143 in *Encyclopedia of Mathematics and its Applications*. Cambridge University Press, Cambridge, 2012. ISBN 9781139003605.
- Mariko Arisawa and Pierre-Louis Lions. On ergodic stochastic control. *Communications in partial differential equations*, 23(11-12):2187–2217, 1998.
- Peter Auer and Ronald Ortner. Logarithmic Online Regret Bounds for Undiscounted Reinforcement Learning. In *Advances in Neural Information Processing Systems*, volume 19. MIT Press, 2006. URL <https://proceedings.neurips.cc/paper/2006/hash/c1b70d965ca504aa751ddb62ad69c63f-Abstract.html>.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47:235–256, 2002.
- Guy Barles and Panagiotis E. Souganidis. Convergence of approximation schemes for fully nonlinear second order equations. *Asymptotic analysis*, 4(3):271–283, 1991.
- Peter L. Bartlett and Ambuj Tewari. Regal: A regularization based algorithm for reinforcement learning in weakly communicating MDPs. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, UAI '09, page 3542, Arlington, VA, 2009. AUAI Press. ISBN 9780974903958.
- Dimitri P Bertsekas. *Dynamic Programming and Optimal Control*, volume II. Athena Scientific, Belmont, MA, 3rd edition, 2011. ISBN 1886529442.

- J. Frédéric Bonnans and Housnaa Zidani. Consistency of Generalized Finite Difference Schemes for the Stochastic HJB Equation. *SIAM Journal on Numerical Analysis*, 41(3):1008–1021, January 2003. URL <https://epubs.siam.org/doi/abs/10.1137/S0036142901387336>. Publisher: Society for Industrial and Applied Mathematics.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: a nonasymptotic theory of independence*. Oxford University Press, Oxford, 1st edition, 2013. ISBN 9780199535255.
- V. V. Buldygin and I. V. Kozachenko. *Metric characterization of random variables and random processes*, volume 188 of *Translations of mathematical monographs*. American Mathematical Society, Providence, RI, 2000. ISBN 9780821805336.
- Hong Chen and David D Yao. *Fundamentals of queueing networks: Performance, asymptotics, and optimization*, volume 46 of *Stochastic Modelling and Applied Probability*. Springer, New York, NY, 2001. ISBN 9781441928962.
- Alon Cohen, Tomer Koren, and Yishay Mansour. Learning linear-quadratic regulators efficiently with only \sqrt{T} regret. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1300–1309. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/cohen19b.html>.
- Rama Cont and Peter Tankov. *Financial modelling with jump processes*. Chapman & Hall/CRC financial mathematics series. Chapman & Hall/CRC, Boca Raton, FL, 2004. ISBN 978-1-58488-413-2.
- Joaquin Fernandez-Tapia, Olivier Guéant, and Jean-Michel Lasry. Optimal Real-Time Bidding Strategies. *Applied Mathematics Research eXpress*, September 2016. ISSN 1687-1200, 1687-1197. URL <https://academic.oup.com/amrx/article-lookup/doi/10.1093/amrx/abw007>.
- D. Gilbarg and N. S. Trudinger. *Elliptic partial differential equations of second order*, volume 224 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, 2nd edition, 1983. ISBN 354013025-X.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for Reinforcement Learning. *Journal of Machine Learning Research*, 11(51):1563–1600, 2010. ISSN 1533-7928.
- Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. *Transactions of the ASME—Journal of Basic Engineering*, 82(Series D):35–45, 1960.
- Rafail Khasminskii. *Stochastic Stability of Differential Equations*, volume 66 of *Stochastic Modelling and Applied Probability*. Springer, Berlin, Heidelberg, 2012. ISBN 9783642232794.
- Peter Knabner and Lutz Angermann. *Numerical methods for elliptic and parabolic partial differential equations*. Number 44 in Texts in applied mathematics. Springer, New York, NY, 2003. ISBN 9780387954493.

- Harold J. Kushner. *Probability Methods for Approximations in Stochastic Control and for Elliptic Equations*, volume 129 of *Mathematics in Science and Engineering*. Elsevier, 1977. ISBN 9780124301405.
- Harold J. Kushner and Paul Dupuis. *Numerical Methods for Stochastic Control Problems in Continuous Time*, volume 24 of *Stochastic Modelling and Applied Probability*. Springer New York, New York, NY, 2001. ISBN 978-1-4612-6531-3 978-1-4613-0007-6.
- Olga A. Ladyzhenskaya and Nina N. Ural'tseva. *Linear and quasilinear elliptic equations*, volume 46 of *Mathematics in Science and Engineering*. Elsevier, 1968. ISBN 9780124328501.
- Peter Lancaster and Leiba Rodman. *Algebraic riccati equations*. Clarendon press, 1995. ISBN 9780198537953.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020. ISBN 9781108571401.
- Michel Ledoux and Michel Talagrand. *Probability in Banach Spaces*, volume 23 of *Classics in Mathematics*. Springer, Berlin, Heidelberg, 1991. ISBN 9783642202117.
- Pierre-Louis Lions. Optimal control of diffusion processes and Hamilton-Jacobi-Bellman equations part 2: viscosity solutions and uniqueness. *Communications in partial differential equations*, 8(11):1229–1276, 1983.
- Anna A. Obizhaeva and Jiang Wang. Optimal trading strategy and supply/demand dynamics. *Journal of Financial Markets*, 16(1):1–32, 2013.
- Ian Osband and Benjamin Van Roy. Model-based reinforcement learning and the eluder dimension. In *Proceedings of the 27th International Conference on Neural Information Processing Systems-Volume 1*, pages 1466–1474, 2014. URL <https://proceedings.neurips.cc/paper/2014/hash/1141938ba2c2b13f5505d7c424ebae5f-Abstract.html>.
- Martin L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. Wiley series in probability and statistics. Wiley-Interscience, Hoboken, NJ, 2005. ISBN 9780471727828.
- Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *Proceedings of the 26th International Conference on Neural Information Processing Systems-Volume 2*, pages 2256–2264, 2013. URL <https://proceedings.neurips.cc/paper/2013/hash/41bfd20a38bb1b0bec75acf0845530a7-Abstract.html>.
- Csaba Szepesvári. *Algorithms for reinforcement learning*. Number 4 in Synthesis lectures on artificial intelligence and machine learning. Morgan & Claypool Publishers, 2010.

Appendices

Appendix A. Preliminaries

A.1. Organisation of Appendices

We prove the results one by one, starting with stability, then learning, planning, and finally concluding with the regret proof of Theorem 4.

In Appendix B, we go over the probabilistic properties of our problem and show several bounds on the stability of the process, in the sense of high probability and moment boundedness. In particular, the main objective of this appendix is to prove Proposition 2.

In Appendix C, we show a generalisation of the existing theory of learning with NLLS to the case of unbounded functions on unbounded domains. The key results are Propositions 5 and 6

In Appendix D, we provide a characterisation of the control part of the RL problem we analyse, including the diffusion limit approximation, namely Propositions 7 to 9.

In Appendix E, we perform regret analysis and collect the last few results used to prove the regret bound of Theorem 4. This includes the treatment of the lazy update scheme.

The remainder of Appendix A is devoted to notations and short-hands used throughout, but each appendix is meant to be as notationally stand-alone as possible.

A.2. General notation

The set of natural numbers including 0 is denoted \mathbb{N} , while $\mathbb{N}^* := \mathbb{N} \setminus \{0\}$ denotes the set of (strictly) positive integers. For $n \in \mathbb{N}^*$, we use $[n]$ to denote the set of positive integers up to and including n , i.e. $[n] := \{1, \dots, n\}$. Let \mathbb{R} denote the set of real numbers and define $\mathbb{R}_+ := (0, +\infty)$ and $\mathbb{R}_*^d := \mathbb{R}^d \setminus \{0\}$. The space of sequences taking values in S will be denoted by $S^{\mathbb{N}}$. For $S \subset \mathbb{R}^d$, we also denote the complement of S by $S^c := \mathbb{R}^d \setminus S$, we use the same notation for the complement of a probability event.

We denote by $\langle \cdot | \cdot \rangle$ the inner product on \mathbb{R}^d , by $\|\cdot\|$ the Euclidean norm on \mathbb{R}^d , and by $\|\cdot\|_{\text{op}}$ the associated operator norm on $\mathbb{R}^{d \times d}$. For $R \in \mathbb{R}_+$ and $x \in \mathbb{R}^d$, we denote the Euclidean ball of radius R centred at x by $\mathcal{B}_2(x, R)$, and when $x = 0$ we use the shorthand $\mathcal{B}_2(R)$ for $\mathcal{B}_2(0, R)$.

For $d \geq 1$, $\mathcal{D} \subset \mathbb{R}^d$ and $\mathcal{D}' \subset \mathbb{R}$, we denote the space of continuous functions from \mathcal{D} to \mathcal{D}' by $\mathcal{C}^0(\mathcal{D}; \mathcal{D}')$. For any $k \in \mathbb{N}^*$, we denote $\mathcal{C}^k(\mathcal{D}; \mathcal{D}')$ the subset of $\mathcal{C}^0(\mathcal{D}; \mathcal{D}')$ containing all functions which are continuously differentiable up to order k .

A.3. Problem dependent notation

The space of càdlàg (rcll) functions from $[0, +\infty)$ to \mathbb{R}^d , for $d \in \mathbb{N}^*$, is denoted \mathbb{D} and \mathbb{P} is a probability measure on $\Omega := \mathbb{D}$. $(N_t)_{t \in \mathbb{R}_+}$ denotes a marked \mathbb{P} -compound Poisson process of intensity $\varepsilon^{-1} > 1$, $(\tau_n)_{n \in \mathbb{N}}$ denotes the sequence of its arrival times, with $\tau_0 := 0$, and $(\xi_n)_{n \in \mathbb{N}}$ denotes the sequence of its marks. Namely, the sequences $(\tau_n)_{n \in \mathbb{N}}$ and $(\xi_n)_{n \in \mathbb{N}}$ are independent, $(\tau_{n+1} - \tau_n)_{n \in \mathbb{N}}$ is i.i.d. with exponential distribution of parameter ε and $(\xi_n)_{n \in \mathbb{N}}$ is i.i.d. with standard Gaussian measure on \mathbb{R}^d , which we denoted by ν .

For $t \in [0, +\infty)$, $\mathcal{F}_t := \sigma((\tau_n, \xi_n)_{\tau_n \leq t})$ and the filtration \mathbb{F} is the completion of $(\mathcal{F}_t)_{t \in \mathbb{R}_+}$. The set of \mathbb{F} -adapted \mathbb{A} -valued processes, which we consider as admissible controls, is denoted \mathcal{A} . For any $(x_0, \alpha, \theta) \in \mathbb{R}^d \times \mathcal{A} \times \Theta$, $X^{\alpha, \theta}$ is the solution of

$$\begin{cases} X_{\tau_n}^{\alpha, \theta} = X_{\tau_{n-1}}^{\alpha, \theta} + \mu_\theta(X_{\tau_{n-1}}^{\alpha, \theta}, \alpha_{\tau_{n-1}}) + \Sigma \xi_n \\ X_{\tau_0}^{\alpha, \theta} = x_0 \end{cases} . \quad (22)$$

When specifying the dependence on the initial condition $x_0 \in \mathbb{R}^d$ is necessary, we write $X^{x_0, \alpha, \theta}$. This process is defined for any $t \in [0, +\infty)$ by considering its trajectories as piece-wise constant on any interval of the form $[\tau_{n-1}, \tau_n)$ for $n \in \mathbb{N}^*$. For any $(x_0, \alpha, \theta) \in \mathbb{R}^d \times \mathcal{A} \times \Theta$, the control problem is denoted by

$$\rho_\theta^*(x_0) := \sup_{\alpha \in \mathcal{A}} \rho_\theta^\alpha(x_0) \text{ in which } \rho_\theta^\alpha(x_0) := \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\sum_{n=1}^{N_T} r(X_{\tau_n}^{x_0, \alpha, \theta}, \alpha_{\tau_n}) \right] .$$

We denote by W a \mathbb{P} -Wiener process (a.k.a Brownian motion), by $\bar{\mathbb{F}}$ the \mathbb{P} -augmentation of the filtration it generates, and by $\bar{\mathcal{A}}$ the collection of \mathbb{A} -valued and $\bar{\mathbb{F}}$ -predictable processes. For any $(x_0, \bar{\alpha}, \theta) \in \mathbb{R}^d \times \bar{\mathcal{A}} \times \Theta$, we denote by $\bar{X}^{\bar{\alpha}, \theta}$ (or $\bar{X}^{x_0, \bar{\alpha}, \theta}$ if specifying the initial condition) the solution of

$$\begin{cases} d\bar{X}_t^{\bar{\alpha}, \theta} = \bar{\mu}_\theta(\bar{X}_t^{\bar{\alpha}, \theta}, \bar{\alpha}_t) dt + \bar{\Sigma} dW_t \\ \bar{X}_0^{\bar{\alpha}, \theta} = x_0 \end{cases} . \quad (23)$$

The associated control problem is denoted by

$$\bar{\rho}_\theta^*(x_0) := \sup_{\bar{\alpha} \in \bar{\mathcal{A}}} \bar{\rho}_\theta^{\bar{\alpha}}(x_0) \text{ in which } \bar{\rho}_\theta^{\bar{\alpha}}(x_0) := \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\int_0^T r(\bar{X}_t^{x_0, \bar{\alpha}, \theta}, \bar{\alpha}_t) dt \right] .$$

According to Propositions 7 and 8, we defined the constants $\rho_\theta^* := \rho_\theta^*(0)$ and $\bar{\rho}_\theta^* := \bar{\rho}_\theta^*(0)$. For $\theta \in \Theta$, $\bar{\pi}_\theta^*$ denotes a policy in $\bar{\mathcal{A}}$ (the set of measurable maps from \mathbb{R}^d to \mathbb{A}) which maximises the right-hand side of the HJB equation (16) associated to $\bar{\rho}_\theta^*$ (see Proposition 8). Throughout, we use the same notation for policies and the Markov controls they induce, provided there is no ambiguity.

We use ϖ to denote the control process output of Algorithm 1 mathematically. For any $\omega \in \Omega$, the trajectory generated by Algorithm 1 is therefore defined as in (22) by $X^{\varpi, \theta^*}(\omega)$. By definition of Algorithm 1, in its k^{th} episode (i.e. for $t \in [\tau_{n_k}, \tau_{n_k+1})$), $\varpi_t = \pi_k(X_t^{\varpi, \theta^*})$, with $\pi_k := \bar{\pi}_{\theta_k}^*$.

Throughout these appendices, we will use the shorthand $\psi_\theta^\varepsilon(x, a) := x + \varepsilon \bar{\mu}_\theta(x, a)$, for any $(x, a, \theta) \in \mathbb{R}^d \times \mathbb{A} \times \Theta$.

Appendix B. State Process Stability

A key aspect of our setting is that both the state process $X^{\alpha,\theta}$, for any $(\alpha, \theta) \in \mathcal{A} \times \Theta$, and the drift μ itself are unbounded. This can lead to an exponential blow-up of the state process, which can be harmful to both the learning and control aspects. In order to avoid this difficulty we imposed Assumption 2, which corresponds to a stochastic Lyapunov condition, and ensures that the state will not explode in expectation. We reinforce this result by leveraging concentration theory to obtain the high-probability bound of Proposition 2. Appendix B.1 is dedicated to its proof, and it will be used in the proofs of learning results and high-probability regret bounds (Appendices C and E).

Proposition 2 *Under Assumptions 1 and 2, there is a function $H_\delta(n) = \mathcal{O}(\sqrt{\log(n\delta^{-1})})$ such that for any $\delta \in (0, 1)$, $\alpha \in \mathcal{A}$, $x_0 \in \mathbb{R}^d$, and $\theta \in \Theta$ we have*

$$\mathbb{P} \left(\sup_{t \in \mathbb{R}_+} \frac{\|X_t^{\alpha,\theta}\|}{H_\delta(N_t)} \geq 1 \right) \leq \delta. \quad (9)$$

Unlike learning and regret, the analysis of the control task is done in expectation via the HJB equation. Here the unbounded drift will materialise as higher moments of $X^{\alpha,\theta}$. The counterpart of Proposition 2 in this case is a moment result, given by Lemma 15, which is proved in Appendix B.2 and will then be used in Appendix D.

Lemma 15 *Under Assumptions 1 and 2, for any $p \geq 2$, there is a constant $\mathbf{c}'_p > 0$ independent of ε such that*

$$\mathbb{E} \left[\|X_t^{x_0, \alpha, \theta}\|^p \right] \leq \frac{1}{\ell_\gamma^p} \left(L_\gamma^p e^{-\frac{\mathbf{c}'_\gamma}{4}t} \|x_0\|^p + \frac{4\mathbf{c}'_p}{\mathbf{c}'_\gamma} \left(1 - e^{-\frac{\mathbf{c}'_\gamma}{4}t} \right) \right),$$

for any $(x_0, \alpha, \theta) \in \mathbb{R}^d \times \mathcal{A} \times \Theta$ and $t \in [0, +\infty)$.

B.1. Proof of Proposition 2

This appendix is dedicated to the proof of Proposition 2 which is a high probability bound on the state process. This proof follows the Chernoff method. Thus, we will derive an exponential moment bound for the state process in Lemma 12. We will first obtain a stochastic stability condition in expectation in Lemma 11. In what follows, let $R_\varepsilon := \sqrt{8d \log(1/\varepsilon)}$ and $\xi \sim \nu$.

Lemma 11 *Under Assumptions 1 and 2,*

(i.) *for any $(\eta, x, a, \theta) \in \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{A} \times \Theta$, we have*

$$\mathcal{V}(\psi_\theta^\varepsilon(x, a) - \sqrt{\varepsilon}\eta) \leq (1 - \varepsilon\mathbf{c}_\gamma)\mathcal{V}(x - \sqrt{\varepsilon}\eta) + \varepsilon M_\gamma L_0(1 + \|\eta\|); \quad (24)$$

(ii.) *and, for any $(a, \theta) \in \mathbb{A} \times \Theta$, and any $x \notin \mathcal{B}_2(\varepsilon^{\frac{1}{2}} \|\bar{\Sigma}\|_{\text{op}} R_\varepsilon)$ we have*

$$\mathbb{E}[\mathcal{V}(\psi_\theta^\varepsilon(x, a) + \Sigma\xi)] \leq (1 - \varepsilon\mathbf{c}_\gamma)\mathcal{V}(x) + \varepsilon\mathbf{c}'_\gamma$$

in which \mathbf{c}'_γ is a constant independent of ε .

Proof

(i.) By Lipschitzness of \mathcal{V} and (4), for any $(\eta, x, a, \theta) \in \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{A} \times \Theta$, we have

$$\begin{aligned} \mathcal{V}(\psi_\theta^\varepsilon(x, a) - \sqrt{\varepsilon}\eta) &= \mathcal{V}(\psi_\theta^\varepsilon(x, a) - \psi_\theta^\varepsilon(\sqrt{\varepsilon}\eta, a) + \varepsilon\bar{\mu}(\sqrt{\varepsilon}\eta, a)) \\ &\leq \mathcal{V}(\psi_\theta^\varepsilon(x, a) - \psi_\theta^\varepsilon(\sqrt{\varepsilon}\eta, a)) + M_{\mathcal{V}}\varepsilon \|\bar{\mu}(\sqrt{\varepsilon}\eta, a)\| \\ &\leq (1 - \varepsilon\mathbf{c}_{\mathcal{V}})\mathcal{V}(x - \sqrt{\varepsilon}\eta) + M_{\mathcal{V}}\varepsilon \|\bar{\mu}(\sqrt{\varepsilon}\eta, a)\|, \end{aligned}$$

from which (24) follows by using Assumption 1, which implies $\|\bar{\mu}(\sqrt{\varepsilon}\eta, a)\| \leq L_0(1 + \sqrt{\varepsilon}\|\eta\|) \leq L_0(1 + \|\eta\|)$ since $\varepsilon \in (0, 1)$.

