

---

# Provably Efficient Reinforcement Learning for Sparse Dynamical Systems with Non-Gaussian Noise

---

Davide Maran<sup>\*,1</sup>

Gianmarco Tedeschi<sup>\*,1</sup>

Enea Gusmeroli<sup>1</sup>

Marcello Restelli<sup>1</sup>

<sup>1</sup>Politecnico di Milano

\*Equal Contribution

## Abstract

The recent development of sparse methods for identifying nonlinear dynamical systems has opened new avenues for efficient and interpretable model-based reinforcement learning (RL). In this work, we study online RL in environments where the system dynamics, modeled as  $s' = f(s, a) + \text{noise}$ , is assumed to be sparse with respect to a big feature map, a structural idea inspired by the SINDY framework. We introduce an optimistic algorithm that combines online sparse regression with confidence set construction to guide exploration and planning. On the theoretical side, we provide the first regret bounds for sparse nonlinear dynamics, showing that regret scales with the sparsity level  $d_0$ . This result holds even when relaxing standard Gaussian noise assumptions by allowing a much more general, non-parametric, family of densities and when the model is misspecified. The algorithm achieving the regret bound is not computationally efficient, as it relies on a very computationally intensive online regression method. To bridge this gap, we propose a practical variant that draws inspiration from theoretical principles but incorporates more scalable components. We adopt SINDY for sparse system identification algorithm and couple it with SAC in a Dyna-style planning framework. Empirical results on classic continuous control tasks demonstrate the practical viability and robustness of our approach.

acts with a system to maximize a cumulative reward over a given number of time-steps. The main challenge to perform this task, which makes it fundamentally different from the bandit problem (Lattimore and Szepesvári, 2020), is that actions can influence the state of the system according to an unknown *transition function*. Algorithms that try to estimate this transition function to reduce RL to the planning problem are called *model-based*. Model-based Reinforcement Learning and classical continuous control share a fundamental similarity: both seek to design decision-making or feedback policies based on an internal representation of how the system evolves over time (Recht, 2019). In continuous control, one typically starts with a (potentially simplified) physics-based model of the system, whose parameters are calibrated to converge to the true model and allow for the design of an improved controller (Simpkins, 2012). In model-based RL, the approach likewise involves learning or approximating the transition function from data so that a policy can be optimized against this learned model. However, unlike continuous control, early RL research often focused on tabular or small discrete environments where learning the dynamics model from scratch was feasible without explicit structural assumptions. This gap between the two communities has been greatly reduced recently due to the ever-increasing interest in high-dimensional environments for RL, which has led to the development of model-based algorithms that go well beyond the tabular setting (Chua et al., 2018). Estimating the transition function in these complex domains requires an overwhelming sample complexity if done as in tabular environments. This motivates the need for system identification methods that are both parsimonious, remaining statistically tractable, and powerful, to achieve a satisfactory approximation of the true dynamics.

A noticeable example of parsimonious modeling in the control theory community is the Sparse Identification of Nonlinear Dynamics (SINDY) framework (Brunton et al., 2016). SINDY states that the dynamics of a complex system can often be accurately represented by a small number of terms coming from a redundant feature map (e.g., polynomials, trigonometric terms). In formulas, this corresponds to assuming that the next state  $s_{h+1}$  of the system behaves ac-

## 1 INTRODUCTION

*Reinforcement learning* (RL) (Sutton and Barto, 2018) is a paradigm of artificial intelligence in which the agent inter-

ording to

$$s_{h+1} = W_h^* \phi_h(s_h, a_h) + \eta_h, \quad (1)$$

for a noise variable  $\eta_h$ , where matrix  $W_h^*$  is *sparse*, so that only a few elements of the dictionary (often called *feature map*)  $\phi_h(\cdot)$  are in fact relevant. This framework has proven remarkably effective in domains ranging from fluid mechanics (Proctor et al., 2016) to biological systems (Champion et al., 2019). By leveraging sparsity, SINDY provides enhanced interpretability and robustness, traits that are increasingly valued in control applications.