(ii.) For any $x \in \mathbb{R}^d$, by the symmetry of the law of $\bar{\Sigma}\xi$, by (24) applied for $\eta = \bar{\Sigma}\xi$, and by taking the expectation, we have

$$\begin{aligned} \mathbb{E}[\mathcal{V}(\psi_\theta^\varepsilon(x, a) + \Sigma\xi)] &= \mathbb{E}[\mathcal{V}(\psi_\theta^\varepsilon(x, a) - \sqrt{\varepsilon}\bar{\Sigma}\xi)] \\ &\leq (1 - \varepsilon\mathbf{c}_{\mathcal{V}})\mathbb{E}[\mathcal{V}(x - \sqrt{\varepsilon}\bar{\Sigma}\xi)] + \varepsilon M_{\mathcal{V}}L_0(1 + \|\bar{\Sigma}\|_{\text{op}}\mathbb{E}[\|\xi\|]). \end{aligned} \quad (25)$$

Since ξ is a standard Gaussian, $\|\xi\|^2$ is a random variable following a χ^2 distribution with d degrees of freedom, thus $\mathbb{E}[\|\xi\|^2] = d$, and by Jensen's inequality $\mathbb{E}[\|\xi\|] \leq \sqrt{d}$. Thus the second term is bounded by $\varepsilon M_{\mathcal{V}}L_0(1 + \|\bar{\Sigma}\|_{\text{op}}\sqrt{d})$.

We now focus on bounding $\mathbb{E}[\mathcal{V}(x - \Sigma\xi)]$. We would like to use a Taylor expansion, but care needs to be taken to handle the non-differentiability of \mathcal{V} at 0. Under the expectation, we distinguish two events: the event on which $\|\xi\| < R_\varepsilon$, which supports the main mass of ν , and the event on which $\|\xi\| \geq R_\varepsilon$, corresponding to the tails.

(a) For the first event we consider (on which $\|\xi\| < R_\varepsilon$), for any $x \notin \mathcal{B}_2(\|\Sigma\|_{\text{op}}R_\varepsilon)$, we must have $0 \notin \mathcal{B}_2(x, \|\Sigma\xi\|)$, and thus $0 \notin (x + \Delta\Sigma\xi)_{\Delta \in [0, 1]}$. Since this line segment doesn't contain 0 (the only point at which \mathcal{V} is not continuously differentiable), we can perform a second-order Taylor expansion of \mathcal{V} to obtain

$$\begin{aligned} \mathbb{E}[\mathcal{V}(x + \Sigma\xi)\mathbf{1}_{\{\|\xi\| < R_\varepsilon\}}] \\ \leq \mathbb{E}\left[\left(\mathcal{V}(x) + \xi^\top \Sigma^\top \nabla \mathcal{V}(x) + \frac{1}{2} \text{Tr}[\Sigma\xi\xi^\top \Sigma^\top \nabla^2 \mathcal{V}(\hat{x})]\right)\mathbf{1}_{\{\|\xi\| < R_\varepsilon\}}\right] \end{aligned}$$

for some $\hat{x} \in (x + \Delta\Sigma\xi)_{\Delta \in [0, 1]}$. By the Cauchy-Schwartz inequality and the derivative bounds of Assumption 2, we obtain

$$\begin{aligned} \mathbb{E}[\mathcal{V}(x + \Sigma\xi)\mathbf{1}_{\{\|\xi\|_2 < R_\varepsilon\}}] &\leq \mathcal{V}(x) + \mathbb{E}[\xi^\top \mathbf{1}_{\{\|\xi\| < R_\varepsilon\}}]\Sigma^\top \nabla \mathcal{V}(x) + \frac{\varepsilon}{2}M'_{\mathcal{V}}\|\bar{\Sigma}\|_{\text{op}}^2 \\ &\leq \mathcal{V}(x) + \frac{\varepsilon}{2}M'_{\mathcal{V}}\|\bar{\Sigma}\|_{\text{op}}^2, \end{aligned}$$

since $\mathbb{E}[\xi^\top \mathbf{1}_{\{\|\xi\| < R_\varepsilon\}}] = 0$ by the rotational invariance property of a truncated Gaussian.

(b) On the second event (on which $\|\xi\| \geq R_\varepsilon$), we cannot use a Taylor expansion. Instead, we use the Lipschitzness of \mathcal{V} followed by the Cauchy-Schwartz inequality,

and then apply a sub-Gaussian concentration inequality (see e.g. (Ledoux and Talagrand, 1991, (3.5))):

$$\begin{aligned}
 \mathbb{E}[\mathcal{V}(x + \Sigma\xi)\mathbf{1}_{\{\|\xi\|\geq R_\varepsilon\}}] &\leq \mathcal{V}(x) + M_\mathcal{V} \|\Sigma\|_{\text{op}} \mathbb{E}[\|\xi\| \mathbf{1}_{\{\|\xi\|\geq R_\varepsilon\}}] \\
 &\leq \mathcal{V}(x) + M_\mathcal{V} \|\Sigma\|_{\text{op}} \sqrt{\mathbb{E}[\|\xi\|^2] \mathbb{P}(\|\xi\| \geq R_\varepsilon)} \\
 &\leq \mathcal{V}(x) + M_\mathcal{V} \|\Sigma\|_{\text{op}} \sqrt{4de^{-\frac{R_\varepsilon^2}{8d}}} \\
 &\leq \mathcal{V}(x) + 2\varepsilon M_\mathcal{V} \|\bar{\Sigma}\|_{\text{op}} \sqrt{d}.
 \end{aligned}$$

To complete the proof, we combine both cases in (25), and let

$$\mathbf{c}'_\mathcal{V} := M_\mathcal{V} L_0(1 + \|\bar{\Sigma}\|_{\text{op}} \sqrt{d}) + 2M_\mathcal{V} \|\bar{\Sigma}\|_{\text{op}} \sqrt{d} + \frac{M'_\mathcal{V}}{2} \|\bar{\Sigma}\|_{\text{op}}^2.$$

■

Lemma 12 *Under Assumptions 1 and 2, for any $(x_0, \alpha, \theta) \in \mathbb{R}^d \times \mathcal{A} \times \Theta$ and any $\lambda \in \mathbb{R}_+$, we have*

$$\mathbb{E} \left[e^{\lambda \mathcal{V}(X_{\tau_n}^{x_0, \alpha, \theta})} \right] \leq (n+1) \exp \left(\lambda \left(\frac{\mathbf{c}'_\mathcal{V}}{\mathbf{c}_\mathcal{V}} + L_\mathcal{V} (\varepsilon^{\frac{1}{2}} \|\bar{\Sigma}\|_{\text{op}} R_\varepsilon + \|x_0\|) \right) + \frac{\lambda^2 M_\mathcal{V}^2 \|\bar{\Sigma}\|_{\text{op}}^2}{2\mathbf{c}_\mathcal{V}} \right),$$

for any $n \in \mathbb{N}$.

Proof For $n \in \mathbb{N}^*$, let us define the following events for $i < n$: $E_{i, n-1} := \{i = \sup\{j \in \{0, \dots, n-1\} : \|X_{\tau_j}^{\alpha, \theta}\| \leq \|\Sigma\|_{\text{op}} R_\varepsilon\}\}$ and $\bar{E}_{n-1} := \{\min_{j \in \{0, \dots, n-1\}} \|X_{\tau_j}^{\alpha, \theta}\| > \|\Sigma\|_{\text{op}} R_\varepsilon\}$. Note that both these events are $\mathcal{F}_{\tau_{n-1}}$ -measurable and that $\cup_{i \leq n-1} E_{i, n-1} = \bar{E}_{n-1}^c$, so that $\{\bar{E}_{n-1}, E_{0, n-1}, \dots, E_{n-1, n-1}\}$ induces a partition of Ω . We begin by working conditionally on each of these events, and in a second part we will collect them to bound $\mathbb{E}[\exp(\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta}))]$.

For any $0 \leq i < n$, by adding and subtracting $\mathbb{E}[\exp(\mathbb{E}[\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta}) | \mathcal{F}_{\tau_{n-1}}]) \mathbf{1}_{E_{i, n-1}}]$ and by the tower rule, we have

$$\begin{aligned}
 \mathbb{E} \left[e^{\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta})} \mathbf{1}_{E_{i, n-1}} \right] &= \mathbb{E} \left[\mathbb{E} \left[e^{\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta})} | \mathcal{F}_{\tau_{n-1}} \right] \mathbf{1}_{E_{i, n-1}} \right] \\
 &= \mathbb{E} \left[\exp \left(\mathbb{E}[\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta}) | \mathcal{F}_{\tau_{n-1}}] \right) \mathbf{1}_{E_{i, n-1}} \right. \\
 &\quad \left. \times \mathbb{E} \left[\exp \left(\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta}) - \mathbb{E}[\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta}) | \mathcal{F}_{\tau_{n-1}}] \right) | \mathcal{F}_{\tau_{n-1}} \right] \right].
 \end{aligned}$$

Using a result for Lipschitz functions of Gaussian random variables (see e.g. Boucheron et al. (2013, Thm 5.5)) applied to \mathcal{V} and ξ , we obtain

$$\begin{aligned}
 \mathbb{E} \left[e^{\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta})} \mathbf{1}_{E_{i, n-1}} \right] &\leq e^{\frac{\lambda^2}{2} M_\mathcal{V}^2 \|\Sigma\|_{\text{op}}^2} \mathbb{E} \left[\exp \left(\mathbb{E}[\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta}) | \mathcal{F}_{\tau_{n-1}}] \right) \mathbf{1}_{E_{i, n-1}} \right] \\
 &= e^{\frac{\lambda^2}{2} M_\mathcal{V}^2 \|\Sigma\|_{\text{op}}^2} \mathbb{E} \left[\exp \left(\mathbb{E}[\lambda \mathcal{V}(\psi_{\bar{\theta}}^\varepsilon(X_{\tau_{n-1}}^{\alpha, \theta}, \alpha_{\tau_{n-1}}) + \Sigma\xi_n) | \mathcal{F}_{\tau_{n-1}}] \right) \mathbf{1}_{E_{i, n-1}} \right].
 \end{aligned} \tag{26}$$

If $i = n - 1$, $\|X_{\tau_{n-1}}^{\alpha, \theta}\| \leq \|\Sigma\|_{\text{op}} R_\varepsilon$ on the event $E_{i, n-1}$, and thus we have

$$\begin{aligned} \mathbb{E} \left[\lambda \mathcal{V}(\psi_\theta^\varepsilon(X_{\tau_{n-1}}^{\alpha, \theta}, \alpha_{\tau_{n-1}}) + \Sigma \xi_n) | \mathcal{F}_{\tau_{n-1}} \right] &\leq \mathbb{E} \left[\lambda L_\gamma \left\| X_{\tau_{n-1}}^{\alpha, \theta} + \mu(X_{\tau_{n-1}}^{\alpha, \theta}, \alpha_{\tau_{n-1}}) + \Sigma \xi \right\| | \mathcal{F}_{\tau_{n-1}} \right] \\ &\leq \lambda L_\gamma \left((1 + L_0) \|\Sigma\|_{\text{op}} R_\varepsilon + 1 + \|\Sigma\|_{\text{op}} \sqrt{d} \right) \end{aligned}$$

by using the fact that $\mathbb{E}[\|\xi\|] \leq \sqrt{\mathbb{E}[\|\xi\|^2]} = \sqrt{d}$, as $\xi \sim \nu$. Noticing that $\sup_{\varepsilon \in (0,1)} \varepsilon^{\frac{1}{2}} R_\varepsilon = \sqrt{8de^{-1}}$, let us introduce

$$C_H := L_\gamma \left((1 + L_0) \|\bar{\Sigma}\|_{\text{op}} \sqrt{8de^{-1}} + 1 + \|\bar{\Sigma}\|_{\text{op}} \sqrt{d} \right). \quad (27)$$

Combining this with (26) yields

$$\mathbb{E}[e^{\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta})} \mathbf{1}_{E_{i, n-1}}] \leq \exp \left(\frac{\lambda^2}{2} M_\gamma^2 \|\Sigma\|_{\text{op}}^2 + \lambda C_H \right), \quad (28)$$

in the case $i = n - 1$.

If $i < n - 1$, we can apply the same methodology, and continuing from (26) apply Lemma 11 to obtain

$$\begin{aligned} \mathbb{E} \left[e^{\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta})} \mathbf{1}_{E_{i, n-1}} \right] &\leq e^{\frac{\lambda^2}{2} M_\gamma^2 \|\Sigma\|_{\text{op}}^2} \mathbb{E} \left[\exp \left(\mathbb{E} \left[\lambda \mathcal{V}(\psi_\theta^\varepsilon(X_{\tau_{n-1}}^{\alpha, \theta}, \alpha_{\tau_{n-1}}) + \Sigma \xi_n) | \mathcal{F}_{\tau_{n-1}} \right] \right) \right. \\ &\quad \left. \times \mathbf{1}_{\{X_{\tau_{n-1}}^{\alpha, \theta} > \|\Sigma\|_{\text{op}} R_\varepsilon\}} \mathbf{1}_{E_{i, n-2}} \right], \quad (29) \\ &\leq e^{\frac{\lambda^2}{2} M_\gamma^2 \|\Sigma\|_{\text{op}}^2 + \lambda \varepsilon \mathbf{c}'_\gamma} \mathbb{E}[\exp((1 - \varepsilon \mathbf{c}_\gamma) \lambda \mathcal{V}(X_{\tau_{n-1}}^{\alpha, \theta})) \mathbf{1}_{E_{i, n-2}}]. \end{aligned}$$

It remains to use an induction argument in n down to $n = i + 1$ and use the fact that $\|X_{\tau_i}^{\alpha, \theta}\| \leq \|\Sigma\|_{\text{op}} R_\varepsilon$ on $E_{i, i}$, to obtain

$$\begin{aligned} &\mathbb{E} \left[e^{\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta})} \mathbf{1}_{E_{i, n-1}} \right] \\ &\leq \exp \left(\lambda C_H + \lambda \varepsilon \mathbf{c}'_\gamma \sum_{k=0}^{n-1-i} (1 - \varepsilon \mathbf{c}_\gamma)^k + \frac{\lambda^2 M_\gamma^2 \|\Sigma\|_{\text{op}}^2}{2} \sum_{k=0}^{n-1-i} (1 - \varepsilon \mathbf{c}_\gamma)^{2k} \right) \\ &\leq \exp \left(\lambda C_H + \lambda \frac{\mathbf{c}'_\gamma}{\mathbf{c}_\gamma} + \frac{\lambda^2 M_\gamma^2 \|\bar{\Sigma}\|_{\text{op}}^2}{2 \mathbf{c}_\gamma} \right). \quad (30) \end{aligned}$$

On the event \bar{E}_{n-1} , that is if the process is never in the ball $\mathcal{B}_2(\|\Sigma\|_{\text{op}} R_\varepsilon)$ before time τ_n , we use the fact that (29) is valid with \bar{E}_{n-1} and \bar{E}_{n-2} in place of $E_{i, n-1}$ and $E_{i, n-2}$. Applying the induction, we obtain

$$\mathbb{E} \left[e^{\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta})} \mathbf{1}_{\bar{E}_{n-1}} \right] \leq \exp \left(\lambda L_\gamma \|x_0\| + \lambda \frac{\mathbf{c}'_\gamma}{\mathbf{c}_\gamma} + \frac{\lambda^2 M_\gamma^2 \|\bar{\Sigma}\|_{\text{op}}^2}{2 \mathbf{c}_\gamma} \right). \quad (31)$$

Using our partition and combining (28), (30), and (31) we can thus write, for any $n \in \mathbb{N}$

$$\begin{aligned} \mathbb{E} \left[e^{\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta})} \right] &\leq \mathbb{E} \left[e^{\lambda \mathcal{V}(X_{\tau_n}^{\alpha, \theta})} \left(\mathbf{1}_{\bar{E}_{n-1}} + \sum_{i=0}^{n-1} \mathbf{1}_{E_{i, n-1}} \right) \right] \\ &\leq (n + 1) \exp \left(\lambda \left(\frac{\mathbf{c}'_\gamma}{\mathbf{c}_\gamma} + C_H + L_\gamma \|x_0\| \right) + \frac{\lambda^2 M_\gamma^2 \|\bar{\Sigma}\|_{\text{op}}^2}{2 \mathbf{c}_\gamma} \right) \end{aligned}$$

which concludes the proof. \blacksquare

With these two lemmas, we can now prove Proposition 2, the main result of this section. First, let us give the exact definition of $H_\delta(n)$:

$$H_\delta(n) := \frac{1}{\ell_\gamma} (C_H + L_\gamma \|x_0\|) + \frac{\mathbf{c}'_\gamma}{\ell_\gamma \mathbf{c}_\gamma} + \frac{M_\gamma}{\ell_\gamma} \|\bar{\Sigma}\|_{\text{op}} \sqrt{\frac{2}{\mathbf{c}_\gamma} \log\left(\frac{\pi^2(n+1)^3}{6\delta}\right)} \quad (32)$$

in which C_H is defined in (27), so that $H_\delta(n) = \mathcal{O}(\sqrt{\log(n\delta^{-1})})$.

Proposition 2 *Under Assumptions 1 and 2, there is a function $H_\delta(n) = \mathcal{O}(\sqrt{\log(n\delta^{-1})})$ such that for any $\delta \in (0, 1)$, $\alpha \in \mathcal{A}$, $x_0 \in \mathbb{R}^d$, and $\theta \in \Theta$ we have*

$$\mathbb{P}\left(\sup_{t \in \mathbb{R}_+} \frac{\|X_t^{\alpha, \theta}\|}{H_\delta(N_t)} \geq 1\right) \leq \delta. \quad (9)$$

Proof Fix $n \in \mathbb{N}$, by Markov's inequality and Assumption 2, for any $u > 0$, we have

$$\mathbb{P}\left(\|X_{\tau_n}^{\alpha, \theta}\| > u\right) \leq \mathbb{E}\left[e^{\lambda \ell_\gamma \|X_{\tau_n}^{\alpha, \theta}\|}\right] e^{-\lambda \ell_\gamma u} \leq \mathbb{E}\left[e^{\lambda \gamma (X_{\tau_n}^{\alpha, \theta})}\right] e^{-\lambda \ell_\gamma u},$$

which implies that

$$\begin{aligned} & \mathbb{P}\left(\|X_{\tau_n}^{\alpha, \theta}\| - \frac{\mathbf{c}'_\gamma}{\ell_\gamma \mathbf{c}_\gamma} - \frac{C_H}{\ell_\gamma} - \frac{L_\gamma}{\ell_\gamma} \|x_0\| > u\right) \\ & \leq \mathbb{E}\left[e^{\lambda \gamma (X_{\tau_n}^{\alpha, \theta})}\right] \exp\left(-\lambda \ell_\gamma \left(u + \frac{\mathbf{c}'_\gamma}{\ell_\gamma \mathbf{c}_\gamma} + \frac{C_H}{\ell_\gamma} + \frac{L_\gamma}{\ell_\gamma} \|x_0\|\right)\right). \end{aligned}$$

Applying Lemma 12, and taking $\lambda = \mathbf{c}_\gamma \ell_\gamma u / (M_\gamma^2 \|\bar{\Sigma}\|_{\text{op}}^2)$, we obtain

$$\begin{aligned} & \mathbb{P}\left(\|X_{\tau_n}^{\alpha, \theta}\| > u + \frac{\mathbf{c}'_\gamma}{\ell_\gamma \mathbf{c}_\gamma} + \varepsilon^{\frac{1}{2}} \frac{L_\gamma}{\ell_\gamma} \|\bar{\Sigma}\|_{\text{op}} R_\varepsilon + \frac{L_\gamma}{\ell_\gamma} \|x_0\|\right) \\ & \leq (n+1) \exp\left(-\lambda \ell_\gamma u + \lambda^2 \frac{M_\gamma^2 \|\bar{\Sigma}\|_{\text{op}}^2}{2\mathbf{c}_\gamma}\right) \\ & = (n+1) \exp\left(-\frac{\mathbf{c}_\gamma \ell_\gamma^2}{2M_\gamma^2 \|\bar{\Sigma}\|_{\text{op}}^2} u^2\right). \end{aligned}$$

Letting $u = M_\gamma \|\bar{\Sigma}\|_{\text{op}} \ell_\gamma^{-1} \sqrt{2\mathbf{c}_\gamma^{-1} \log((n+1)/\delta')}$, yields

$$\mathbb{P}\left(\|X_{\tau_n}^{\alpha, \theta}\| \geq \frac{C_H}{\ell_\gamma} + \frac{L_\gamma}{\ell_\gamma} \|x_0\| + \frac{\mathbf{c}'_\gamma}{\ell_\gamma \mathbf{c}_\gamma} + \frac{M_\gamma}{\ell_\gamma} \|\bar{\Sigma}\|_{\text{op}} \sqrt{\frac{2}{\mathbf{c}_\gamma} \log\left(\frac{n+1}{\delta'}\right)}\right) \leq \delta'.$$

Setting $\delta' = 6\delta/\pi^2(n+1)^2$, and taking a union bound over $n \in \mathbb{N}$ yields

$$\mathbb{P}\left(\sup_{t \in \mathbb{R}_+} \frac{X_t^{\alpha, \theta}}{H_\delta(N_t)} \geq 1\right) = \mathbb{P}\left(\bigcup_{n \in \mathbb{N}} \{\|X_{\tau_n}^{\alpha, \theta}\| \geq H_\delta(n)\}\right) \leq \delta,$$

which implies the result since $\delta \in (0, 1)$ implies $\log(n^3/\delta) \leq \log(n^3/\delta^3) = 3 \log(n/\delta)$. \blacksquare

B.2. Expectation Bounds of Higher Orders

In this appendix, we will focus on higher moment conditions of the state process, which will be used in the control results of Appendix D. In Lemma 13 and Theorem 14 we work to raise the stochastic stability condition from Lemma 11 to a power $p \geq 2$. Lemma 15, the main result of this section, will follow from this by arguments of Abeille et al. (2022).