**Contribution.** In this work, we analyze in an online model-based RL setting, environments where the system dynamics is assumed to be sparse with respect to a big feature map. Our contribution is both theoretical and practical. We first present in Section 3 Dynamic Planner with Confidence-Based Estimation (DPCBE) an optimistic algorithm for which we provide a regret bound in case of sparse nonlinear dynamics. The regret shows how the bound scales with sparsity rather than the ambient dimensionality. The analysis is further extended to general sub-gaussian noise and approximately sparse models. Based on the intuition of DPCBE we present in Section 4 Optimistic Simulation with Confidence-bAll Regression (OSCAR), a practical solution that uses SINDY for sparse system identification and couple it with Soft Actor-Critic (SAC) in a Dyna-style loop. While lacking formal guarantees, this approach inherits DPCBE insights showing optimal performance on standard control benchmarks.

**Paper structure.** We first introduce the problem setting and our core assumptions in Section 2. We then present a theoretically grounded algorithm with regret guarantees under sparsity in Section 3. This is followed by an extension to misspecified models. Finally, we describe a practical variant based on SINDY and SAC, and validate our approach through experiments on standard control benchmarks in Section 4. Finally, we report the related works in Section 5.

## 2 SETTING

We consider the setting of online optimization of a *sparse* dynamical system. Given  $\mathcal{S} \subset \mathbb{R}^{p_S}, \mathcal{A} \subset \mathbb{R}^{p_A}$  (we set  $p = p_S + p_A$ ) the state/action spaces respectively, the system evolves as

$$s_{h+1} = f_h(s_h, a_h) + \eta_h, \quad (2)$$

where  $f_h : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  is an arbitrary, possibly non-linear function and  $\eta_h$  is random noise with independent realizations at each time-step  $h$ . The interaction with this system is divided into  $K \in \mathbb{N}$  episodes, with each episode lasting for a total of  $H \in \mathbb{N}$  time steps. In each episode, the goal of the agent is to choose actions  $\{a_h\}_{h=1}^H$  in order to maximize a cumulative reward function, given by  $\sum_{h=1}^H r_h(s_h, a_h), r_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ .

**Value functions.** To better formalize this goal, we call *policy* the selection rule according to which the agent selects their actions at any time step of any episode. Formally, at each episode  $k \in [K] := \{1, \dots, K\}$ , the agent chooses a policy  $\pi^k = \{\pi_h^k\}_{h=1}^H$ , which is a sequence of mappings from  $\mathcal{S}$  to the probability distributions over  $\mathcal{A}$ . For each stage  $h \in [H]$ , the action is chosen according to  $a_h \sim \pi_h^k(\cdot | s_h)$ , the agent gains reward  $r_h(s_h, a_h)$  and the environment transitions according to Equation 2. To measure the expected reward that a policy can achieve, we define the *state value function*, i.e., the function  $V_h^\pi(s) := \mathbb{E}_\pi \left[ \sum_{\ell=h}^H r_\ell(s_\ell, a_\ell) \middle| s_h = s \right]$ . This will be compared to the *optimal state value function*:  $V_h^*(s) = \sup_\pi V_h^\pi(s)$ . Similarly, we will denote with  $Q_h^\pi(s, a)/Q_h^*(s, a)$  the state-action counterparts, where one fixes both the initial state and the initial action.

**Regret.** As a measure of performance of an agent across the  $K$  episodes of interaction with the environment, we define the *regret* as  $R_K := \sum_{k=1}^K V_1^*(s_1^k) - V_1^{\pi^k}(s_1^k)$ , where  $s_h^k$  denotes the state visited at step  $h$  of episode  $k$ . Note that if  $R_K = o(K)$  with some probability, then the value function of the "average played policy" will tend to that of the optimal policy as  $K \rightarrow \infty$ . As the value function is the expected value of the cumulative reward, this corresponds to learning the optimal behavior in the environment.