Lemma 13 *Under Assumptions 1 and 2, for $p \geq 2$, there is a function $g : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$ and a constant $C_p > 0$ independent of ε satisfying*

$$g(x, \eta) \leq \varepsilon C_p (1 + \mathcal{V}(x - \sqrt{\varepsilon}\eta)^{p-1}) (1 + \|\eta\|^p),$$

for any $(\eta, x) \in \mathbb{R}^d \times \mathbb{R}^d$, such that

$$\mathcal{V}(\psi_\theta^\varepsilon(x, a) - \sqrt{\varepsilon}\eta)^p \leq (1 - \varepsilon c_\gamma) \mathcal{V}(x - \sqrt{\varepsilon}\eta)^p + g(x, \eta). \quad (33)$$

for any $(\eta, x, a, \theta) \in \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{A} \times \Theta$.

Proof We first raise both sides of (24) to the power p

$$\mathcal{V}(\psi_\theta^\varepsilon(x, a) - \sqrt{\varepsilon}\eta)^p \leq \left((1 - \varepsilon c_\gamma) \mathcal{V}(x - \sqrt{\varepsilon}\eta) + \varepsilon M_\gamma L_0 (1 + \|\eta\|) \right)^p.$$

We will now expand the right-hand side. Let $a = (1 - \varepsilon c_\gamma) \mathcal{V}(x - \sqrt{\varepsilon}\eta)$ and $b = \varepsilon M_\gamma L_0 (1 + \|\eta\|)$, by the binomial theorem we have

$$\begin{aligned} (a + b)^p &= \sum_{k=0}^p \binom{p}{k} a^k b^{p-k} = a^p + b \sum_{k=0}^{p-1} \binom{p}{k} a^k b^{p-1-k} \\ &\leq a^p + b(1 + b)^{p-1} (1 + a)^{p-1} \sum_{k=0}^{p-1} \binom{p}{k}. \end{aligned}$$

Since $(1 - \varepsilon c_\gamma) \in (0, 1)$, $\varepsilon \leq 1$, $b \leq 1 + b$, and $\sum_{k=0}^{p-1} \binom{p}{k} \leq 2^p$, by using the binomial identity $(1 + a)^q \leq 2^{q-1} (1 + a^q)$ for $(a, q) \in [0, +\infty) \times [1, +\infty)$, we have

$$\begin{aligned} \mathcal{V}(\psi_\theta^\varepsilon(x, a) - \sqrt{\varepsilon}\eta)^p &\leq (1 - \varepsilon c_\gamma) \mathcal{V}(x - \sqrt{\varepsilon}\eta)^p \\ &\quad + \varepsilon (1 + M_\gamma L_0 (1 + \|\eta\|))^p (1 + \mathcal{V}(x - \sqrt{\varepsilon}\eta)^{p-1}) 2^{p-2+p}. \end{aligned} \quad (34)$$

Finally, we have

$$\begin{aligned} (1 + M_\gamma L_0 (1 + \|\eta\|))^p &= (1 + M_\gamma L_0 + M_\gamma L_0 \|\eta\|)^p \\ &\leq (1 + M_\gamma L_0 + (1 + M_\gamma L_0) \|\eta\|)^p \\ &= (1 + M_\gamma L_0)^p (1 + \|\eta\|)^p \\ &\leq (1 + M_\gamma L_0)^p (1 + \|\eta\|^p) 2^{p-1}. \end{aligned} \quad (35)$$

Combining (34) and (35) leads to the required result. ■

Recall that $\xi \sim \nu$ is a centred standard Gaussian random variable.

Corollary 14 *Under Assumptions 1 and 2, for any $p \geq 2$, there is a constant $\mathfrak{c}_p > 0$ independent of ε such that*

$$\mathbb{E}[\mathcal{V}(\psi_\theta^\varepsilon(x, a) + \Sigma\xi)^p] \leq \left(1 - \varepsilon \frac{\mathfrak{c}_\mathcal{V}}{2}\right) \mathbb{E}[\mathcal{V}(x - \sqrt{\varepsilon}\xi)^p] + \varepsilon \mathfrak{c}_p$$

for any $(x, a, \theta) \in \mathbb{R}^d \times \mathbb{A} \times \Theta$.

Proof

- i. Taking the expectation of the bound on g from Lemma 13 and applying Hölder's inequality yields

$$\begin{aligned} \mathbb{E}[g(x, \xi)] &\leq \varepsilon C_p \mathbb{E}[(1 + \mathcal{V}(x - \sqrt{\varepsilon}\xi)^{p-1})(1 + \|\xi\|^p)] \\ &\leq \varepsilon C_p \mathbb{E}\left[(1 + \mathcal{V}(x - \sqrt{\varepsilon}\xi)^{p-1})^{\frac{p}{p+1}}\right]^{\frac{p+1}{p}} \mathbb{E}[(1 + \|\xi\|^p)^{p+1}]^{\frac{1}{p+1}} \\ &\leq 4\varepsilon C_p \mathbb{E}\left[1 + \mathcal{V}(x - \sqrt{\varepsilon}\xi)^{\frac{(p-1)(p+1)}{p}}\right] \mathbb{E}[(1 + \|\xi\|^p)^{p+1}]^{\frac{1}{p+1}}, \end{aligned}$$

by using the identities: for $(u, v) \in \mathbb{R}_+^2$, $(1 + u)^{(p+1)/p} \leq 4(1 + u^{(p+1)/p})$ and $(1 + v)^{p/(p+1)} \leq 1 + v$. Since ξ has bounded moments of any order,

$$C'_p := 4C_p \mathbb{E}[(1 + \|\xi\|^p)^{p+1}]^{\frac{1}{p+1}}$$

is a finite constant and we have

$$\mathbb{E}[g(x, \xi)] \leq \varepsilon C'_p \mathbb{E}\left[1 + \mathcal{V}(x - \sqrt{\varepsilon}\xi)^{p-\frac{1}{p}}\right].$$

- ii. Recalling Lemma 13, we have

$$\begin{aligned} \mathbb{E}[\mathcal{V}(\psi_\theta^\varepsilon(x, a) + \Sigma\xi)^p] &\leq (1 - \varepsilon \mathfrak{c}_\mathcal{V}) \mathbb{E}[\mathcal{V}(x - \sqrt{\varepsilon}\xi)^p] + \mathbb{E}[g(x, \xi)] \\ &\leq \left(1 - \varepsilon \frac{\mathfrak{c}_\mathcal{V}}{2}\right) \mathbb{E}[\mathcal{V}(x - \sqrt{\varepsilon}\xi)^p] \\ &\quad + \varepsilon \mathbb{E}\left[C'_p(1 + \mathcal{V}(x - \sqrt{\varepsilon}\xi)^{p-\frac{1}{p}}) - \frac{\mathfrak{c}_\mathcal{V}}{2} \mathcal{V}(x - \sqrt{\varepsilon}\xi)^p\right]. \quad (36) \end{aligned}$$

- iii. Note that, for any $p \geq 2$, the function

$$z \in \mathbb{R}^d \mapsto \frac{\|z\|^{p-\frac{1}{p}}}{1 + \|z\|^p} \in \mathbb{R}_+$$

is bounded, so there exists a constant $C''_p > 0$ such that, for any $z \in \mathbb{R}^d$,

$$C'_p \mathcal{V}(z)^{p-\frac{1}{p}} - \frac{\mathfrak{c}_\mathcal{V}}{2} \mathcal{V}(z)^p \leq C''_p.$$

Applying this to the expectation in (36), we have

$$\mathbb{E}[\mathcal{V}(\psi_\theta^\varepsilon(x, a) + \Sigma\xi)^p] \leq \left(1 - \varepsilon \frac{\mathfrak{c}_\mathcal{V}}{2}\right) \mathbb{E}[\mathcal{V}(x + \sqrt{\varepsilon}\xi)^p] + \varepsilon(C''_p + C'_p).$$

Letting $\mathfrak{c}_p := C'_p + C''_p$ completes the proof.

■

Lemma 15 *Under Assumptions 1 and 2, for any $p \geq 2$, there is a constant $\mathbf{c}'_p > 0$ independent of ε such that*

$$\mathbb{E} \left[\|X_t^{x_0, \alpha, \theta}\|^p \right] \leq \frac{1}{\ell_{\mathcal{V}}^p} \left(L_{\mathcal{V}}^p e^{-\frac{\mathbf{c}_{\mathcal{V}}}{4}t} \|x_0\|^p + \frac{4\mathbf{c}'_p}{\mathbf{c}_{\mathcal{V}}} \left(1 - e^{-\frac{\mathbf{c}_{\mathcal{V}}}{4}t} \right) \right),$$

for any $(x_0, \alpha, \theta) \in \mathbb{R}^d \times \mathcal{A} \times \Theta$ and $t \in [0, +\infty)$.

Proof Recall from Theorem 14 that we have

$$\mathbb{E} [\mathcal{V}(\psi_{\theta}^{\varepsilon}(x, a) + \Sigma\xi)^p] \leq \left(1 - \varepsilon \frac{\mathbf{c}_{\mathcal{V}}}{2} \right) \mathbb{E} [\mathcal{V}(x + \Sigma\xi)^p] + \varepsilon \mathbf{c}_p \quad (37)$$

for any $(x, a, \theta) \in \mathbb{R}^d \times \mathbb{A} \times \Theta$. We begin by eliminating the $\Sigma\xi$ from the right-hand side so that we have a proper Lyapunov contraction property on the generator. We expand $\mathcal{V}^p \in \mathcal{C}^2(\mathbb{R}^d; [0, +\infty))$ and use the fact that $\mathbb{E}[\xi] = 0$ to obtain

$$\begin{aligned} \mathbb{E} [\mathcal{V}(x + \Sigma\xi)^p] &= \mathcal{V}(x)^p + \varepsilon p \mathbb{E} \left[\mathcal{V}(x + \Delta\Sigma\xi)^{p-1} \text{Tr}[\xi \bar{\Sigma} \bar{\Sigma}^{\top} \xi^{\top} \nabla^2 \mathcal{V}(x + \Delta\Sigma\xi)] \right] \\ &\quad + \varepsilon p(p-1) \mathbb{E} \left[\mathcal{V}(x + \Delta\Sigma\xi)^{p-2} \text{Tr} \left[\xi \bar{\Sigma} \bar{\Sigma}^{\top} \xi^{\top} \nabla \mathcal{V}(x + \Delta\Sigma\xi) \nabla \mathcal{V}(x + \Delta\Sigma\xi)^{\top} \right] \right] \end{aligned}$$

for some random variable Δ taking value in $[0, 1]$. This is now upper-bounded by using the Lipschitzness of \mathcal{V} and the Cauchy-Schwartz inequality

$$\begin{aligned} \mathbb{E} [\mathcal{V}(x + \Sigma\xi)^p] &\leq \mathcal{V}(x)^p + \varepsilon p M'_{\mathcal{V}} \|\bar{\Sigma}\|_{\text{op}}^2 \mathbb{E} \left[(\mathcal{V}(x) + M_{\mathcal{V}} \Delta \|\xi\|)^{p-1} \|\xi\|^2 \right] \\ &\quad + \varepsilon p(p-1) (M_{\mathcal{V}})^2 \|\bar{\Sigma}\|_{\text{op}}^2 \mathbb{E} \left[(\mathcal{V}(x) + M_{\mathcal{V}} \Delta \|\xi\|)^{p-2} \|\xi\|^2 \right]. \end{aligned}$$

By the binomial theorem as in the proof of Lemma 13, and as $|\Delta| \leq 1$, we have

$$\begin{aligned} \mathbb{E} [\mathcal{V}(x + \Sigma\xi)^p] &\leq \mathcal{V}(x)^p + \varepsilon \left(p M'_{\mathcal{V}} \|\bar{\Sigma}\|_{\text{op}}^2 \mathbb{E} \left[\|\xi\|^2 \sum_{k=0}^{p-1} \binom{p-1}{k} \mathcal{V}(x)^k (M_{\mathcal{V}} \|\Sigma\|_{\text{op}} \|\xi\|)^{p-1-k} \right] \right. \\ &\quad \left. + p(p-1) (M_{\mathcal{V}} \|\bar{\Sigma}\|_{\text{op}})^2 \mathbb{E} \left[\sum_{k=0}^{p-2} \binom{p-2}{k} \mathcal{V}(x)^k (M_{\mathcal{V}} \|\Sigma\|_{\text{op}} \|\xi\|)^{p-2-k} \right] \right). \end{aligned}$$

Since $\|\xi\|$ is a sub-Gaussian random variable it has moments of all orders, and we can express the interior of the bracket above as a polynomial in $\mathcal{V}(x)$ of order $p-1$ with finite coefficients $\{a_k\}_{k=0}^{p-1} \subset \mathbb{R}_+$. Recalling (37), we thus have

$$\begin{aligned} \mathbb{E} [\mathcal{V}(\psi_{\theta}^{\varepsilon}(x, a) + \Sigma\xi)^p] &\leq (1 - \varepsilon \mathbf{c}_{\mathcal{V}}) \left(\mathcal{V}(x)^p + \varepsilon \sum_{k=0}^{p-1} a_k \mathcal{V}(x)^k \right) + \varepsilon \mathbf{c}_p \\ &\leq \left(1 - \varepsilon \frac{\mathbf{c}_{\mathcal{V}}}{4} \right) \mathcal{V}(x)^p + \varepsilon \left(\mathbf{c}_p - \frac{\mathbf{c}_{\mathcal{V}}}{4} \mathcal{V}(x)^p + \sum_{k=0}^{p-1} a_k \mathcal{V}(x)^k \right) \end{aligned}$$

As in part iii. of the proof of Theorem 14, the interior of the second bracket is a continuous function which goes to $-\infty$ as $\|x\| \rightarrow +\infty$, so there must be a constant $\mathbf{c}'_p \in \mathbb{R}_+$ (independent of ε) such that

$$\mathbf{c}_p + \sup_{x \in \mathbb{R}^d} \left(-\frac{\mathbf{c}_\gamma}{4} \mathcal{V}(x)^p + \sum_{k=0}^{p-1} a_k \mathcal{V}(x)^k \right) \leq \mathbf{c}'_p < +\infty.$$

Therefore, we have the desired Lyapunov generator condition

$$\mathbb{E} [\mathcal{V}(\psi_\theta^\varepsilon(x, a) + \Sigma \xi)^p] \leq \left(1 - \varepsilon \frac{\mathbf{c}_\gamma}{4} \right) \mathcal{V}(x)^p + \varepsilon \mathbf{c}'_p,$$

which is equivalently written for any $(x, a) \in \mathbb{R}^d \times \mathbb{A}$ as

$$\frac{1}{\varepsilon} \int (\mathcal{V}(\psi_\theta^\varepsilon(x, a) + \Sigma e)^p - \mathcal{V}(x)^p) \nu(\mathrm{d}e) \leq -\frac{\mathbf{c}_\gamma}{4} \mathcal{V}(x)^p + \mathbf{c}'_p. \quad (38)$$

By Itô's Lemma, (38), and a localisation argument, we have, for any $t \geq t_0 \geq 0$, that

$$\begin{aligned} \mathbb{E} \left[\mathcal{V}(X_t^{x_0, \alpha, \theta})^p \right] &= \mathbb{E} \left[\mathcal{V}(X_{t_0}^{x_0, \alpha, \theta})^p \right] \\ &\quad + \mathbb{E} \left[\int_{t_0}^t \frac{1}{\varepsilon} \int (\mathcal{V}(\psi_\theta^\varepsilon(X_s^{x_0, \alpha, \theta}, \alpha_s) + \Sigma e)^p - \mathcal{V}(X_s^{x_0, \alpha, \theta})^p) \nu(\mathrm{d}e) \mathrm{d}s \right] \\ &\leq \mathbb{E} \left[\mathcal{V}(X_{t_0}^{x_0, \alpha, \theta})^p \right] - \frac{\mathbf{c}_\gamma}{4} \int_{t_0}^t \mathbb{E} \left[\mathcal{V}(X_s^{x_0, \alpha, \theta})^p \right] \mathrm{d}s + (t - t_0) \mathbf{c}'_p. \end{aligned}$$

By a simple comparison argument for ODEs, we then obtain

$$\mathbb{E} \left[\mathcal{V}(X_t^{x_0, \alpha, \theta})^p \right] \leq e^{-\frac{\mathbf{c}_\gamma}{4} t} \mathcal{V}(x_0)^p + \frac{4\mathbf{c}'_p}{\mathbf{c}_\gamma} \left(1 - e^{-\frac{\mathbf{c}_\gamma}{4} t} \right).$$

Using now Assumption 2, we obtain

$$\mathbb{E} \left[\|X_t^{x_0, \alpha, \theta}\|^p \right] \leq \frac{1}{\ell_\gamma^p} \left(L_\gamma^p e^{-\frac{\mathbf{c}_\gamma}{4} t} \|x_0\|^p + \frac{4\mathbf{c}'_p}{\mathbf{c}_\gamma} \left(1 - e^{-\frac{\mathbf{c}_\gamma}{4} t} \right) \right).$$

■

Appendix C. Concentration Inequality and Online Prediction Error

The key result of this section, Proposition 5, builds heavily on (Russo and Van Roy, 2013, Prop. 5). Proposition 5 differs from this existing result in three ways. First, it is *any-time* i.e. does not require *a priori* knowledge of a time horizon. This is a minor technical refinement, but it is of practical importance. Second, it applies to a pure-jump process defined on \mathbb{R}_+ . This apparent complexity vanishes when the filtration of the pure-jump process is chosen correctly, as the state process is piece-wise constant. Third, and most important, it applies to learning in a function class (\mathcal{F}_Θ) of unbounded drifts for an unbounded process $X^{\alpha, \theta}$, which is an inherent difficulty in handling continuous state RL problems.

This third extension is non-trivial and leads us to significantly reshuffle the proof structure of (Russo and Van Roy, 2013) and to incorporate some self-normalised inequality arguments as well as high-probability bounds on the state from Appendix B. While many of the original ideas are still used, the way they link together has changed and thus we will include, in Appendix C.1, a complete derivation for the sake of clarity. In this spirit, we will prove a generic result (Theorem 18), which itself implies Proposition 5.

Proposition 5 (Adapted from Osband and Van Roy (2014, Prop. 5)) *Under Assumptions 1 and 2, for any $x_0 \in \mathbb{R}^d$, and $\delta > 0$,*

$$\mathbb{P} \left(\left\{ \theta^* \in \bigcap_{n=1}^{\infty} \mathcal{C}_n(\delta) \right\} \cap \left\{ \sup_{n \in \mathbb{N}^*} \frac{\|X_{\tau_n}^{\varpi, \theta^*}\|}{H_\delta(n)} \leq 1 \right\} \right) \geq 1 - \delta, \quad (13)$$

Proposition 5 ensures that the sets $(\mathcal{C}_n(\delta))_{n \in \mathbb{N}}$ defined in (6) are valid confidence sets. In order to bound the regret, we need to go further and to bound the online prediction error of functions within these confidence sets along the trajectory (see. (58)).