**Environment transition and sparsity.** To propose an algorithm that is able to deal with an environment of this form, we need some assumptions on its transition (Equation 2).

**Assumption 2.1.** (*Sparse feature map decomposition*) We assume that

$$f_h(s, a) = W_h^* \phi_h(s, a),$$

for some feature  $\phi_h : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$  and matrix  $W_h^* \in \mathbb{R}^{p_S \times d}$ . Moreover, each row of  $W_h^*$  has at most  $d_0$  non-zero entries, with  $p_S \ll d_0 \ll d$ .

This assumption is the core of the paper, as it represents what is usually meant by a sparse dynamical system (Brunton et al., 2016). Moreover, the assumption that  $p_S \ll d_0 \ll d$  holds whenever the feature map is complex enough to capture general classes of nonlinear functions in  $s, a$ . Indeed, polynomials/trigonometric functions of degree  $N$ , which have been used as feature maps for nonlinear continuous RL for their approximation power (Maran et al. (2024a) and Maran et al. (2024b) respectively), have a length scaling as  $d = \binom{N+p}{p} \approx N^p$ . Instead,  $d_0$  is usually smaller than  $d$ , but just by a polynomial factor (e.g.  $d \approx d_0^2$ ).

**SINDY.** SINDY is a model discovery framework that identifies the governing equations of a dynamical system directly from measured data. Specifically, SINDY can extract differential equations or data given a collection of possible features  $\phi_h$ . SINDY sparsity refers to selecting a minimal

set of functions to represent the dynamics, while its non-linearity indicates that these functions can include nonlinear combinations of state variables. The sparsity assumption (Assumption 2.1) motivates the choice of SINDY for the identification of the dynamics, since it relies on the Sequentially Threshold Least Squares with RIDGE (STLRidge) algorithm which iteratively performs RIDGE regression forcing to zero coefficients below a certain threshold. This process is repeated until convergence, progressively refining the model by discarding less relevant terms. To improve robustness in the presence of noise and enhance performance on real-world data, Ensemble-SINDY (Fasel et al., 2022, E-SINDY) has been proposed as an extension of the standard algorithm.

### 3 ALGORITHM: THEORETICAL FOUNDATION

As we said in the introduction, the central challenge in model-based RL is to balance accurate system identification with effective control. On the one hand, a good estimation for matrix  $W_h^*$  (Assumption 2.1) to predict the dynamics of the system. On the other hand, only focusing on estimating the transition function may lead to a bad regret, as it ignores the reward function. In our approach, we leverage the sparsity on the principle of *optimism in the face of uncertainty* (Auer et al., 2008) to balance these two needs. We are going to use some method from online regression (Strehl and Littman, 2007) to maintain a confidence set in the parameter space, then plan optimistically within that set. This idea draws on principles from both sparse regression and optimistic exploration in model-based RL. Below, we give a high-level overview of the three main ingredients needed to build our algorithm.

First, we need a regression method that, from collected transitions  $(s_h^k, a_h^k, s_{h+1}^k)$ , and through an appropriate sparsity-enforcing regularization, yields an estimate  $\hat{s}_{h+1}^k$  of  $\mathbb{E}[s_{h+1}^k]$ . Crucially, this algorithm must provide a theoretical guarantee that exploits the sparsity of the problem, so that assumption 2.1 is used at its full. Second, we need a confidence ball that contains  $W_h^*$  with high probability. To do so, we can use the predictions  $\hat{s}_{h+1}^k$  by the regression algorithm in order to build an estimation  $\widehat{W}_h$  of the transition matrix, and the statistical complexity guarantee of the same algorithm prove that  $\|\widehat{W}_h - W_h^*\| \leq \beta(\delta)$ , where  $\delta$  is the error probability and  $\|\cdot\|$  is some norm. In this way, we can define our confidence sets as  $\{W : \|\widehat{W}_h - W\| \leq \beta(\delta)\}$ . Lastly, once the confidence sets for the transition are fixed, we use the reward, which we assume to be known, to plan the best actions. In this way, we can ensure that (as long as the confidence sets hold), the agent plays with the best policy according to the most optimistic transition function in the confidence balls. In this way, in any episode, we either achieve a good return (that means the model is accu-

rate) or we gather information allowing to reduce the size of the confidence balls. The first critical part is to provide a regression method that can deal with the sparse regression case.