For any $n \in \mathbb{R}$, let $d_{E,n}$ denotes the $2\sqrt{\varepsilon/n}$ -eluder dimension of the model class restricted to the set $B_n := \mathcal{B}_2(\sup_{s \leq \tau_n} \|X_s^{\varpi, \theta^*}\|)$, i.e. $d_{E,n} := \dim_E(\{f|_{B_n}\}_{f \in \mathcal{F}_\Theta}, 2\sqrt{\varepsilon/n})$. In Appendix C.2, we derive a general result (Proposition 22) from which Proposition 5 follows.

Proposition 6 *Under Assumptions 1 and 2, for any $\delta \in (0, 1)$, $\alpha \in \mathcal{A}$, $x_0 \in \mathbb{R}^d$, and $t \in \mathbb{R}_+$, we have with probability at least $1 - \delta$*

$$\sum_{n=1}^{N_t} \left\| \mu_{\hat{\theta}_n}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) - \mu_{\theta^*}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) \right\| \leq \tilde{O} \left(\sqrt{\varepsilon d_{E,N_t} \log(\mathcal{N}_{N_t}^\varepsilon)} N_t + \varepsilon d_{E,N_t} \right), \quad (14)$$

and

$$\sum_{n=1}^{N_t} \left\| \mu_{\hat{\theta}_n}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) - \mu_{\theta^*}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) \right\|^2 \leq \tilde{O} \left(d_{E,N_t} \log(\mathcal{N}_{N_t}^\varepsilon) \right). \quad (15)$$

C.1. Confidence sets

In this section, we work in a generic online learning framework, so that our results can be more easily compared and contrasted with (Russo and Van Roy, 2013; Osband and

(Van Roy, 2014) and others. We, therefore, introduce some dedicated notation and a stand-alone assumption for this section.

Consider a set of functions \mathcal{F} from $\mathbb{R}^d \rightarrow \mathbb{R}^d$, and fix $f^* \in \mathcal{F}$. We will study pairs of (random) \mathbb{R}^d -valued sequences $((X_i)_{i \in \mathbb{N}}, (Y_i)_{i \in \mathbb{N}})$ generated as

$$Y_i = f^*(X_i) + \xi_i$$

for $(\xi_i)_{i \in \mathbb{N}}$ a stochastic process in some filtered probability space $(\Omega', \mathcal{H}_\infty, \mathbb{H}, \mathbb{P})$, with each ξ_i independent of everything else up to time i . We take \mathcal{H}_i as the completion of $\sigma(\{\xi_j\}_{j \leq i})$, for $i \in \mathbb{N}$, and we let $\mathbb{H} = (\mathcal{H}_i)_{i \geq 0}$.

Given some \mathbb{R}^d -valued and \mathbb{H} -adapted sequences $(Z_i)_{i \in \mathbb{N}}$ and $(Z'_i)_{i \in \mathbb{N}}$, and some $n \in \mathbb{N}^*$, let us define

$$\langle Z|Z' \rangle_n := \sum_{i=0}^{n-1} \langle Z_i|Z'_i \rangle \text{ and } \|Z\|_n := \sqrt{\langle Z|Z \rangle_n}.$$

While $\|\cdot\|_n$ is not a norm, it plays this role and we follow here the notational convention of (Russo and Van Roy, 2013). We will extend the definitions of $\langle \cdot | \cdot \rangle_n$ and $\|\cdot\|_n$ to $n = 0$ by simply taking the empty sum to be 0, i.e. $\langle Z, Z' \rangle_0 := 0$.

To simplify notation, we will drop the sequence $(X_i)_{i \in \mathbb{N}}$ when it is an argument to a function inside $\|\cdot\|_n$ or $\langle \cdot | \cdot \rangle_n$: i.e. $\|f\|_n$ stands for $\|(f(X_i))_{i \in \mathbb{N}}\|_n$. With this notation in mind, for any $n \in \mathbb{N}$, we define \hat{f}_n as an arbitrary element of

$$\operatorname{argmin}_{f \in \mathcal{F}} \|Y - f\|_n^2.$$

In other words \hat{f}_n is a non-linear least-square fit in \mathcal{F} using the first n points of $(X_i, Y_i)_{i \in \mathbb{N}}$. In this generic setting, we introduce Assumption 3, which in our end-goal application subsumes Assumptions 1 and 2 and Proposition 2.

Assumption 3 *There is $(L, \Gamma) \in \mathbb{R}_+^2$ and a function $H_\delta : \mathbb{N} \rightarrow \mathbb{R}_+$ such that*

$$\sup_{f \in \mathcal{F}} \sup_{x \in \mathbb{R}^d} \frac{\|f(x)\|}{1 + \|x\|} \leq L,$$

and for all $i \in \mathbb{N}^$, ξ_i is an \mathcal{H}_{i-1} -conditionally Γ^2 -sub-Gaussian random variable, ξ_0 is Γ^2 -sub-Gaussian, and the sequence $(X_i)_{i \in \mathbb{N}}$ satisfies*

$$\mathbb{P} \left(\sup_{n \in \mathbb{N}} \frac{\|X_n\|}{H_\delta(n)} > 1 \right) < \delta$$

for all $\delta \in (0, 1)$.

Let $(\mathcal{C}_n^\Gamma)_{n \in \mathbb{N}^*}$ denote a deterministic sequence of finite covers of \mathcal{F} , whose cardinalities are respectively given by $(\mathcal{N}_n^\Gamma)_{n \in \mathbb{N}^*}$, such that for all $n \in \mathbb{N}^*$

$$\sup_{f \in \mathcal{F}} \min_{g \in \mathcal{C}_n^\Gamma} \sup_{x \in \mathcal{B}(H_\delta(n))} \|f(x) - g(x)\| \leq \frac{\Gamma^2}{n}.$$

The definition of this cover corresponds to one used in [Russo and Van Roy \(2013\)](#) with a domain restricted to lie in the high-probability region of the state process instead of the whole domain. This ensures the cover remains finite for all $n \in \mathbb{N}^*$.

For any $\delta \in (0, 1)$, $n \in \mathbb{N}^*$, and $f \in \mathcal{F}$ let us define the quantities

$$\begin{aligned} L_n^1(\delta) &:= \log((\Gamma^2 + 8nL^2(1 + \sup_{i \leq n} \|X_i\|_2^2)) \mathcal{A}_n^\Gamma \delta^{-1}), \\ L_n^0(\delta) &:= L_n^1(6\delta\pi^{-2}n^{-2}), \\ C_n^1(f) &:= \Gamma^2 + \|f - f^*\|_n^2 \\ C_n^2(f) &:= \sup_{i \leq n} \|f(X_i) - \hat{f}_n(X_i)\|, \end{aligned}$$

and the event

$$\begin{aligned} \mathcal{E}_n^0(\delta) &:= \left\{ \|\hat{f}_n - f^*\|_n \leq 2\Gamma \sqrt{L_n^1\left(\frac{3\delta}{\pi^2 n^2}\right)} \right. \\ &\quad \left. + 2\sqrt{\Gamma^2 + 2\Gamma \left(n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g) \sqrt{2 \log\left(\frac{4\pi^2 n^3}{3\delta}\right)} + \sqrt{2n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g) L_n^1\left(\frac{3\delta}{\pi^2 n^2}\right)} \right)} \right\}. \end{aligned} \quad (39)$$

Building upon the proof method of [Russo and Van Roy \(2013\)](#), the cornerstone of this section is Lemma 17, which shows that, with high-probability, f^* is contained in all the elements of a sequence of confidence sets, each centred at \hat{f}_n in the $\|\cdot\|_n$ norm.

Lemma 17 *Under Assumption 3, for $n \in \mathbb{N}^*$ and $\delta \in (0, 1)$, we have*

$$\mathbb{P} \left(\bigcap_{n \in \mathbb{N}^*} \mathcal{E}_n^0(\delta) \right) \geq 1 - \delta.$$

We begin the proof of Lemma 17 by giving the concentration inequality of Lemma 16.

Lemma 16 *Under Assumption 3, for all $n \in \mathbb{N}^*$, $\delta \in (0, 1)$, and $f \in \mathcal{F}$*

$$\mathbb{P} \left(|\langle \xi | f - f^* \rangle_n| \geq \Gamma \sqrt{2(\Gamma^2 + \|f - f^*\|_n) \log\left(\frac{\Gamma^2 + \|f - f^*\|_n}{\delta}\right)} \right) \leq \delta.$$

Proof This proof relies on extensively studied arguments for self-normalised inequalities, but we include it for completeness because it uses non-standard constants. Let us begin by fixing $f \in \mathcal{F}$. For all $n \in \mathbb{N}$, let

$$Z_n(f) := \langle \xi | f - f^* \rangle_n.$$

For any $\lambda \in \mathbb{R}$, let us define the process $(M_n^\lambda(f))_{n \in \mathbb{N}}$ defined by

$$M_n^\lambda(f) := \exp \left(\lambda Z_n(f) - \frac{\lambda^2 \Gamma^2}{2} \|f - f^*\|_n^2 \right).$$

Let us show that $M_n^\lambda(f)$ is a conditional supermartingale. For any $n \in \mathbb{N}$, we have

$$\mathbb{E} \left[M_{n+1}^\lambda(f) | \mathcal{H}_n \right] = M_n^\lambda(f) \mathbb{E} \left[\exp \left(\lambda \langle \xi_{n+1} | f(X_n) - f^*(X_n) \rangle_n \right) \middle| \mathcal{H}_n \right] e^{-\frac{\lambda^2 \Gamma^2}{2} \|f(X_n) - f^*(X_n)\|_n^2}. \quad (40)$$

By the Cauchy-Schwartz inequality

$$|\langle \xi_n | f(X_n) - f^*(X_n) \rangle_n| \leq \|\xi_n\|_n \|f(X_n) - f^*(X_n)\|_n$$

and thus, since ξ_n is conditionally Γ^2 -subgaussian with variance Γ^2 , $\|\xi_n\|$ is Γ^2 -subgaussian. Therefore

$$\mathbb{E} \left[\exp \left(\lambda \langle \xi_n | f(X_n) - f^*(X_n) \rangle_n - \frac{\lambda^2 \Gamma^2}{2} \|f(X_n) - f^*(X_n)\|_n^2 \right) \middle| \mathcal{H}_n \right] \leq 1$$

and thus, by (40), $M_n^\lambda(f)$ is a supermartingale. By definition of $\langle \cdot | \cdot \rangle_0$ and $\|\cdot\|_0$, $M_0^\lambda(f) = 1$, so that $\mathbb{E}[M_n^\lambda(f)] \leq 1$ for all $n \in \mathbb{N}$.

We now perform a Laplace trick. Let Φ be the Gaussian measure of mean 0 and variance Γ^{-4} on \mathbb{R} , and let us define the process $(M_n(f))_{n \in \mathbb{N}}$ by

$$\begin{aligned} M_n(f) &:= \int M_n^\lambda(f) \Phi(d\lambda) \\ &= \int \exp \left(\lambda Z_n(f) - \frac{\lambda^2 \Gamma^2}{2} \|f - f^*\|_n^2 \right) \Phi(d\lambda) \\ &= \frac{1}{\Gamma^2 + \|f - f^*\|_n^2} \exp \left\{ \frac{Z_n^2(f)}{2\Gamma^2(\Gamma^2 + \|f - f^*\|_n^2)} \right\}. \end{aligned}$$

By Markov's inequality, $\mathbb{P}(M_n(f) \geq \delta^{-1}) \leq \delta$, and thus

$$\mathbb{P} \left(Z_n(f) \geq \Gamma \sqrt{2(\Gamma^2 + \|f - f^*\|_n^2) \log \left(\frac{\Gamma^2 + \|f - f^*\|_n^2}{\delta} \right)} \right) \leq \delta. \quad \blacksquare$$

We will turn to the proof of Lemma 17. Recall (39), which defined for $\delta \in (0, 1)$ and $n \in \mathbb{N}^*$, the event

$$\begin{aligned} \mathcal{E}_n^0(\delta) &:= \left\{ \|\hat{f}_n - f^*\|_n \leq 2\Gamma \sqrt{L_n^1 \left(\frac{3\delta}{\pi^2 n^2} \right)} \right. \\ &\quad \left. + 2\sqrt{\Gamma^2 + 2\Gamma \left(n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g) \sqrt{2 \log \left(\frac{4\pi^2 n^3}{3\delta} \right)} + \sqrt{2n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g) L_n^1 \left(\frac{3\delta}{\pi^2 n^2} \right)} \right)} \right\}. \end{aligned}$$

Lemma 17 *Under Assumption 3, for $n \in \mathbb{N}^*$ and $\delta \in (0, 1)$, we have*

$$\mathbb{P} \left(\bigcap_{n \in \mathbb{N}^*} \mathcal{E}_n^0(\delta) \right) \geq 1 - \delta.$$

Proof The proof builds on elements of [Russo and Van Roy \(2013\)](#). We begin by giving two small auxiliary results which we will use.

- i. Let $n \in \mathbb{N}^*$, and $\delta \in (0, 1)$, by a union bound over the family of conditionally sub-Gaussian random variables $(\|\xi_i\|)_{i \in [n]}$, we have

$$\mathbb{P} \left(\sup_{i \leq n} \|\xi_i\| \leq \Gamma \sqrt{2 \log \left(\frac{2n}{\delta} \right)} \right) \geq 1 - \delta \quad (41)$$

- ii. For any $f \in \mathcal{F}$, and $n \in \mathbb{N}^*$ we have

$$\begin{aligned} \|f^* - Y\|_n^2 - \|f - Y\|_n^2 &= \langle f^* - Y | f^* - Y \rangle_n - \langle f - f^* + f^* - Y | f - f^* + f^* - Y \rangle_n \\ &= \langle f^* - Y | f^* - Y \rangle_n - \langle f - f^* | f - f^* \rangle_n \\ &\quad + 2 \langle Y - f^* | f - f^* \rangle_n - \langle Y - f^* | Y - f^* \rangle_n \\ &= -\|f - f^*\|_n^2 + 2 \langle \xi | f - f^* \rangle_n. \end{aligned} \quad (42)$$

Applying (42) with $f := \hat{f}_n$, the n -point non-linear least-square fit, leads to a non-negative left-hand side and thus

$$\left\| \hat{f}_n - f^* \right\|_n^2 \leq 2 |\langle \xi | f - f^* \rangle_n|.$$

At the same time, for all $n \in \mathbb{N}^*$, by definition of \mathcal{C}_n^Γ , it holds that for all $g \in \mathcal{C}_n^\Gamma$

$$\begin{aligned} \left\| \hat{f}_n - f^* \right\|_n^2 &\leq 2 |\langle \xi | g - f^* \rangle_n| + 2 \left| \langle \xi | \hat{f}_n - g \rangle_n \right| \\ &\leq 2 |\langle \xi | g - f^* \rangle_n| + 2n \sup_{i \leq n} \|\xi_i\|_2 C_n^2(g). \end{aligned} \quad (43)$$

Combining (41) and (43), we obtain, for all $\delta \in (0, 1)$, $n \in \mathbb{N}^*$, and $g \in \mathcal{C}_n^\Gamma$, that

$$\mathbb{P} \left(\left\| \hat{f}_n - f^* \right\|_n^2 \geq 2 |\langle \xi | g - f^* \rangle_n| + 2n C_n^2(g) \Gamma \sqrt{2 \log \left(\frac{2n}{\delta} \right)} \right) \leq \delta \quad (44)$$

Let us now provide two bounds on $C_n^1(g)$ we will use. For all $n \in \mathbb{N}^*$, $\delta \in (0, 1)$ and $g \in \mathcal{C}_n^\Gamma$, let

$$C_n^1(g) \leq \Gamma^2 + 8nL^2(1 + \sup_{i \leq n} \|X_i\|^2). \quad (45)$$

$$C_n^1(g) \leq \Gamma^2 + \left\| \hat{f}_n - f^* \right\|_n^2 + \left\| g - \hat{f}_n \right\|_n^2 \leq C_n^1(\hat{f}_n) + n C_n^2(g), \quad (46)$$

Applying Lemma 16 for each $g \in \mathcal{C}_n^\Gamma$, by a union bound over $g \in \mathcal{C}_n^\Gamma$, we have for any $\delta_0(n) \in (0, 1)$ (to be fixed at the end), that

$$\delta_0(n) \geq \mathbb{P} \left(\sup_{g \in \mathcal{C}_n^\Gamma} |\langle \xi | g - f^* \rangle_n| \geq \Gamma \sqrt{2 \sup_{g \in \mathcal{C}_n^\Gamma} C_n^1(g) \log \left(\frac{\sup_{g \in \mathcal{C}_n^\Gamma} C_n^1(g) \mathcal{N}_n^\Gamma}{\delta_0(n)} \right)} \right).$$

Applying (45) and (46) this becomes

$$\begin{aligned} \delta_0(n) &\geq \mathbb{P} \left(\sup_{g \in \mathcal{C}_n^\Gamma} |\langle \xi | g - f^* \rangle_n| \right. \\ &\quad \left. \geq \Gamma \sqrt{2(C_n^1(\hat{f}_n) + n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g)) \log \left(\frac{(\Gamma^2 + 8nL^2(1 + \sup_{i \leq n} \|X_i\|^2)) \mathcal{N}_n^\Gamma}{\delta_0(n)} \right)} \right) \end{aligned}$$

and thus

$$\delta_0(n) \geq \mathbb{P} \left(\sup_{g \in \mathcal{C}_n^\Gamma} |\langle \xi | g - f^* \rangle_n| \geq \Gamma \sqrt{2L_n^1(\delta_0(n))} \left(\sqrt{C_n^1(\hat{f}_n)} + \sqrt{n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g)} \right) \right). \quad (47)$$

Combining (44) and (47) by a union bound gives us

$$\begin{aligned} \delta_0(n) &\geq \mathbb{P} \left(\left\| \hat{f}_n - f^* \right\|_n^2 \geq 2\Gamma \sqrt{2L_n^1 \left(\frac{\delta_0(n)}{2} \right)} \left(\sqrt{C_n^1(\hat{f}_n)} + \sqrt{n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g)} \right) \right. \\ &\quad \left. + 2nC_n^2(g)\Gamma \sqrt{2 \log \left(\frac{4n}{\delta_0(n)} \right)} \right). \end{aligned}$$

For all $n \in \mathbb{N}^*$, on the complement of this event (whose probability is at least $1 - \delta_0(n)$), we have

$$C_n^1(\hat{f}_n) \leq \Gamma^2 + \Gamma \sqrt{2C_n^1(\hat{f}_n)L_n^1(\delta_0(n)/2)} + h_n^\Gamma, \quad (48)$$

in which

$$h_n^\Gamma := 2\Gamma \left(n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g) \sqrt{2 \log \left(\frac{4n}{\delta_0(n)} \right)} + \sqrt{2n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g) L_n^1 \left(\frac{\delta_0(n)}{2} \right)} \right).$$

Viewing (48) as a second order polynomial in $\sqrt{C_n^1(\hat{f}_n)}$, we obtain via its roots that

$$\begin{aligned} \sqrt{C_n^1(\hat{f}_n)} &\leq \Gamma \sqrt{L_n^1(\delta_0(n)/2)} + \sqrt{\left(\Gamma \sqrt{L_n^1(\delta_0(n)/2)} \right)^2 + 4(\Gamma^2 + h_n^\Gamma)} \\ &\leq 2\Gamma \sqrt{L_n^1(\delta_0(n)/2)} + 2\sqrt{\Gamma^2 + h_n^\Gamma}. \end{aligned}$$

Since $\|\hat{f}_n - f^*\|_n \leq \sqrt{C_n^1(\hat{f}_n)}$ by definition of $C_n^1(\hat{f}_n)$, we have

$$\|\hat{f}_n - f^*\|_n \leq 2\sqrt{\Gamma^2 + 2\Gamma \left(n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g) \sqrt{2 \log \left(\frac{4n}{\delta_0(n)} \right)} + \sqrt{2n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g) L_n^1 \left(\frac{\delta_0(n)}{2} \right)} \right)} + 2\Gamma \sqrt{L_n^1(\delta_0(n)/2)}.$$

Therefore, letting

$$\mathcal{E}_n^1(\delta) := \left\{ \|\hat{f}_n - f^*\|_n \leq 2\Gamma \sqrt{L_n^1 \left(\frac{\delta}{2} \right)} + 2\sqrt{\Gamma^2 + 2\Gamma \left(n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g) \sqrt{2 \log \left(\frac{4n}{\delta} \right)} + \sqrt{2n \sup_{g \in \mathcal{C}_n^\Gamma} C_n^2(g) L_n^1 \left(\frac{\delta}{2} \right)} \right)} \right\},$$

we have, for all $n \in \mathbb{N}^*$, that $\mathbb{P}(\mathcal{E}_n^1(\delta_0(n))) \geq \delta_0(n)$. Letting $\delta_0(n) = \frac{6}{\pi^2 n^2} \delta$, by a union bound we obtain

$$\mathbb{P} \left(\bigcap_{n \in \mathbb{N}^*} \mathcal{E}_n^1(\delta_0(n)) \right) \geq 1 - \delta \frac{6}{\pi^2} \sum_{n=1}^{\infty} \frac{1}{n^2} = 1 - \delta.$$

Noting that $\mathcal{E}_n^0(\delta) = \mathcal{E}_n^1(\delta_0(n))$ for all $\delta \in (0, 1)$ and $n \in \mathbb{N}^*$ completes the proof. \blacksquare

In the proof of Lemma 17, we used self-normalised inequalities to generalise the results of Russo and Van Roy (2013) to unbounded states. We now incorporate the high probability bound of Assumption 3 and formalise confidence sets, which will prove Theorem 18. Theorem 18 can then be specified for our setting by merging it with the results of Appendix B in Proposition 5.