**Regression algorithm.** Selecting an appropriate regression algorithm is non-trivial. Indeed, the need to work for an arbitrary sequence of input values that may depend on past realizations of the noise (as they depend on the state and the actions) rules out standard Ridge Regression approaches. We adopt algorithm SEQSEW, originally introduced in (Gerchinovitz, 2013). While we do not delve into its intricate structure here, we highlight its regret bound. This result, which belongs to the original paper but was restated in a more convenient form by (Lattimore and Szepesvári, 2020) (Theorem 23.6), ensures what follows

**Theorem 3.1.** (Theorem 23.6 from Lattimore and Szepesvári (2020)) *Let  $\mathbf{x}_t \in \mathbb{R}^d, y_t \in \mathbb{R}$  be arbitrary sequences such that  $\|\mathbf{x}_t\|_2 \leq X, |y_t| \leq Y$  for  $t = 1, \dots, T$ . Let  $\hat{y}_t$  be the output of SEQSEW at step  $t$  (therefore, knowing  $\mathbf{x}_t$  and all previous data for  $\tau < t$ ). Then, for any  $\theta \in \mathbb{R}^d$  with  $\|\theta\|_0 = \theta_0$ , we have*

$$\begin{aligned} \rho_T(\theta) &:= \sum_{t=1}^T (y_t - \hat{y}_t)^2 - (y_t - \mathbf{x}_t^\top \theta)^2 \\ &\leq \rho_{\max} = \tilde{\mathcal{O}}(X^2 \theta_0), \end{aligned} \quad (3)$$

where  $\tilde{\mathcal{O}}$  hides logarithmic factors in  $X, Y, T, \theta_0, \|\theta\|_2$ .

In our case, the pairs  $\{\mathbf{x}_t, y_t\}$  correspond to  $\mathbf{x}_t = \phi_h(s_h^k, a_h^k)$ , and  $y_t = s_{h+1, \iota}^k$ , where  $\iota$  defines the  $\iota$ -th component of the vector. This is done since  $s_{h+1}^k$  is multivariate and we treat its prediction as  $p_S$  single regressions. Our algorithm is going to use these predictions by SEQSEW in order to build some confidence balls for the parameters.

**Confidence ball.** The result from Theorem 3.1 can be easily transformed into the definition of a confidence ball.

**Theorem 3.2.** *Fix an error probability  $\delta > 0$ ,  $\iota \in \{1, \dots, p_S\}$  and  $h \in \{1, \dots, H\}$ . Let  $\hat{y}_{h, \iota}^k$  be the output of SEQSEW for the sequences  $\phi_h(s_h^k, a_h^k), s_{h+1, \iota}^k$ , that are indexed by  $k$ . Define the confidence balls*

$$\begin{aligned} \mathcal{C}_{h, \iota}^k &:= \left\{ \theta \in \mathbb{R}^d : \sum_{\tau=1}^k (\hat{y}_{h, \iota}^\tau - \phi_h(s_h^\tau, a_h^\tau)^\top \theta)^2 \right. \\ &\quad \left. \leq \beta(\delta) + \|W_{h, \iota}^*\|_2^2 \right\} \end{aligned} \quad (5)$$

where  $W_{h, \iota}^*$  is the  $\iota$ -th row of matrix  $W_h^*$ , and  $\beta(\delta) = 1 + 2\rho_{\max} + 32 \log\left(\frac{\sqrt{8} + \sqrt{1 + \rho_{\max}}}{\delta}\right) = \tilde{\mathcal{O}}(\rho_{\max} + \log(1/\delta))$  (where  $\rho_{\max}$  comes from equation 4). Then,

$$\mathbb{P}(\exists k : W_{h, \iota}^* \notin \mathcal{C}_{h, \iota}^k) \leq \delta.$$