For $\delta \in (0, 1)$, let $\beta_0 \in \mathbb{R}_+$ and let us define the sequence $(\mathcal{C}_n(\delta))_{n \in \mathbb{N}}$ in which

$$\mathcal{C}_n(\delta) := \left\{ f \in \mathcal{F} : \|f - \hat{f}_n\|_n \leq \beta_n \right\} \quad (49)$$

with

$$\beta_n(\delta) := \beta_0 \vee 2\Gamma \left(\sqrt{1 + 2 \left(\sqrt{2\Gamma \log \left(\frac{8n}{\delta} \right)} + \sqrt{2L_n^0 \left(\frac{\delta}{4} \right)} \right)} + \sqrt{L_n^0 \left(\frac{\delta}{4} \right)} \right). \quad (50)$$

Theorem 18 *Under Assumption 3, we have for all $\delta \in (0, 1)$*

$$\mathbb{P} \left(\left\{ \bigcap_{n \in \mathbb{N}^*} \{f^* \in \mathcal{C}_n(\delta)\} \right\} \cap \left\{ \sup_{n \in \mathbb{N}^*} \frac{\|X_n\|}{H_\delta(n)} \leq 1 \right\} \right) \leq \delta$$

Proof Fix $\delta \in (0, 1)$, and assume $\omega \in \{\omega' \in \Omega : \sup_{n \in \mathbb{N}^*} \|X_n(\omega')\|_2 / H_\delta(n) \leq 1\}$. In this case, we have the following bound, for all $n \in \mathbb{N}^*$

$$2n \min_{g \in \mathcal{C}_n^\Gamma} C_n^2(g) \leq 2\Gamma^2$$

by definition of \mathcal{C}_n^Γ as a $\Gamma^2 n^{-1}$ cover on $\mathcal{B}_2(H_\delta(n))$. Therefore, the event

$$\left\{ \bigcap_{n \in \mathbb{N}^*} \mathcal{E}_n^0(\delta) \right\} \cap \left\{ \sup_{n \in \mathbb{N}^*} \frac{\|X_n\|}{H_\delta(n)} \leq 1 \right\}$$

is contained in the event

$$\mathcal{E}^0(\delta) := \left\{ \bigcap_{n \in \mathbb{N}^*} \left\{ \|f^* - \hat{f}_n\|_n \leq \beta_n(2\delta) \right\} \right\} \cap \left\{ \sup_{n \in \mathbb{N}^*} \frac{\|X_n\|}{H_\delta(n)} \leq 1 \right\}.$$

By Lemma 17, Assumption 3, and a union bound, $\mathbb{P}(\mathcal{E}^0(\delta)) \geq 1 - 2\delta$, and we obtain the result by (49) and (50), i.e. by definition of $\mathcal{C}_n(\delta)$. \blacksquare

Proposition 5 (Adapted from Osband and Van Roy (2014, Prop. 5)) *Under Assumptions 1 and 2, for any $x_0 \in \mathbb{R}^d$, and $\delta > 0$,*

$$\mathbb{P} \left(\left\{ \theta^* \in \bigcap_{n=1}^{\infty} \mathcal{C}_n(\delta) \right\} \cap \left\{ \sup_{n \in \mathbb{N}^*} \frac{\|X_{\tau_n}^{\varpi, \theta^*}\|}{H_\delta(n)} \leq 1 \right\} \right) \geq 1 - \delta, \quad (13)$$

Proof The proof follows by applying Theorem 18 to this setting. Where $(X_i)_{i \in \mathbb{N}} := ((X_{\tau_i}^{\varpi, \theta^*}, \varpi_{\tau_i}))_{i \in \mathbb{N}}$, $(Y_i)_{i \in \mathbb{N}} := (X_{\tau_{i+1}}^{\varpi, \theta^*} - X_{\tau_i}^{\varpi, \theta^*})_{i \in \mathbb{N}}$, $\mathcal{F} := \mathcal{F}_\Theta$ and with $(\xi_{n+1})_{n \in \mathbb{N}}$ and $(\beta_n(\delta))_{n \in \mathbb{N}^*}$ as defined in Section 2 and (12) respectively. This sets $\Gamma = \|\Sigma\|_{\text{op}} = \varepsilon^{\frac{1}{2}} \|\bar{\Sigma}\|_{\text{op}}$. The only subtlety is that the process X^{ϖ, θ^*} is measured at random times, but since these times are independent of anything else, and the process is almost surely constant between them, they do not affect the proof. \blacksquare

C.2. Widths of confidence sets

In Appendix C.1, we showed how to design confidence sets along a trajectory of $X^{\alpha, \theta}$ for learning μ by using NLLS to minimise a fit error of the form

$$\sum_{n=1}^N \left\| \mu_1(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) - \mu_2(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) \right\|,$$

for $(\mu_1, \mu_2) \in \mathcal{C}_N(\delta)$ and $N \in \mathbb{N}^*$. When analysing the regret of such a learning algorithm this is not sufficient: instead of the fit error, we need to control a prediction error of the form

$$\sum_{n=1}^N \left\| \mu_{\theta_n}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) - \mu_{\theta^*}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) \right\|,$$

for $(\mu_{\theta_n})_{n \in \mathbb{N}} \subset \mathcal{F}_\theta$ such that $\mu_{\theta_n} \in \mathcal{C}_n(\delta)$ for all $n \in \mathbb{N}$. The difference is that μ_{θ_n} changes over time so that the sum counts the errors in predicting the next state made by the sequence $(\mu_{\theta_n})_{n \in \mathbb{N}}$.

Since we will want to implement lazy updates, we will need a more general result where the μ_{θ_n} are not all in their respective $\mathcal{C}_n(\delta)$ but rather are from a piece-wise constant sequence with $\mu_{\theta_n} := \mu_{\theta_{k(n)}} \in \mathcal{C}_{k(n)}(\delta)$, where $k(n) \leq n$ for all $n \in \mathbb{N}$. Therefore, as in Appendix C.1, we begin by showing a general result in the learning framework of [Russo and Van Roy \(2013\)](#) (Proposition 19), then apply it to our setting to prove Proposition 6. Using the notation of Appendix C.1, let \mathcal{F} be a function class of functions from $\mathbb{R}^d \rightarrow \mathbb{R}^d$, and recall the arbitrary sequence $(X_n)_{n \in \mathbb{N}} \subset \mathbb{R}^d$.

The Ξ -eluder dimension of a function class \mathcal{F} , for $\Xi \in \mathbb{R}_+$, introduced in [Russo and Van Roy \(2013\)](#) is designed for converting fit errors into prediction errors. Unlike [Russo and Van Roy \(2013\)](#), we must adapt our eluder dimension to work with unbounded functions on unbounded processes. Failing to do so would lead our results to be largely vacuous since the eluder dimension of \mathcal{F} might be infinite for any Ξ .

We work with a modified eluder dimension, which takes three arguments: a function class \mathcal{F} whose elements have for domain a set $\mathcal{X} \subset \mathbb{R}^d$; a set $S \subset \mathcal{X}$; and $\Xi \in \mathbb{R}_+$. Our modified eluder dimension is the Ξ -eluder dimension of $\{f|_S : f \in \mathcal{F}\}$, the class containing the restrictions to S of elements of \mathcal{F} , which we denote by $\dim_{\Xi}^S(\mathcal{F}, \Xi)$. In this way, the eluder dimension of [Russo and Van Roy \(2013\)](#) is $\dim_{\Xi}^{\mathcal{X}}(\mathcal{F}, \Xi)$. For $n \in \mathbb{N}^*$, let $B_n := \mathcal{B}_2(\sup_{i \in [n]} \|X_i\|)$ and, for any $u \in \mathbb{R}_+$, let us define the sequence $(d_{\mathbb{E},n}^{\mathcal{F}}(u))_{n \in \mathbb{N}^*}$, in which

$$d_{\mathbb{E},n}^{\mathcal{F}}(u) := \dim_{\mathbb{E}}^{B_n} \left(\mathcal{F}, \frac{2u}{\sqrt{n}} \right)$$

for all $n \in \mathbb{N}^*$ and $u \in \mathbb{R}_+$.

Proposition 19 *Let $(\tilde{\beta}_i)_{i \in \mathbb{N}}$ be a non-decreasing positive real-valued sequence, $(\tilde{f}_i)_{i \in \mathbb{N}}$, and $(\mathcal{F}_i)_{i \in \mathbb{N}}$ be a sequence of subsets of \mathcal{F} of the form $\mathcal{F}_i := \{f \in \mathcal{F} : \|f - \tilde{f}_i\|_i \leq \tilde{\beta}_i\}$. Under Assumption 3, for any $n \in \mathbb{N}^*$, we have*

$$\sum_{i=1}^n \sup_{(f,f') \in \mathcal{F}_n^2} \|f(X_i) - f'(X_i)\| \leq 2\tilde{\beta}_n \sqrt{d_{\mathbb{E},n}^{\mathcal{F}}(\tilde{\beta}_0)n} + d_{\mathbb{E},n}^{\mathcal{F}}(\tilde{\beta}_0)2L(1 + \sup_{i \in [n]} \|X_i\|), \quad (51)$$

and

$$\begin{aligned} \sum_{i=1}^n \sup_{(f,f') \in \mathcal{F}_n^2} \|f(X_i) - f'(X_i)\|^2 &\leq 4\tilde{\beta}_n^2 d_{\mathbb{E},n}^{\mathcal{F}}(\tilde{\beta}_0) \left(3 + \log \left(\frac{n8L^2(1 + \sup_{i \in [n]} \|X_i\|)}{16\tilde{\beta}_n^4 (d_{\mathbb{E},n}^{\mathcal{F}}(\tilde{\beta}_0))^2} \right) \right) \\ &\quad + 2d_{\mathbb{E},n}^{\mathcal{F}}(1 + 2\tilde{\beta}_n^2 d_{\mathbb{E},n}^{\mathcal{F}}(\tilde{\beta}_0))(1 + 8L^2(1 + \sup_{i \in [n]} \|X_i\|)). \end{aligned} \quad (52)$$

To prove Proposition 19, the key result of [Russo and Van Roy \(2013\)](#) we leverage is Lemma 20 which we combine with two functional inequalities given in Lemma 21.

For a function class \mathcal{F} with domain $\mathcal{X} \subset \mathbb{R}^d$, and any $x \in \mathcal{X}$, let us define

$$\Lambda(\mathcal{F}; x) = \sup_{(f_1, f_2) \in \mathcal{F}^2} \|f_1(x) - f_2(x)\|.$$

The quantity $\Lambda(\mathcal{F}, x)$ is the maximal prediction gap at x between two functions in \mathcal{F} . Bounding the prediction error along $(X_i)_{i \in \mathbb{N}}$ of a sequence of function classes $(\mathcal{F}_i)_{i \in \mathbb{N}} \subset \mathcal{F}$ means bounding $\sum_{i=1}^n \Lambda(\mathcal{F}_i, X_i)$ in terms of $n \in \mathbb{N}$.

Lemma 20 [*Russo and Van Roy (2013, Prop.3)*] Let $(\tilde{f}_i)_{i \in \mathbb{N}}$ be a sequence of elements of \mathcal{F} , $(\mathcal{F}_i)_{i \in \mathbb{N}}$ be a sequence of subsets of \mathcal{F} of the form $\mathcal{F}_i := \{f \in \mathcal{F} : \|f - \tilde{f}_i\|_i \leq \tilde{\beta}_i\}$. For any $\Xi \in (0, 1)$ and $n \in \mathbb{N}$, one has

$$\sum_{i=1}^n \mathbb{1}_{\{\Lambda(\mathcal{F}_i; X_i) > \Xi\}} \leq \left(\frac{4\tilde{\beta}_n^2}{\Xi^2} + 1 \right) \dim_{\mathbb{E}^n}^{B^n}(\mathcal{F}, \Xi).$$

Proof Following the proof of (*Russo and Van Roy, 2013, Prop.3*), the only modification involves the bound $\|\bar{f} - \underline{f}\|_n \leq \tilde{\beta}_n$, for any $(\bar{f}, \underline{f}) \in \mathcal{F}_n^2$, which holds by assumption. \blacksquare

Lemma 21 Let $(x_i)_{i \in \mathbb{N}} \in \mathbb{R}_+^{\mathbb{N}^*}$. Assume there is a family of positive sequences $((\zeta_n^\vartheta)_{n \in \mathbb{N}})_{\vartheta \in \mathbb{R}_+}$ and a family of positive constants $(\chi^\vartheta)_{\vartheta \in \mathbb{R}_+}$ such that, for any $n \in \mathbb{N}^*$ and $\vartheta > 0$,

$$\sum_{i=1}^n \mathbb{1}_{\{x_i > \vartheta\}} \leq \frac{\zeta_n^\vartheta}{\vartheta^2} + \chi^\vartheta \quad (53)$$

then the following two inequalities hold

$$\sum_{i=1}^n x_i \leq 2\sqrt{n\zeta_n^\vartheta} + \chi^\vartheta \sup_{i \in [n]} x_i \quad (54)$$

$$\sum_{i=1}^n x_i^2 \leq \zeta_n^\vartheta \left(3 + \log \left(\frac{n \sup_{i \in [n]} x_i^2}{(\zeta_n^\vartheta)^2} \right) \right) + \chi^\vartheta (2 + \zeta_n^\vartheta) (1 + \sup_{i \in [n]} x_i^2). \quad (55)$$

Proof

i. For $\vartheta > 0$, we have by (53)

$$\begin{aligned} \sum_{i=1}^n (x_i - \vartheta) \mathbb{1}_{\{x_i > \vartheta\}} &= \sum_{i=1}^n \int_{\vartheta}^{x_i} \mathbb{1}_{\{x_i > u\}} du \\ &\leq \int_{\vartheta}^{\sup_{i \in [n]} x_i} \sum_{i=1}^n \mathbb{1}_{\{x_i > u\}} du \\ &\leq \int_{\vartheta}^{\sup_{i \in [n]} x_i} \frac{\zeta_n^\vartheta}{u^2} + \chi^\vartheta du \\ &= \chi \sup_{i \in [n]} x_i - \frac{\zeta_n^\vartheta}{\sup_{i \in [n]} x_i} - \chi^\vartheta \vartheta + \frac{\zeta_n^\vartheta}{\vartheta}, \end{aligned}$$

and thus

$$\sum_{i=1}^n (x_i - \vartheta) \mathbb{1}_{\{x_i > \vartheta\}} \leq \frac{\zeta_n^\vartheta}{\vartheta} + \chi^\vartheta \sup_{i \in [n]} x_i. \quad (56)$$

Combining (56) with

$$\sum_{i=1}^n (x_i - \vartheta) \leq \sum_{i=1}^n (x_i - \vartheta) \mathbb{1}_{\{x_i > \vartheta\}}$$

yields

$$\sum_{i=1}^n x_i \leq n\vartheta + \frac{\zeta_n^\vartheta}{\vartheta} + \chi^\vartheta \sup_{i \in [n]} x_i.$$

Setting $\vartheta = \sqrt{\zeta_n^\vartheta/n}$ yields (54).

ii. To prove (55), we iterate the bound (56)

$$\begin{aligned} \sum_{i=1}^n (x_i - \vartheta)^2 \mathbf{1}_{\{x_i > \vartheta\}} &= 2 \sum_{i=1}^n \int_{\varepsilon}^{x_i} (x_i - u) \mathbf{1}_{\{x_i > u\}} du \\ &\leq 2 \sum_{i=1}^n \int_{\vartheta}^{\sup_{i \in [n]} x_i} (x_i - u) \mathbf{1}_{\{x_i > u\}} du \\ &\leq 2 \int_{\vartheta}^{\sup_{i \in [n]} x_i} \frac{\zeta_n^\vartheta}{\vartheta} + \chi^\vartheta \sup_{i \in [n]} x_i du \\ &\leq 2 \left(\chi \left(\sup_{i \in [n]} x_i^2 - \sup_{i \in [n]} x_i \vartheta \right) + \zeta_n^\vartheta \log \left(\frac{\sup_{i \in [n]} x_i}{\vartheta} \right) \right) \\ &\leq 2 \zeta_n^\vartheta \log \left(\frac{\sup_{i \in [n]} x_i}{\vartheta} \right) + 2 \chi^\vartheta \sup_{i \in [n]} x_i^2. \end{aligned}$$

Now, by some algebraic manipulations of $\sum_{i=1}^n x_i^2$, completing the square, discarding negative terms, and using (56) in the third step, we get

$$\begin{aligned} \sum_{i=1}^n x_i^2 &\leq \sum_{i=1}^n x_i^2 \mathbf{1}_{\{x_i > \vartheta\}} + \vartheta^2 \sum_{i=1}^n \mathbf{1}_{\{x_i > \vartheta\}} \\ &\leq \sum_{i=1}^n (x_i - \vartheta)^2 \mathbf{1}_{\{x_i > \vartheta\}} + 2\vartheta \sum_{i=1}^n x_i \mathbf{1}_{\{x_i > \vartheta\}} + n\vartheta^2 \\ &\leq 2 \zeta_n^\vartheta \log \left(\frac{\sup_{i \in [n]} x_i}{\vartheta} \right) + 2 \chi^\vartheta \sup_{i \in [n]} x_i^2 + \vartheta \left(\frac{\zeta_n^\vartheta}{\vartheta} + \chi^\vartheta \sup_{i \in [n]} x_i + \vartheta n \right) + n\vartheta^2. \end{aligned}$$

Taking $\vartheta = \zeta_n^\vartheta/\sqrt{n}$ and factoring, using also $u \leq 1 + u^2$ for $u \in \mathbb{R}$, yields

$$\sum_{i=1}^n x_i^2 \leq \zeta_n^\vartheta \left(3 + \log \left(\frac{n \sup_{i \in [n]} x_i^2}{(\zeta_n^\vartheta)^2} \right) \right) + \chi^\vartheta (2 + \zeta_n^\vartheta) (1 + \sup_{i \in [n]} x_i^2).$$

■

Proposition 19 *Let $(\tilde{\beta}_i)_{i \in \mathbb{N}}$ be a non-decreasing positive real-valued sequence, $(\tilde{f}_i)_{i \in \mathbb{N}}$, and $(\mathcal{F}_i)_{i \in \mathbb{N}}$ be a sequence of subsets of \mathcal{F} of the form $\mathcal{F}_i := \{f \in \mathcal{F} : \|f - \tilde{f}_i\|_i \leq \tilde{\beta}_i\}$. Under Assumption 3, for any $n \in \mathbb{N}^*$, we have*

$$\sum_{i=1}^n \sup_{(f, f') \in \mathcal{F}_n^2} \|f(X_i) - f'(X_i)\| \leq 2\tilde{\beta}_n \sqrt{d_{\mathcal{E}, n}^{\mathcal{F}}(\tilde{\beta}_0) n} + d_{\mathcal{E}, n}^{\mathcal{F}}(\tilde{\beta}_0) 2L(1 + \sup_{i \in [n]} \|X_i\|), \quad (51)$$

and

$$\begin{aligned} \sum_{i=1}^n \sup_{(f,f') \in \mathcal{F}_n^2} \|f(X_i) - f'(X_i)\|^2 &\leq 4\tilde{\beta}_n^2 d_{E,n}^{\mathcal{F}}(\tilde{\beta}_0) \left(3 + \log \left(\frac{n8L^2(1 + \sup_{i \in [n]} \|X_i\|)}{16\tilde{\beta}_n^4 (d_{E,n}^{\mathcal{F}}(\tilde{\beta}_0))^2} \right) \right) \\ &\quad + 2d_{E,n}^{\mathcal{F}}(1 + 2\tilde{\beta}_n^2 d_{E,n}^{\mathcal{F}}(\tilde{\beta}_0))(1 + 8L^2(1 + \sup_{i \in [n]} \|X_i\|^2)). \end{aligned} \quad (52)$$

Proof The proof consists in applying Lemma 21 to Lemma 20, with $x_i = \Lambda(\mathcal{F}_i, X_i)$, $\zeta_n^\vartheta = 4\tilde{\beta}_n^2 \dim_E^{B_n}(\mathcal{F}, \Xi)$ ($B_n := \mathcal{B}_2(\sup_{i \in [n]} \|X_i\|)$), and $\chi^\vartheta = \dim_E^{B_n}(\mathcal{F}, \Xi)$. When we set the value of ϑ in the proof of Lemma 21, χ^ϑ becomes

$$\dim_E^{B_n} \left(\mathcal{F}, \sqrt{\frac{4\tilde{\beta}_n^2}{n}} \right) \leq \dim_E^{B_n} \left(\mathcal{F}, \sqrt{\frac{4\tilde{\beta}_0^2}{n}} \right)$$

as $(\tilde{\beta}_n)_{n \in \mathbb{N}}$ is non-decreasing and the eluder dimension is decreasing in its third argument. An analogous remark holds for ζ_n^ϑ . We can thus substitute $\zeta_n^\vartheta = 4\tilde{\beta}_n^2 d_{E,n}^{\mathcal{F}}(\tilde{\beta}_0)$ and $\chi^\vartheta = d_{E,n}^{\mathcal{F}}(\tilde{\beta}_0)$ in (54) and (55), which gives the result. \blacksquare

We now apply Proposition 19 to our setting. For $n \in \mathbb{N}^*$, let us recall the shorthand notation

$$d_{E,n} := \dim_E^{B_n} \left(\mathcal{F}_\Theta, 2\sqrt{\frac{\varepsilon}{n}} \right) \quad (57)$$

in which we extended the notation from $(X_i)_{i \in \mathbb{N}}$ to $X^{\alpha, \theta}$ in the self-evident manner.

Proposition 22 *Under Assumptions 1 and 2, for any $(\alpha, \theta) \in \mathcal{A} \times \Theta$ and $t \in \mathbb{R}_+$, any non-decreasing positive real-valued sequence $(\tilde{\beta}_n)_{n \in \mathbb{N}}$, any $(\tilde{\mu}_n)_{n \in \mathbb{N}} \subset \mathcal{F}_\Theta$, and any sequence $(\mathcal{F}_n)_{n \in \mathbb{N}}$ of subsets of \mathcal{F}_θ of the form*

$$\mathcal{F}_n = \left\{ \mu \in \mathcal{F}_\Theta : \sqrt{\sum_{i=0}^{n-1} \left\| \mu_n(X_{\tau_i}^{\alpha, \theta}, \alpha_{\tau_i}) - \tilde{\mu}_n(X_{\tau_i}^{\alpha, \theta}, \alpha_{\tau_i}) \right\|_2^2} \leq \tilde{\beta}_n \right\},$$

we have

$$\sum_{n=1}^{N_t} \sup_{(\mu_1, \mu_2) \in \mathcal{F}_n} \left\| \mu_1(X_{\tau_n}^{\alpha, \theta}, \alpha_{\tau_n}) - \mu_2(X_{\tau_n}^{\alpha, \theta}, \alpha_{\tau_n}) \right\| \leq 2\beta_{N_t} \sqrt{d_{E, N_t}} + d_{E, N_t} 2\varepsilon L_0 (1 + \sup_{s \leq t} \|X_s^{\alpha, \theta}\|), \quad (58)$$

and

$$\begin{aligned} &\sum_{n=1}^{N_t} \sup_{(\mu_1, \mu_2) \in \mathcal{F}_n} \left\| \mu_1(X_{\tau_n}^{\alpha, \theta}, \alpha_{\tau_n}) - \mu_2(X_{\tau_n}^{\alpha, \theta}, \alpha_{\tau_n}) \right\|^2 \\ &\leq 4\beta_{N_t}^2 d_{E, N_t} \left(3 + \log \left(\frac{N_t 8\varepsilon^2 L_0^2 (1 + \sup_{s \leq t} \|X_s^{\alpha, \theta}\|)}{16\beta_{N_t}^4 d_{E, N_t}^2} \right) \right) \\ &\quad + 2d_{E, N_t} (1 + 2\beta_{N_t}^2 d_{E, N_t}) \left(1 + 8\varepsilon^2 L_0^2 (1 + \sup_{s \leq t} \|X_s^{\alpha, \theta}\|^2) \right). \end{aligned} \quad (59)$$

Proof Immediate by applying Proposition 19 to our setting, as we did in the proof of Proposition 5. \blacksquare

Under the event of Proposition 5, which ensures that $\theta^* \in \cap_{n \in \mathbb{N}} \mathcal{C}_n(\delta)$, we can derive from Proposition 22 a bound on the prediction error relative to the true dynamics X^{α, θ^*} generated by the control $\alpha \in \mathcal{A}$, in particular, we are interested in $\alpha = \varpi$.

Proposition 6 *Under Assumptions 1 and 2, for any $\delta \in (0, 1)$, $\alpha \in \mathcal{A}$, $x_0 \in \mathbb{R}^d$, and $t \in \mathbb{R}_+$, we have with probability at least $1 - \delta$*

$$\sum_{n=1}^{N_t} \left\| \mu_{\hat{\theta}_n}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) - \mu_{\theta^*}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) \right\| \leq \tilde{\mathcal{O}} \left(\sqrt{\varepsilon d_{\mathbb{E}, N_t} \log(\mathcal{N}_{N_t}^\varepsilon) N_t} + \varepsilon d_{\mathbb{E}, N_t} \right), \quad (14)$$

and

$$\sum_{n=1}^{N_t} \left\| \mu_{\hat{\theta}_n}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) - \mu_{\theta^*}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) \right\|^2 \leq \tilde{\mathcal{O}} \left(d_{\mathbb{E}, N_t} \log(\mathcal{N}_{N_t}^\varepsilon) \right). \quad (15)$$

Proof This follows from Proposition 22 by choosing $(\tilde{\beta}_n)_{n \in \mathbb{N}} = (\beta_n(\delta))_{n \in \mathbb{N}}$ and $(\tilde{\mathcal{F}}_n)_{n \in \mathbb{N}} = (\mathcal{C}_n(\delta))_{n \in \mathbb{N}}$, i.e. choosing $(\tilde{\mu}_n)_{n \in \mathbb{N}} = (\mu_{\hat{\theta}_n})_{n \in \mathbb{N}}$, the NLLS fit on n points. It is key to notice that these choices of $(\tilde{\beta}_n)_{n \in \mathbb{N}}$, $(\tilde{\mathcal{F}}_n)_{n \in \mathbb{N}}$, and $(\tilde{\mu}_n)_{n \in \mathbb{N}}$ are adapted to \mathbb{F} , and therefore we can apply Proposition 22 on the event of Proposition 5 without issues. This yields

$$\sum_{n=1}^{N_t} \left\| \mu_{\hat{\theta}_n}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) - \mu_{\theta^*}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) \right\| \leq 2\beta_{N_t}(\delta) \sqrt{d_{\mathbb{E}, N_t}} + 2\varepsilon L_0 d_{\mathbb{E}, N_t} (1 + H_\delta(N_T)),$$

and

$$\begin{aligned} \sum_{n=1}^{N_t} \left\| \mu_{\hat{\theta}_n}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) - \mu_{\theta^*}(X_{\tau_n}^{\alpha, \theta^*}, \alpha_{\tau_n}) \right\|^2 \\ \leq 4\beta_{N_T}(\delta)^2 d_{\mathbb{E}, N_t} \left(3 + \log \left(\frac{8\varepsilon^2 L_0^2 N_t (1 + H_\delta(N_T))}{16\beta_{N_t}(\delta)^4 d_{\mathbb{E}, N_t}^2} \right) \right) \\ + 2d_{\mathbb{E}, N_t} (1 + 2\beta_{N_t}(\delta)^2 d_{\mathbb{E}, N_t}) (1 + 8\varepsilon^2 L_0^2 (1 + H_\delta(N_t))^2). \end{aligned}$$

To obtain the estimates of (14)–(15), it suffices to recall the definitions of $\beta_n(\delta)$ (i.e. (12)) and $H_\delta(n)$ (i.e. (32)). \blacksquare

Appendix D. Planning and Diffusive Limit Approximation

Our work builds upon (Abeille et al., 2022), but with specialised results for our setting. This paper recovers the key results of this section (Propositions 7 to 9) under a stronger and more abstract set of assumptions. For the comfort of the reader, we present the necessary steps to extend their results to our assumptions. Since our assumptions do not directly subsume theirs, we exhibit in each case from Assumptions 1 and 2 how to recover the keystone results which underpin the technical arguments of (Abeille et al., 2022).

We begin with the well-posedness results for the pure jump case (Proposition 7) and the diffusive limit case (Proposition 8) and then focus on the approximation result linking the two regimes (Proposition 9). In Abeille et al. (2022), Proposition 7 corresponds to Theorem 2.3. and Remark 2.4. In Appendix D.1, we show how it follows from Assumptions 1 and 2 by proving the two intermediary results used in Abeille et al. (2022) to prove the result.

Proposition 7 (Adapted from (Abeille et al., 2022, Thm. 2.3 & Rem. 2.4.))

Under Assumptions 1 and 2, there is $L_W \in \mathbb{R}_+$, independent of ε , such that for any $\theta \in \Theta$

- (i.) The map $x \mapsto \rho_\theta^*(x)$ is constant, taking only one value which we denote by $\rho_\theta^* \in \mathbb{R}$;
- (ii.) There is an L_W -Lipschitz function W_θ^* such that

$$\varepsilon \rho_\theta^* = \max_{a \in \mathbb{A}} \{ \mathbb{E}[W_\theta^*(x + \mu_\theta(x, a) + \Sigma \xi)] - W_\theta^*(x) + r(x, a) \} \quad \forall x \in \mathbb{R}^d; \quad (16)$$

- (iii.) There is $\pi_\theta^* \in \mathcal{A}$, such that for all $x \in \mathbb{R}^d$, $\pi_\theta^*(x)$ maximises the right hand side in (16), and $\pi_\theta^* \circ X^{\pi_\theta^*, \theta}$ is an optimal Markov control, i.e. $\rho_\theta^{\pi_\theta^*}(\cdot) \equiv \rho_\theta^*$.

In Abeille et al. (2022), Proposition 8 corresponds to Theorem 3.4. In Appendix D.2, we show that it also follows from Assumptions 1 and 2 by proving that (Abeille et al., 2022, Assumption 5) holds under Assumptions 1 and 2.

Proposition 8 (Adapted from Abeille et al. (2022, Thm. 3.4.))

Under Assumptions 1 and 2, for any $\theta \in \Theta$,

- (i.) The map $x \mapsto \bar{\rho}_\theta^*(x)$ is constant, taking only one value which we denote by $\bar{\rho}_\theta^* \in \mathbb{R}$.
- (ii.) There is an L_W -Lipschitz function $\bar{W}_\theta^* \in \mathcal{C}^2(\mathbb{R}^d; \mathbb{R})$ such that

$$\bar{\rho}_\theta^* = \max_{a \in \mathbb{A}} \left\{ \bar{\mu}_\theta(x, a)^\top \nabla \bar{W}_\theta^*(x) + \bar{r}(x, a) \right\} + \frac{1}{2} \text{Tr}[\bar{\Sigma} \bar{\Sigma}^\top \nabla^2 \bar{W}_\theta^*(x)], \quad \forall x \in \mathbb{R}^d. \quad (18)$$

- (iii.) There is $\bar{\pi}_\theta^* \in \mathcal{A}$ such that, for all $x \in \mathbb{R}^d$, $\bar{\pi}_\theta^*(x)$ maximises the right hand side in (18), and $\bar{\pi}_\theta^* \circ \bar{X}^{\bar{\pi}_\theta^*, \theta}$ is an optimal Markov control, i.e. $\bar{\rho}_\theta^{\bar{\pi}_\theta^*}(\cdot) \equiv \bar{\rho}_\theta^*$.

Remark 23 Proposition 8.(iii.) is not stated as is in (Abeille et al., 2022, Thm. 3.4), but it follows from it by the same arguments as (Abeille et al., 2022, Remark 2.4).

Propositions 7 and 8 together ensure that both the prelimit and limit regimes are well posed, while Proposition 9 gives the rate of convergence of the control problems along this limit. This result is essentially contained in the proof of (Abeille et al., 2022, Thm. 3.6), but since its statement is different, we include a proof for completeness in Appendix D.3.

Proposition 9 (Adapted from Abeille et al. (2022, Thm. 3.6.))

Under Assumptions 1 and 2, for any $\gamma \in (0, 1)$, there is a constant $C_\gamma > 0$, independent of ε , such that, for any $\theta \in \Theta$,

$$|\bar{\rho}_\theta^* - \rho_\theta^*| \leq C_\gamma \varepsilon^{\frac{\gamma}{2}} \text{ and } \rho_\theta^* - \rho_\theta^*(0) \leq C_\gamma \varepsilon^{\frac{\gamma}{2}}. \quad (19)$$

Moreover, there is a function $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ such that,

$$\varepsilon \rho_\theta^*(0) = \mathbb{E}[\bar{W}_\theta^*(x + \mu_\theta(x, a) + \Sigma\xi)] - \bar{W}_\theta^*(x) + r(x, \bar{\pi}_\theta^*(x)) + e_\theta(x), \forall x \in \mathbb{R}^d \quad (20)$$

and there is $C'_\gamma > 0$, independent of ε , such that $|e_\theta(x)| \leq C'_\gamma \varepsilon^{1+\frac{\gamma}{2}}(1 + \|x\|^3)$ for all $x \in \mathbb{R}^d$.

D.1. Proof of Proposition 7

In Abeille et al. (2022), Theorem 2.3 and Remark 2.4 follow from Lemmas A.1 and A.2, which respectively give a mixing condition and a moment bound for $X^{\alpha, \theta}$. We already proved (Abeille et al., 2022, Lemma A.2) in Lemma 15. Moreover, Lemma 24 which reproduced (Abeille et al., 2022, Lemmas A.1) holds with only minor modifications of the proof from Abeille et al. (2022).

Lemma 15 Under Assumptions 1 and 2, for any $p \geq 2$, there is a constant $\mathfrak{c}'_p > 0$ independent of ε such that

$$\mathbb{E} \left[\|X_t^{x_0, \alpha, \theta}\|^p \right] \leq \frac{1}{\ell_\gamma^p} \left(L_\gamma^p e^{-\frac{\mathfrak{c}_\gamma}{4}t} \|x_0\|^p + \frac{4\mathfrak{c}'_p}{\mathfrak{c}_\gamma} \left(1 - e^{-\frac{\mathfrak{c}_\gamma}{4}t} \right) \right),$$

for any $(x_0, \alpha, \theta) \in \mathbb{R}^d \times \mathcal{A} \times \Theta$ and $t \in [0, +\infty)$.

Lemma 24 For any $(x, x') \in \mathbb{R}^d \times \mathbb{R}^d$, $\theta \in \Theta$, and $\alpha \in \mathcal{A}$,

$$\mathbb{E} \left[\|X_t^{x, \alpha, \theta} - X_t^{x', \alpha, \theta}\| \right] \leq \frac{L_\gamma}{\ell_\gamma} \|x - x'\| e^{-\mathfrak{c}_\gamma t}$$

for any $t \in [0, +\infty)$.

Proof We can follow the proof of Abeille et al. (2022) using Assumption 2 directly without resorting to the higher order Lyapunov function ζ which they use. ■

D.2. Proof of Proposition 8

Proposition 8, such as it is stated in (Abeille et al., 2022, Thm 3.4.) relies on their Assumption 5. This assumption contains two conditions, which we will show respectively in Lemmata 25 and 26.

As detailed in (Abeille et al., 2022, Remark 3.2.(i)), the first condition can be shown by proving an analogue of (Abeille et al., 2022, Lemma A.1) for the diffusive limit process (23). In terms of arguments of the proof, this analogue requires only a change in the stochastic generator used in Itô's Lemma⁶. In the proof of Lemma 25, we, therefore, show how to adapt (Abeille et al., 2022, Lemma A.1) to the generator of the diffusion under Assumptions 1 and 2.

In the proof of (Abeille et al., 2022, Lemma A.1), there are two key steps. First, study the discounted version of the control problem, and show that it is equi-Lipschitz continuous in the discount, which rests on the result in Lemma 25. Then one takes the vanishing discount limit in the HJB equation using the theory of viscosity solutions to complete the proof.

Lemma 25 *For any $(x_0, x'_0) \in \mathbb{R}^d \times \mathbb{R}^d$, $\theta \in \Theta$, $\alpha \in \mathcal{A}$,*

$$\mathbb{E} \left[\left\| \bar{X}_t^{x, \alpha, \theta} - \bar{X}_t^{x', \alpha, \theta} \right\| \right] \leq \frac{L_{\mathcal{V}}}{\ell_{\mathcal{V}}} \|x - x'\| e^{-c_{\mathcal{V}} t}$$

for any $t \in [0, +\infty)$.

Proof If $x_0 = x'_0$, this is trivially true by pathwise-uniqueness, so we suppose $x_0 \neq x'_0$. Let us consider $(x_1, x_2) \in \mathbb{R}^d \times \mathbb{R}^d$ with $x_1 \neq x_2$. By a Taylor expansion in (4), we obtain as $\varepsilon \rightarrow 0$

$$(\bar{\mu}(x_1, a) - \bar{\mu}(x_2, a))^\top \nabla \mathcal{V}(x_1 - x_2) \leq -c_{\mathcal{V}} \mathcal{V}(x_1 - x_2). \quad (60)$$

The Lyapunov function \mathcal{V} is not differentiable at 0, so we will construct an approximating sequence for it. Let erf denote the error function and let $\mathcal{V}_\iota := \mathcal{V} \operatorname{erf}(\iota \mathcal{V})$ for $\iota > 0$. Note that $\mathcal{V}_\iota \in \mathcal{C}^1(\mathbb{R}^d; \mathbb{R}_+)$ and \mathcal{V}_ι is Lipschitz, let us show that it satisfies (60) everywhere.

Let $z := x_1 - x_2$. Since $z \neq 0$ we have

$$\nabla \mathcal{V}_\iota(z) = \nabla \mathcal{V}(z) \left(\operatorname{erf}(\iota \mathcal{V}(z)) + \frac{2\iota}{\sqrt{\pi}} \mathcal{V}(z) e^{-\iota^2 \mathcal{V}^2(z)} \right).$$

By Assumption 2, this implies that

$$\begin{aligned} (\bar{\mu}_\theta(x_1, a) - \bar{\mu}_\theta(x_2, a))^\top \nabla \mathcal{V}_\iota(z) &\leq -c_{\mathcal{V}} \mathcal{V}(z) \operatorname{erf}(\iota \mathcal{V}(z)) - \frac{2\iota}{\sqrt{\pi}} c_{\mathcal{V}} \mathcal{V}(z)^2 e^{-\iota^2 \mathcal{V}^2(z)} \\ &\leq -c_{\mathcal{V}} \mathcal{V}_\iota(z). \end{aligned} \quad (61)$$

Since $\nabla \mathcal{V}_\iota$ is continuous in z , and so is the right-hand side, we can let $\|z\| \rightarrow 0$ and conclude the bound also holds for $x_1 = x_2$.

6. For a general overview of this sort of stability results and of Stochastic Lyapunov conditions in the diffusive case, see e.g. (Khasminskii, 2012, § 5.7).