For the proof, see the Appendix C. This confidence set satisfies a very useful property, which is going to be one of the main backbones of the regret bound:

**Lemma 3.1.** *For every  $h, k, \iota$ , let  $\Lambda_h^k := I + \sum_{\tau=1}^k \phi_h(s_h^\tau, a_h^\tau) \phi_h(s_h^\tau, a_h^\tau)^\top$ , where  $I$  stands for the identity matrix of size  $d$ . Moreover, let  $\widehat{W}_{h,\iota}^k := (\Lambda_h^k)^{-1} \sum_{\tau=1}^k \widehat{y}_{h,\iota}^k \phi_h(s_h^\tau, a_h^\tau)$ . Then,*

$$\mathcal{C}_{h,\iota}^k \subseteq \left\{ \theta \in \mathbb{R}^d : \|\theta - \widehat{W}_{h,\iota}^k\|_{\Lambda_h^k}^2 \leq \beta(\delta) + \|W_{h,\iota}^*\|_2^2 \right\}.$$

The previous lemma allows us to give a much better idea of what the confidence ball looks like. In fact, if we see  $\widehat{W}_{h,\iota}^k$  as a replacement for the standard OLS solution, the result of Lemma 3.1 is similar to standard confidence bound for linear bandits/MDPs (Jin et al., 2020; Abbasi-Yadkori et al., 2011), just that it scales with the sparsity parameter  $d_0$ , instead of  $d$ .

**Planning.** To design Algorithm 1, the only missing part is the choice of the policy  $\pi_h^k$  to be played at stage  $h$  of episode  $k$ . Here, the idea is to leverage our construction of the confidence balls to build an optimistic estimate of the  $Q_h^*$ -function of the MDP, which we can obtain iteratively at each step  $h$ , starting from  $h = H$ . This procedure builds on the well-known Value Iteration VI algorithm (Sutton and Barto, 2018). In fact, if the matrix  $W^*$  were known, we could explicitly implement VI as follows:

$$V_{h+1}^*(s) = \max_{a \in \mathcal{A}} Q_{h+1}^*(s, a) \quad (6)$$

$$Q_h^*(s, a) = r_h(s, a) + \quad (7)$$

$$+ \underbrace{\int_{\mathcal{S}} V_{h+1}^*(W_h^* \phi_h(s, a) + \eta) d\mu_{h+1}(\eta)}_{\mathbb{E}[V_{h+1}^k(s')]} \quad (8)$$

This procedure cannot be directly implemented, as we do not know  $W^*$ . Still, knowing that  $W^* \in \mathcal{C}_h^k$  with high probability, we can perform the same algorithm with respect to the most promising matrix  $W \in \mathcal{C}_h^k$ . This is precisely what we do in lines 11 and 12 of Algorithm 1. The only difference w.r.t. VI (Equation 8) is in line 12 where we replace the  $\mathbb{E}[V_{h+1}^k(s')]$  with the expected next state under  $W$ , that is  $\mathbb{E}[V_{h+1}^k(s')|W] = \int_{\mathcal{S}} V_{h+1}^k(W \phi_h(s, a) + \eta) d\mu_{h+1}(\eta)$ . Then, to ensure optimism, we maximize over  $W \in \mathcal{C}_h^k$ , which, as we said, contains the true value  $W^*$  with high probability. In this way, we are able to ensure optimism, as we prove in the appendix C.