We now apply Itô's lemma for the process $\bar{X}^{x,\alpha,\theta} - \bar{X}^{x',\alpha,\theta}$ to \mathcal{V}_ι . Using (61), this yields, for $t \geq t_0 \geq 0$,

$$\begin{aligned} & \mathbb{E} \left[\mathcal{V}_\iota \left(\bar{X}_t^{x_0,\alpha,\theta} - \bar{X}_t^{x'_0,\alpha,\theta} \right) \right] \\ & \leq \mathbb{E} \left[\mathcal{V}_\iota \left(\bar{X}_{t_0}^{x_0,\alpha,\theta} - \bar{X}_{t_0}^{x'_0,\alpha,\theta} \right) \right] \\ & + \mathbb{E} \left[\int_{t_0}^t \left(\bar{\mu}_\theta \left(\bar{X}_s^{x_0,\alpha,\theta}, \alpha_s \right) - \bar{\mu}_\theta \left(\bar{X}_s^{x'_0,\alpha,\theta}, \alpha_s \right) \right)^\top \nabla \mathcal{V}_\iota \left(\bar{X}_s^{x_0,\alpha,\theta} - \bar{X}_s^{x'_0,\alpha,\theta} \right) ds \right] \\ & \leq \mathbb{E} \left[\mathcal{V}_\iota \left(\bar{X}_{t_0}^{x_0,\alpha,\theta} - \bar{X}_{t_0}^{x'_0,\alpha,\theta} \right) \right] - \int_{t_0}^t \mathbf{c}_\mathcal{V} \mathbb{E} \left[\mathcal{V}_\iota \left(X_s^{x_0,\alpha,\theta} - X_s^{x'_0,\alpha,\theta} \right) \right] ds. \end{aligned}$$

We conclude by the same ODE comparison argument as in the proof of Lemma 15 and then pass to the limit as $\iota \rightarrow 0$ to obtain the claimed result using Assumption 2.(i). \blacksquare

While Lemma 25 showed that (Abeille et al., 2022, Assumption 5.(i)) is implied by Assumptions 1 and 2. It remains now to verify their Assumption 5.(ii). Note that by (Abeille et al., 2022, Remark 3.2.(ii)), an equation of the form of their (3.3) is sufficient to do so. Lemma 26 gives exactly this result with (62), by noting that (Abeille et al., 2022, (3.4)) holds by Assumption 2.

Lemma 26 *Under Assumptions 1 and 2, for any $p \geq 2$ there are $(\bar{\mathbf{c}}_p, \bar{\mathbf{c}}'_p) \in \mathbb{R}_+^2$ such that*

$$\bar{\mu}_\theta(x, a)^\top \nabla \mathcal{V}(x)^p + \text{Tr}[\bar{\Sigma} \bar{\Sigma}^\top \nabla^2 \mathcal{V}(x)^p] \leq -\bar{\mathbf{c}}_p \mathcal{V}(x)^p + \bar{\mathbf{c}}'_p \quad (62)$$

for any $(x, a, \theta) \in \mathbb{R}^d \times \mathbb{A} \times \Theta$.

Proof Let us take $(x, x') \in \mathbb{R}^d \times \mathbb{R}^d$ such that $\|x - x'\| \geq \varepsilon/(1 - \varepsilon L_0)$, which implies $\|x - x' + \Delta(\mu_\theta(x, a) - \mu_\theta(x', a))\| > 0$ for any $\Delta \in [0, 1]$ and for all $(a, \theta) \in \mathbb{A} \times \Theta$ and we can expand (4), which gives

$$\begin{aligned} -\varepsilon \mathbf{c}_\mathcal{V} \mathcal{V}(x - x') & \geq (\mu_\theta(x, a) - \mu_\theta(x', a))^\top \nabla \mathcal{V}(x - x') \\ & + \frac{1}{2} (\mu_\theta(x, a) - \mu_\theta(x', a))^\top \nabla^2 \mathcal{V}(\hat{x}) (\mu_\theta(x, a) - \mu_\theta(x', a)), \end{aligned}$$

in which $\hat{x} = x + \hat{\Delta}(x' - x)$ for some $\hat{\Delta} \in [0, 1]$. Thus

$$\begin{aligned} & (\bar{\mu}_\theta(x, a) - \bar{\mu}_\theta(x', a))^\top \nabla \mathcal{V}(x - x') \\ & \leq -\mathbf{c}_\mathcal{V} \mathcal{V}(x - x') - \frac{\varepsilon}{2} (\bar{\mu}_\theta(x, a) - \bar{\mu}_\theta(x', a))^\top \nabla^2 \mathcal{V}(\hat{x}) (\bar{\mu}_\theta(x, a) - \bar{\mu}_\theta(x', a)). \end{aligned}$$

Letting $\varepsilon \rightarrow 0$, the constraint on (x, x') vanishes as well as the second term (on compact sets), and we recover

$$(\bar{\mu}_\theta(x, a) - \bar{\mu}_\theta(x', a))^\top \nabla \mathcal{V}(x - x') + \frac{1}{2} \text{Tr}[\bar{\Sigma} \bar{\Sigma}^\top \nabla^2 \mathcal{V}(x - x')] \leq -\mathbf{c}_\mathcal{V} \mathcal{V}(x - x') + \frac{d}{2} \|\bar{\Sigma}\|_{\text{op}}^2 M'_\mathcal{V}.$$

Taking $x' = 0$ implies that

$$\bar{\mu}_\theta(x, a)^\top \nabla \mathcal{V}(x) + \frac{1}{2} \text{Tr}[\bar{\Sigma} \bar{\Sigma}^\top \nabla^2 \mathcal{V}(x)] \leq -\mathbf{c}_\gamma \mathcal{V}(x) + C$$

for all $(x, a) \in \mathbb{R}_*^d \times \mathbb{A}$, in which $C := d \|\bar{\Sigma}\|_{\text{op}}^2 M'_\gamma / 2 + L_0 M_\gamma$.

Notice that, since $\mathcal{V} \in \mathcal{C}^2(\mathbb{R}_*^d; \mathbb{R}_+)$ and vanishes at 0 (see Assumption 1), $\mathcal{V}(\cdot)^p$ can be extended by continuity at 0 so that $\mathcal{V}(\cdot)^p \in \mathcal{C}^2(\mathbb{R}^d; \mathbb{R}_+)$. For any $(x, a, \theta) \in \mathbb{R}^d \times \mathbb{A} \times \Theta$, let

$$\begin{aligned} k(x, a) &:= \bar{\mu}_\theta(x, a)^\top \nabla \mathcal{V}(x)^p + \frac{1}{2} \text{Tr} \left[\bar{\Sigma} \bar{\Sigma}^\top \nabla^2 \mathcal{V}(x)^p \right] \\ &= p \bar{\mu}_\theta(x, a)^\top \nabla \mathcal{V}(x) \mathcal{V}(x)^{p-1} \\ &\quad + \frac{1}{2} \text{Tr} \left[\bar{\Sigma} \bar{\Sigma}^\top \left(p \mathcal{V}(x)^{p-1} \nabla^2 \mathcal{V}(x) + p(p-1) \mathcal{V}(x)^{p-2} \nabla \mathcal{V}(x) \nabla^\top \mathcal{V}(x) \right) \right] \\ &= p \mathcal{V}^{p-1}(x) \left(\bar{\mu}_\theta(x, a)^\top \nabla \mathcal{V}(x) + \frac{1}{2} \text{Tr}[\bar{\Sigma} \bar{\Sigma}^\top \nabla^2 \mathcal{V}(x)] \right) \\ &\quad + \frac{p(p-1)}{2} \mathcal{V}(x)^{p-2} \text{Tr}[\bar{\Sigma} \bar{\Sigma}^\top \nabla \mathcal{V}(x) \nabla^\top \mathcal{V}(x)] \\ &\leq -p \mathbf{c}_\gamma \mathcal{V}(x)^p + C p \mathcal{V}(x)^{p-1} + \frac{dp(p-1)}{2} (\|\bar{\Sigma}\|_{\text{op}} M_\gamma)^2 \mathcal{V}(x)^{p-2} \end{aligned}$$

and we can now choose $\bar{\mathbf{c}}_p = -p \mathbf{c}_\gamma / 2$, for which there exists a constant $\bar{\mathbf{c}}'_p$ such that

$$-\bar{\mathbf{c}}_p \mathcal{V}^p(x) + C p \mathcal{V}^{p-1}(x) + \frac{dp(p-1)}{2} (\|\bar{\Sigma}\|_{\text{op}} M_\gamma)^2 \mathcal{V}^{p-2}(x) \leq \bar{\mathbf{c}}'_p$$

for all $x \in \mathbb{R}^d$. ■

D.3. Proof of Proposition 9

The rest of this section is dedicated to showing Proposition 9 using modifications of the proof of (Abeille et al., 2022, Thm. 3.6.) to which it corresponds. Here we produce a self-contained proof in order to clarify how (20) is derived from the proof.

Proposition 9 (Adapted from Abeille et al. (2022, Thm. 3.6.))

Under Assumptions 1 and 2, for any $\gamma \in (0, 1)$, there is a constant $C_\gamma > 0$, independent of ε , such that, for any $\theta \in \Theta$,

$$|\bar{\rho}_\theta^* - \rho_\theta^*| \leq C_\gamma \varepsilon^{\frac{\gamma}{2}} \quad \text{and} \quad \rho_\theta^* - \rho_{\bar{\pi}_\theta^*}^*(0) \leq C_\gamma \varepsilon^{\frac{\gamma}{2}}. \quad (19)$$

Moreover, there is a function $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ such that,

$$\varepsilon \rho_{\bar{\pi}_\theta^*}^*(0) = \mathbb{E}[\bar{W}_\theta^*(x + \mu_\theta(x, a) + \Sigma \xi)] - \bar{W}_\theta^*(x) + r(x, \bar{\pi}_\theta^*(x)) + e_\theta(x), \quad \forall x \in \mathbb{R}^d \quad (20)$$

and there is $C'_\gamma > 0$, independent of ε , such that $|e_\theta(x)| \leq C'_\gamma \varepsilon^{1+\frac{\gamma}{2}} (1 + \|x\|^3)$ for all $x \in \mathbb{R}^d$.

Proof The first part of Proposition 9, i.e. (19), corresponds to Abeille et al. (2022, Thm. 3.6.), which we previously showed holds in our setting by verifying its assumptions. We now prove the second claim. Let

$$\delta r_\theta^\varepsilon(x, a) := \bar{\mu}_\theta(x, a)^\top \nabla \bar{W}_\theta^*(x) + \frac{1}{2} \text{Tr}[\bar{\Sigma} \bar{\Sigma}^\top \nabla^2 \bar{W}_\theta^*(x)] - \frac{1}{\varepsilon} (\mathbb{E} [\bar{W}_\theta^*(\psi_\theta^\varepsilon(x, a) + \Sigma \xi)] - \bar{W}_\theta^*(x)).$$

From (18), and Proposition 8.(iii.) we have

$$\begin{aligned} \bar{\rho}_\theta^* &= \max_{a \in \mathbb{A}} \left\{ \bar{\mu}_\theta(x, a)^\top \nabla \bar{W}_\theta^* + \frac{1}{2} \text{Tr}[\bar{\Sigma} \bar{\Sigma}^\top \nabla^2 \bar{W}_\theta^*(x)] + \bar{r}(x, a) \right\} \\ &= \bar{\mu}_\theta(x, \bar{\pi}_\theta^*(x))^\top \nabla \bar{W}_\theta^*(x) + \frac{1}{2} \text{Tr}[\bar{\Sigma} \bar{\Sigma}^\top \nabla^2 \bar{W}_\theta^*(x)] + \bar{r}(x, \bar{\pi}_\theta^*(x)) \end{aligned}$$

which implies

$$\varepsilon \rho_\theta^{\bar{\pi}_\theta^*}(0) = \mathbb{E}[\bar{W}_\theta^*(\psi_\theta^\varepsilon(x, \bar{\pi}_\theta^*(x)) + \Sigma \xi)] - \bar{W}_\theta^*(x) + r(x, \bar{\pi}_\theta^*(x)) + \varepsilon(\delta r_\theta^\varepsilon(x, \bar{\pi}_\theta^*(x)) + \bar{\rho}_\theta^* - \rho_\theta^{\bar{\pi}_\theta^*}(0)).$$

Note that $|\delta r_\theta^\varepsilon(x, \bar{\pi}_\theta^*(x))| \leq \sup_{a \in \mathbb{A}} |\delta r_\theta^\varepsilon(x, a)|$, which by Abeille et al. (2022, (3.10)) is bounded by $c_\gamma \varepsilon^{\frac{\gamma}{2}} (1 + \|x\|^3)$ for some constant $c_\gamma > 0$. An application of (19) yields

$$\bar{\rho}_\theta^* - \rho_\theta^{\bar{\pi}_\theta^*}(0) = \bar{\rho}_\theta^* - \rho_\theta^* + \rho_\theta^* - \rho_\theta^{\bar{\pi}_\theta^*}(0) \leq 2C_\gamma \varepsilon^{\frac{\gamma}{2}}$$

and, at the same time, $\bar{\rho}_\theta^* - \rho_\theta^{\bar{\pi}_\theta^*}(0) \geq \bar{\rho}_\theta^* - \rho_\theta^* \geq -C_\gamma \varepsilon^{\frac{\gamma}{2}}$. Therefore, there is a function $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ such that (20) holds, which also satisfies

$$|e_\theta(x)| \leq (2C_\gamma + c_\gamma) \varepsilon^{1+\frac{\gamma}{2}} (1 + \|x\|^3).$$

■

Appendix E. Regret Analysis

In this final appendix, we complete the analysis of the regret of Algorithms 1 and 2. Seeing as the proof of Theorem 4 is directly contained in the proof of Theorem 10, as indicated in the proof sketch of the latter, we will only prove the latter. First, we will give the regret decomposition, and then in the later sections, we will bound terms one by one calling upon the results of the previous appendices.

Theorem 10 *Under Assumptions 1 and 2, for any $\delta \in (0, 1)$, $x_0 \in \mathbb{R}^d$, and $\gamma \in (0, 1)$, there is a pair $(C_\gamma, C) \in \mathbb{R}_+^2$ of constants independent of ε such that Algorithm 2 achieves*

$$R_T(\varpi') \leq 2C_\gamma \varepsilon^{\frac{\gamma}{2}} T + C \sqrt{d_{E, [T\varepsilon^{-1}]} \log \left(\mathcal{N}_{[T\varepsilon^{-1}]}^\varepsilon \right) T \log(T\delta^{-1})} \quad (21)$$

with probability at least $1 - \delta$.

E.1. Regret Decomposition

Recall that we defined $k : n \in \mathbb{N} \mapsto k(n)$ as the map associating to each event n the episode of Algorithm 1 in which they occur. Like in Section 4.3, let us define $\theta_n = \tilde{\theta}_{k(n)}$ for all $n \in \mathbb{N}$. The regret of Algorithm 1, which generates the control $\varpi \in \mathcal{A}$, is

$$\mathcal{R}_T(\varpi) := T\rho_{\theta^*}^* - \sum_{n=1}^{N_T} r(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n})$$

By definition of ϖ in Algorithm 1, $\varpi_{\tau_n} = \bar{\pi}_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*})$, so that

$$\mathcal{R}_T(\varpi) := T\rho_{\theta^*}^* - \sum_{n=1}^{N_T} r(X_{\tau_n}^{\varpi, \theta^*}, \bar{\pi}_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*}))$$

At the heart of the decomposition is the use of the HJB-type equation (20) applied for each n at the point $X_{\tau_n}^{\varpi, \theta^*}$. For clarity, let us introduce for all $n \in \mathbb{N}$ the random variable $\tilde{X}_{\tau_{n+1}}^{\varpi, \theta_n}$ equal in distribution, conditionally on \mathcal{F}_{τ_n} , to the random variable $\psi_{\theta_n}^\varepsilon(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) + \Sigma\xi_{n+1}$. With this notation (20) becomes

$$\varepsilon\rho_{\theta_n}^*(0) = \mathbb{E}[\bar{W}_{\theta_n}^*(\tilde{X}_{\tau_{n+1}}^{\varpi, \theta_n}) | \mathcal{F}_{\tau_n}] - \bar{W}_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*}) + r(X_{\tau_n}^{\varpi, \theta^*}, \bar{\pi}_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*})) + e_{\theta_n}(X_{\tau_n}^{\varpi, \theta^*}). \quad (63)$$

This *imagined* evolution of the system represents the counterfactual induced by a single step transition at time τ_{n+1} , according to the belief in θ_n . With this notation, applying (63) yields

$$\begin{aligned} \mathcal{R}_T(\varpi) &= T\rho_{\theta^*}^* - \sum_{n=1}^{N_T} \varepsilon\rho_{\theta_n}^*(0) + \sum_{n=1}^{N_T} e_{\theta^*}(X_{\tau_n}^{\varpi, \theta^*}) + \sum_{n=1}^{N_T} \mathbb{E}[\bar{W}_{\theta_n}^*(\tilde{X}_{\tau_{n+1}}^{\varpi, \theta_n}) | \mathcal{F}_{\tau_n}] - \bar{W}_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*}). \\ &= (T - \varepsilon N_T)\rho_{\theta^*}^* \end{aligned} \quad (R_1)$$

$$+ \varepsilon \sum_{n=1}^{N_T} (\rho_{\theta^*}^* - \rho_{\theta_n}^*(0)) + \sum_{n=1}^{N_T} e_{\theta^*}(X_{\tau_n}^{\varpi, \theta^*}) \quad (R_2)$$

$$+ \sum_{n=1}^{N_T} \mathbb{E}[\bar{W}_{\theta_n}^*(\tilde{X}_{\tau_{n+1}}^{\varpi, \theta_n}) | \mathcal{F}_{\tau_n}] - \bar{W}_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*}). \quad (64)$$

The first term, (R_1) , quantifies the deviation of the Poisson clock from its mean. On the other hand, (R_2) quantifies both the optimistic nature of Algorithm 1 and the approximation error of its approximate planning. The third term, (64) , resembles a martingale (up to reordering), but it fails to be one on two key counts. First, the element from the family of functions $(\bar{W}_{\theta_n}^*)_{n \in \mathbb{N}}$ used at each step n changes. Second, the expectation terms are with respect to the counterfactual transitions $(\tilde{X}_{\tau_{n+1}}^{\varpi, \theta^*})_{n \in \mathbb{N}}$ while the random terms use the real transitions $(X_{\tau_{n+1}}^{\varpi, \theta^*})_{n \in \mathbb{N}}$.

Note that we can control the difference between the counterfactual and the real trajectory at a one-step time horizon, by using

$$\tilde{X}_{\tau_{n+1}}^{\varpi, \theta} \stackrel{d}{=} X_{\tau_{n+1}}^{\varpi, \theta^*} - \mu_{\theta^*}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) + \mu_{\theta}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}), \quad (65)$$

in which $\stackrel{d}{=}$ denotes equality in the same conditionally distributional sense as above. By adding and subtracting relevant terms to exhibit the key quantities we get:

$$\begin{aligned} \sum_{n=1}^{N_T} \mathbb{E}[\bar{W}_{\theta_n}^*(\tilde{X}_{\tau_{n+1}}^{\varpi, \theta_n}) | \mathcal{F}_{\tau_n}] - \bar{W}_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*}) &\leq \sum_{n=1}^{N_T} \mathbb{E}[\bar{W}_{\theta_n}^*(\tilde{X}_{\tau_{n+1}}^{\varpi, \theta_n}) | \mathcal{F}_{\tau_n}] - \mathbb{E}[\bar{W}_{\theta_n}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}] \\ &+ \sum_{n=1}^{N_T} \mathbb{E}[\bar{W}_{\theta_n}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}] - \mathbb{E}[\bar{W}_{\theta_{n+1}}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}] \\ &+ \sum_{n=1}^{N_T} \mathbb{E}[\bar{W}_{\theta_{n+1}}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}] - \bar{W}_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*}). \end{aligned}$$

Using (65) , and the uniform L_W -Lipschitzness of $(\bar{W}_{\theta_n}^*)_{n \in \mathbb{N}}$, we get for each $n \in \mathbb{N}$

$$\mathbb{E}[\bar{W}_{\theta_n}^*(\tilde{X}_{\tau_{n+1}}^{\varpi, \theta_n}) | \mathcal{F}_{\tau_n}] - \mathbb{E}[\bar{W}_{\theta_n}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}] \leq L_W \left\| \mu_{\theta_n}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) - \mu_{\theta^*}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) \right\|$$

and thus the regret term (64) is bounded by

$$\sum_{n=1}^{N_T} \mathbb{E}[\bar{W}_{\theta_n}^*(\tilde{X}_{\tau_{n+1}}^{\varpi, \theta_n}) | \mathcal{F}_{\tau_n}] - \bar{W}_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*}) \leq R_3 + R_4 + R_5$$

in which

$$R_3 := L_W \sum_{n=1}^{N_T} \left\| \mu_{\theta_n}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) - \mu_{\theta^*}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) \right\| \quad (R_3)$$

$$R_4 := \sum_{n=1}^{N_T} \mathbb{E}[\bar{W}_{\theta_n}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) - \bar{W}_{\theta_{n+1}}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}] \quad (R_4)$$

$$R_5 := \sum_{n=1}^{N_T} \mathbb{E}[\bar{W}_{\theta_{n+1}}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}] - \bar{W}_{\theta_n}^*(X_{\tau_n}^{\varpi, \theta^*}). \quad (R_5)$$

At the end of this decomposition, we have constructed a true martingale in (R_5) , which we bound in Appendix E.6. The first term (R_3) accumulates the fit error described in Proposition 6, up to the lazy updates, which we study in Appendix E.4. The term (R_4) is bounded

by the number of effective updates of θ_n (namely, $\sum_{n=1}^{N_T} \mathbb{1}_{\{\theta_{n+1} \neq \theta_n\}}$) in Appendix E.5. Finally, the bounds on (R_1) and (R_2) are given in Appendices E.2 and E.3 respectively.

To combine the high-probability events used to bound (R_1) and (R_5) , with the event of Proposition 5 used by the other terms, we will perform a union bound. This corresponds to the $\delta/3$ used in the definition of the confidence sets of Algorithm 1.

E.2. Bounding the Poisson clock variation term (R_1)

We bound (R_1) using Lemma 27 which is a standard sub-exponential concentration result, see e.g. (Buldygin and Kozachenko, 2000, Lemma 4.1). It implies

$$\mathbb{P} \left(|T - \varepsilon N_T| \geq 2\sqrt{\varepsilon T \log \left(\frac{6}{\delta} \right)} \vee 2\varepsilon \log \left(\frac{6}{\delta} \right) \right) \leq \frac{\delta}{3}.$$

Lemma 27 For any $T \in \mathbb{R}_+^*$ and $\delta \in (0, 1)$,

$$\mathbb{P} \left(|\varepsilon N_T - T| > 2\sqrt{\varepsilon T \log \left(\frac{2}{\delta} \right)} \vee 2\varepsilon \log \left(\frac{2}{\delta} \right) \right) \leq \delta.$$

Proof Let $v := \varepsilon^{-1}T$. For any $\lambda \in [-1, 1]$, $\mathbb{E}[e^{\lambda(N_T - v)}] = \exp(v(e^\lambda - 1 - \lambda)) \leq e^{\lambda^2 v}$. Therefore, N_T is $(\sqrt{2v}, 1)$ -subexponential (see e.g. (Buldygin and Kozachenko, 2000)) and therefore,

$$\mathbb{P}(|N_T - v| > c) \leq \begin{cases} e^{-\frac{c^2}{4v}} & \text{for } c \in (0, 2v] \\ e^{-\frac{c}{2}} & \text{for } c > 2v \end{cases},$$

which implies

$$\mathbb{P} \left(|N_T - v| > 2\sqrt{v \log \left(\frac{2}{\delta} \right)} \mathbb{1}_{\{\delta \geq e^{-v}\}} + 2 \log \left(\frac{2}{\delta} \right) \mathbb{1}_{\{\delta \leq e^{-v}\}} \right) \leq \delta. \quad \blacksquare$$

E.3. Bounding the optimistic approximation term (R_2)

There are two terms in (R_2) . The second is the most straightforward as it can be bounded by applying the bound on e_{θ^*} of Proposition 9, which yields

$$\sum_{n=1}^{N_T} e_{\theta^*}(X_{\tau_n}^{\varpi, \theta^*}) \leq 2C'_\gamma N_T \varepsilon^{1+\frac{\gamma}{2}} (1 + \sup_{s \leq T} \|X_s^{\varpi, \theta^*}\|^3).$$

We decompose the remaining term of (R_2) into

$$\begin{aligned} \varepsilon \sum_{n=1}^{N_T} (\rho_{\theta^*}^* - \rho_{\theta_n}^*) &= \varepsilon \sum_{n=1}^{N_T} \left(\rho_{\theta^*}^* - \bar{\rho}_{\theta^*}^* + \bar{\rho}_{\theta^*}^* - \bar{\rho}_{\theta_n}^* + \bar{\rho}_{\theta_n}^* - \rho_{\theta_n}^* + \rho_{\theta_n}^* - \rho_{\theta_n}^* \right) \\ &\leq 4N_T C_\gamma \varepsilon^{1+\frac{\gamma}{2}} + \varepsilon \sum_{n=1}^{N_T} \left(\bar{\rho}_{\theta^*}^* - \bar{\rho}_{\theta_n}^* \right) \end{aligned}$$

by applying Proposition 9 to all but the second pair of terms.

On the event of Proposition 5, with $\delta/3$ in place of δ , we have $\theta^* \in \cap_{n \in \mathbb{N}^*} \mathcal{C}_n(\delta/3)$ and thus, by definition of Algorithm 1, $\bar{\rho}_{\theta^*}^* - \bar{\rho}_{\theta_n}^* \leq 0$ for all $n \in \mathbb{N}^*$. Thus, on this event we have

$$\varepsilon \sum_{n=1}^{N_T} (\rho_{\theta^*}^* - \rho_{\theta_n}^*) \leq 4N_T C_\gamma \varepsilon^{1+\frac{\gamma}{2}}.$$

E.4. Bounding the prediction error term (R_3)

Because of the lazy updates, $\mu_{\theta_n} = \mu_{\theta_{k(n)}}$ is chosen within $\mathcal{C}_{k(n)}(\delta/3)$ instead of $\mathcal{C}_n(\delta/3)$ preventing us from using directly Proposition 22. Nevertheless, the lazy update scheme is designed not to degrade the overall learning performance by more than a constant factor. Leveraging (7),

$$\sum_{i=1}^{n-1} \left\| \mu_{\theta_n}(X_{\tau_i}^{\varpi, \theta^*}, \varpi_{\tau_i}) - \mu_{\theta^*}(X_{\tau_i}^{\varpi, \theta^*}, \varpi_{\tau_i}) \right\| \leq \begin{cases} 2\beta_n(\delta/3) & \text{if } n < n_k \\ \beta_n(\delta/3) & \text{if } n = n_k \end{cases} \quad (66)$$

As a result, μ_{θ_n} is chosen within an inflated version of $\mathcal{C}_n(\delta/3)$, defined as in (6) but with $\beta_n(\delta/3)$ replaced by $2\beta_n(\delta/3)$. Thus, we can follow the same arguments as in the proof of Proposition 6, by applying Proposition 22 to the inflated confidence sets, up to the constant factor 2 in the bounds. And therefore on the event of Proposition 5, we have

$$\begin{aligned} R_3 &= L_W \sum_{n=1}^{N_T} \left\| \mu_{\theta_n}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) - \mu_{\theta^*}(X_{\tau_n}^{\varpi, \theta^*}, \varpi_{\tau_n}) \right\| \\ &\leq 6L_W \beta_{N_T}(\delta/3) \sqrt{d_{E, N_t}} + L_W d_{E, N_t} H_{\delta/3}(N_T). \end{aligned}$$

E.5. Bounding the lazy-update term (R_4)

We observe that (R_4) is bounded by

$$\begin{aligned} R_4 &= \sum_{n=1}^{N_T} \mathbb{E}[\bar{W}_{\theta_n}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) - \bar{W}_{\theta_{n+1}}^*(X_{\tau_{n+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_n}] \\ &\leq 2L_W \sum_{n=1}^{N_T} \mathbb{E} \left[\left(1 + \|X_{\tau_{n+1}}^{\varpi, \theta^*}\| \right) \mathbb{1}_{\{\theta_n \neq \theta_{n+1}\}} | \mathcal{F}_{\tau_n} \right] \\ &\leq 2L_W \sum_{n=1}^{N_T} \left((1 + \varepsilon L_0)(1 + \|X_{\tau_n}^{\varpi, \theta^*}\|) + \varepsilon^{\frac{1}{2}} \|\bar{\Sigma}\|_{\text{op}} \mathbb{E}[\|\xi_{n+1}\| | \mathcal{F}_{\tau_n}] \right) \mathbb{1}_{\{\theta_n \neq \theta_{n+1}\}} \\ &\leq 2L_W(1 + \varepsilon L_0) \left(1 + \sup_{s \leq T} \|X_s^{\varpi, \theta^*}\| + \sqrt{d} \varepsilon^{\frac{1}{2}} \|\bar{\Sigma}\|_{\text{op}} \right) \sum_{n=1}^{N_T} \mathbb{1}_{\{\theta_n \neq \theta_{n+1}\}}. \end{aligned}$$

Thus bounding the number of updates with Lemma 28 bounds (R_4).

Lemma 28 *Under Assumptions 1 and 2, Algorithm 1 generates episodes which satisfy for all $T \in \mathbb{R}_+$ and $\delta \in (0, 1)$*

$$\begin{aligned} \sum_{n=1}^{N_T} \mathbf{1}_{\{\theta_n \neq \theta_{n+1}\}} &\leq 4\beta_{N_T}(\delta/3)^2 d_{E, N_T} \left(3 + \log \left(\frac{N_T 8\varepsilon^2 L_0^2 (1 + \sup_{s \leq t} \|X_s^{\varpi, \theta^*}\|)}{16\beta_{N_T}(\delta/3)^4 d_{E, N_T}^2} \right) \right) \\ &\quad + 2d_{E, N_T} (1 + 2\beta_{N_T}(\delta/3)^2 d_{E, N_T}) (1 + 8\varepsilon^2 L_0^2 (1 + \sup_{s \leq t} \|X_s^{\varpi, \theta^*}\|^2)). \end{aligned}$$

Proof Consider $k \in \mathbb{N}^*$, by (7), each time we trigger an update we have

$$\begin{aligned} 2\beta_{n_k}(\delta/3)^2 &< \sup_{\theta \in \mathcal{C}_{n_{k-1}}(\delta)} \left\| \mu_\theta - \mu_{\hat{\theta}_{n_{k-1}}} \right\|_{n_k}^2 \\ &\leq \sup_{\theta \in \mathcal{C}_{n_{k-1}}(\delta)} \left\| \mu_\theta - \mu_{\hat{\theta}_{n_{k-1}}} \right\|_{n_{k-1}}^2 \\ &\quad + \sup_{\theta \in \mathcal{C}_{n_{k-1}}(\delta)} \sum_{n=n_{k-1}+1}^{n_k} \left\| \mu_\theta(X_{\tau_n}^{\varpi, \theta}, \varpi_{\tau_n}) - \mu_{\hat{\theta}_{n_{k-1}}}(X_{\tau_n}^{\varpi, \theta}, \varpi_{\tau_n}) \right\|^2 \\ &\leq \beta_{n_k}(\delta/3)^2 + \sum_{n=n_{k-1}+1}^{n_k} \Lambda(\mathcal{C}_{n_{k-1}}(\delta/3); X_{\tau_n}^{\varpi, \theta}, \varpi_{\tau_n})^2. \end{aligned}$$

Summing over all episodes, since the sequence $(\beta_n(\delta/3))_{n \in \mathbb{N}}$ is non-decreasing, we have that for all $T \in \mathbb{R}_+$

$$\sum_{n=1}^{N_T} \Lambda(\mathcal{C}_{n_k}(\delta/3); (X_{\tau_n}, \varpi_{\tau_n}))^2 \geq \sum_{k=1}^{K_T} \beta_{n_k}(\delta/3)^2 \geq K_T \beta_0(\delta/3)^2,$$

in which $K_T := k(N_T) \in \mathbb{N}$ is the number of episodes by time T . An application of the second part of Proposition 22, i.e. (59) now yields the desired result as $\beta_0(\delta/3)^2 = \varepsilon$. \blacksquare

E.6. Bounding the martingale term (R_5)

Let

$$Z_n := \mathbb{E}[\bar{W}_{\theta_n}^*(X_{\tau_n}^{\alpha, \theta^*}) | \mathcal{F}_{n-1}] - \bar{W}_{\theta_n}^*(X_{\tau_n}^{\alpha, \theta^*}).$$

By definition

$$R_5 = \mathbb{E}[\bar{W}_{\theta_{N_T+1}}^*(X_{\tau_{N_T+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_{N_T}}] + \bar{W}_{\theta_0}^*(x_0) + \sum_{n=1}^{N_T} Z_n.$$

On the one hand, Z_n is a $L_W \|\Sigma\|_{\text{op}}$ -Lipschitz function of ξ_n , which is Gaussian and of mean 0. Therefore, by (Boucheron et al., 2013, Thm 5.5), Z_n is $L_W \|\Sigma\|_{\text{op}}$ -sub-Gaussian and

$$\mathbb{P} \left(\sum_{n=1}^{N_T} Z_n > L_W \|\bar{\Sigma}\|_{\text{op}} \sqrt{2\varepsilon N_T \log \left(\frac{1}{\delta} \right)} \right) \leq \delta. \quad (67)$$

On the other hand, by the uniform Lipschitzness of $(\bar{W}_\theta^*)_{\theta \in \Theta}$, $\bar{W}_{\theta_0}^*(x_0) \leq L_W(1 + \|x_0\|)$ and

$$\begin{aligned} \mathbb{E}[\bar{W}_{\theta_{N_T+1}}^*(X_{\tau_{N_T+1}}^{\varpi, \theta^*}) | \mathcal{F}_{\tau_{N_T}}] &\leq L_W(1 + \mathbb{E}[\|X_{\tau_{N_T+1}}^{\varpi, \theta^*}\| | \mathcal{F}_{\tau_{N_T}}]) \\ &\leq L_W(1 + \varepsilon L_0 + (1 + \varepsilon L_0)\|X_{\tau_{N_T}}^{\varpi, \theta^*}\| + \varepsilon^{\frac{1}{2}}\|\bar{\Sigma}\|_{\text{op}}\mathbb{E}[\|\xi_{N_T+1}\| | \mathcal{F}_{\tau_{N_T}}]) \\ &\leq L_W(1 + \varepsilon L_0) \left(1 + \sup_{s \leq T} \|X_s^{\varpi, \theta^*}\|_2 + \varepsilon^{\frac{1}{2}}\|\bar{\Sigma}\|_{\text{op}}\sqrt{d}L_W \right). \end{aligned} \quad (68)$$

Combining (67) and (68) yields

$$R_5 \leq L_W \|\bar{\Sigma}\|_{\text{op}} \sqrt{2\varepsilon N_T \log\left(\frac{3}{\delta}\right)} + 2L_W(1 + \varepsilon L_0) \left(1 + \sup_{s \leq T} \|X_s^{\varpi, \theta^*}\| + \varepsilon^{\frac{1}{2}}\|\bar{\Sigma}\|_{\text{op}}\sqrt{d}L_W \right) \quad (69)$$

with probability at least $1 - \delta/3$.

E.7. Collecting the bounds

We conclude the proof of Theorem 4 by collecting all the terms from Appendices E.2–E.6 and simplifying them. By a union bound over the events listed in steps Appendices E.2, E.4 and E.6, with probability at least $1 - \delta$

$$\begin{aligned} \mathcal{R}_T(\varpi) &\leq 2L_0 \left(\sqrt{\varepsilon T \log\left(\frac{6}{\delta}\right)} \vee 2\varepsilon \log\left(\frac{6}{\delta}\right) \right) \\ &\quad + 4N_T C_\gamma \varepsilon^{1+\frac{\gamma}{2}} + 2C'_\gamma N_T \varepsilon^{1+\frac{\gamma}{2}} (1 + H_{\delta/3}^3(N_T)) \\ &\quad + 6L_W \beta_{N_T} (\delta/3) \sqrt{d_{\mathbb{E}, N_T}} + 2\varepsilon L_0 L_W d_{\mathbb{E}, N_T} (1 + H_{\delta/3}(N_T)) \\ &\quad + 2L_W(1 + \varepsilon L_0) \left(1 + H_{\delta/3}(N_T) + d\varepsilon^{\frac{1}{2}}\|\bar{\Sigma}\|_{\text{op}} \right) \\ &\quad \times \left(4\beta_{N_T} (\delta/3)^2 d_{\mathbb{E}, N_t} \left(3 + \log\left(\frac{N_t 8\varepsilon^2 L_0^2 (1 + H_{\delta/3}(N_T))}{16\beta_{N_T} (\delta/3)^4 d_{\mathbb{E}, N_T}^2}\right) \right) \right. \\ &\quad \left. + 2d_{\mathbb{E}, N_T} (1 + 2\beta_{N_t} (\delta/3)^2 d_{\mathbb{E}, N_T}) (1 + 8\varepsilon^2 L_0^2 (1 + H_{\delta/3}(N_T)^2)) \right) \\ &\quad + L_W \|\bar{\Sigma}\|_{\text{op}} \sqrt{2\varepsilon N_T \log\left(\frac{3}{\delta}\right)} + 2L_W(1 + \varepsilon L_0) (1 + H_{\delta/3}(N_T) + \varepsilon^{\frac{1}{2}}\|\bar{\Sigma}\|_{\text{op}}\sqrt{d}L_W). \end{aligned}$$

This can be more simply expressed for some constants $C_{\mathcal{R}}^{(i)} \in \mathbb{R}_+$, $i \in [5]$, as

$$\begin{aligned} \mathcal{R}_T(\varpi) &\leq C_{\mathcal{R}}^{(1)}(C_\gamma + C'_\gamma)\varepsilon^{1+\frac{\gamma}{2}}N_T \log(N_T)^3 + C_{\mathcal{R}}^{(2)}\sqrt{d_{\mathbf{E},N_T}\varepsilon N_T \log\left(\frac{N_T(1 + \varepsilon\mathcal{N}_{N_T}^\varepsilon)}{\delta}\right)} \\ &\quad + C_{\mathcal{R}}^{(3)}\left(1 + \varepsilon d_{\mathbf{E},N_T} \log(N_T) \log(N_T(1 + \varepsilon\mathcal{N}_{N_T}^\varepsilon))\right)d_{\mathbf{E},N_T} \log(N_T)^4 \\ &\quad + C_{\mathcal{R}}^{(4)}\sqrt{\varepsilon T \log\left(\frac{1}{\delta}\right)} + C_{\mathcal{R}}^{(5)}\left(1 + \log\left(\frac{1}{\delta}\right)\right) \end{aligned}$$

still with probability at least $1 - \delta$. On this high-probability event, we can write $\mathcal{R}_T(\varpi)$ (up rounding up $T\varepsilon^{-1}$ where necessary and up to a change in the constants) as

$$\begin{aligned} \mathcal{R}_T(\varpi) &\leq C_{\mathcal{R}}^{(1)}(C_\gamma + C'_\gamma)\varepsilon^{\frac{\gamma}{2}}T \log\left(\frac{T}{\varepsilon}\right) + C_{\mathcal{R}}^{(2)}\sqrt{d_{\mathbf{E},T\varepsilon^{-1}}T \log\left(\frac{T\varepsilon^{-1}(1 + \varepsilon\mathcal{N}_{T\varepsilon^{-1}}^\varepsilon)}{\delta}\right)} \\ &\quad + C_{\mathcal{R}}^{(3)}\left(1 + \varepsilon d_{\mathbf{E},T\varepsilon^{-1}} \log(T\varepsilon^{-1}) \log(T\varepsilon^{-1}(1 + \varepsilon\mathcal{N}_{T\varepsilon^{-1}}^\varepsilon))\right)d_{\mathbf{E},T\varepsilon^{-1}} \log(T\varepsilon^{-1})^4 \\ &\quad + C_{\mathcal{R}}^{(4)}\sqrt{\varepsilon T \log\left(\frac{1}{\delta}\right)} + C_{\mathcal{R}}^{(5)}\left(1 + \log\left(\frac{1}{\delta}\right)\right). \end{aligned}$$

Considering only the two dominant terms and ignoring logarithmic factors we get the claimed bound.