We now explain in more depth the structure of Algorithm 1, which combines the previous three steps. The first lines, up to 4 initialize our confidence balls and online regression algorithms SEQSEW (Gerchinovitz, 2013), one for each time-step  $h$  and dimension  $\iota$ . Then, we start iterating over the episode and, we said, we compute the current policy using VI over the most optimistic model  $W$  in lines 11 and

---

**Algorithm 1** DYNAMIC PLANNER WITH CONFIDENCE-BASED ESTIMATION (DPCBE)
 

---

**Require:** Upper bounds  $X, Y$ , Initial radius  $R$ , Reward function  $r_h$ , Failure probability  $\delta$

- 1: **for**  $\iota = 1 \dots p_S$  and  $h = 1, \dots, H$  **do**
- 2:     Initialize  $\mathcal{C}_{h,\iota}^1 := R \cdot B_1(0)$
- 3:     Initialize  $\text{ONLINELEARNER}[h, \iota] \leftarrow \text{SEQSEW}(X, Y)$
- 4: **end for**
- 5: **for**  $k = 1, 2, \dots, K$  **do**
- 6:     **for**  $h = 1, 2, \dots, H$  **do**
- 7:          $\mathcal{C}_h^k = \times_{\iota \in [p_S]} \mathcal{C}_{h,\iota}^k$
- 8:     **end for**
- 9:      $Q_H^k(s, a) = r_H(s, a)$
- 10:     **for**  $h = H - 1, \dots, 1$  **do**
- 11:          $V_{h+1}^k(s) = \arg \max_{a \in \mathcal{A}} Q_{h+1}^k(s, a)$
- 12:          $Q_h^k(s, a) = \max_{W \in \mathcal{C}_{h+1}^k} r_h(s, a) + \mathbb{E}[V_{h+1}^k(s')|W]$
- 13:          $\pi_h^k(s) = \arg \max_{a \in \mathcal{A}} Q_h^k(s, a)$
- 14:     **end for**
- 15:     Receive initial state  $s_1^k$
- 16:     **for**  $h = 1, 2, \dots, H - 1$  **do**
- 17:         Choose action  $a_h^k \sim \pi_h^k(\cdot|s_h^k)$
- 18:         **for**  $\iota = 1, \dots, p_S$  **do**
- 19:             Input  $\phi_h(s_h^k, a_h^k)$  to  $\text{ONLINELEARNER}[h, \iota]$ , receiving  $\widehat{s}_{h+1,\iota}^k$
- 20:             Update  $\mathcal{C}_{h,\iota}^k$  with  $\widehat{s}_{h+1,\iota}^k$  according to equation (5).
- 21:         **end for**
- 22:         Receive next state from the environment  $s_{h+1}^k$
- 23:         **for**  $\iota = 1, \dots, p_S$  **do**
- 24:             Update  $\text{ONLINELEARNER}[h, \iota]$  with  $s_{h+1,\iota}^k$
- 25:         **end for**
- 26:     **end for**
- 27: **end for**

---

12. This policy is then run over the environment (line 17). Here, there is an additional step compared to standard RL algorithms: once chosen the action, we query the online learner (line 19) in order to get a prediction that we are then using to update the confidence ball (line 20). After this, we receive the next state (line 22) from the environment. This state is *not* directly used to update the confidence balls, as it usually happens, but to update the online learners (line 24).

In fact, our policy is a function of confidence balls, which are only updated based on data that we collect from the online learner. The online learner is the only element that directly uses data from the observed transition. As we shall see, this multi-layered structure is the key to achieving our regret bound in the next section.

**Regret bound.** In this section, we present the guarantees for our algorithm. We start from a very general formulation that holds in a wide generality of environments, and then specialize to a few cases of interest. Apart from Assumption 2.1, we have to make an assumption about the distribution of the noise  $\eta_h$  that influences the transition at any time step. Usually, it is assumed that this noise follows from a multivariate Gaussian distribution (Kakade et al., 2020). Here, we relax this assumption by introducing a general class of distributions for the noise.

**Assumption 3.3.** For any  $h$ , the noise  $\eta_h$  is  $\sigma$ -subgaussian. Moreover, it follows a distribution in  $\mathbb{R}^{p_S}$  with density function  $\mu_h(\cdot)$  such that  $\mu_h(\cdot) \in BV(\chi_h)$ , where  $BV(\chi) := \{\|\nabla\mu(\cdot)\|_{\mathcal{M}} \leq \chi\}$ .

Here,  $\nabla$  stands for the (weak) gradient operator, and the norm  $\|\cdot\|_{\mathcal{M}}$  is defined, for any function  $g: \mathbb{R}^p \rightarrow \mathbb{R}^p$ , as  $\|g\|_{\mathcal{M}} = \sup_{f \in \mathcal{C}(\mathbb{R}^p, \mathbb{R}^p)} \int_{\mathbb{R}^p} \langle f(x), g(x) \rangle dx$ . The generality of this call is going to be explained in the next paragraph. For now, we just give the regret bound under this assumption.

**Theorem 3.4.** Assume 2.1 and 3.3. Let  $\chi := \max\{\chi_h : h = 1, \dots, H-1\}$ . Then, with probability at least  $1 - \delta$ , Algorithm 1 achieves the following regret bound

$$R_K \leq \tilde{\mathcal{O}}(H^2 \chi p_S \sqrt{d_0 d K}).$$

The most interesting parameters that emerge from this regret bound are  $d$ ,  $d_0$ , and  $\chi$ , as  $p_S \ll d_0$  and  $H, K$  are standard in the theoretical RL literature. As we can see, even if the system is sparse, the dependency on  $d$  still appears in the regret bound, even if under a square root. This is still an improvement of the state of the art, where the dependency was linear, and cannot be further improved, even in the case of sparse linear bandits (see Section 23.3 in Lattimore and Szepesvári (2020)). In fact, our bound scales as the one by (Abbasi-Yadkori et al., 2012), which works with sparse linear bandits. The other very interesting parameter that appears is  $\chi$ , which is strictly linked to Assumption 3.3 and represents the maximal total variation of the density function of the noise. Even if it is formally difficult to characterize this parameter in terms of something more interpretable, the next paragraph shows that it is bounded in most real cases.

**Bounding  $\chi$  for common noise classes.** Before this work, it was known that a) in case of *Gaussian noise*, a regret bound scaling polynomially in  $H$  in a model with transition of the form  $s_{h+1} = f(s_h, a_h) + \eta$  was possible (Kakade et al., 2020). b) In general, MDPs with this form may suffer an *exponential lower bound* in  $H$  (Maran et al. (2024a) Theorem 2). In particular, the lower bound b) leverages a Dirac-type noise at some time steps, which is highly concentrated on a discrete number of points. We can see our Assumption 3.3 as a way to prevent this scenario to happen and see  $\chi = \max_h \|\mu_h\|_{BV}$  as a measure of how much the density function of the noise concentrates on one or more single points. In fact, we have, fixing  $\mathcal{S} \subset \mathbb{R}$ , for clarity,

- $\mu(s) = \delta_0(s)$  leads to  $\chi = +\infty$  (Dirac's delta is not in BV).
- $\mu(s) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-s^2/(2\sigma^2))$  leads to  $\chi = \frac{2}{\sqrt{2\pi\sigma^2}}$  (Gaussian noise is in BV).
- $\mu(s) = \mathbf{1}_{(-a,a)}(s)$  leads to  $\chi = 1/a$  (uniform noise is in BV).

- Every bounded  $\mu(\cdot)$  with at most  $m$  local maxima is in BV with  $\chi \leq (m+1) \sup_{s \in \mathcal{S}} \mu(s)$ .

Moreover, one crucial property of this class of density functions, which makes it very useful, is that it is closed under convex combinations, as results from the following proposition.

**Proposition 3.5.** For any pair of density functions  $\mu(\cdot) \in BV(\chi)$ ,  $\tilde{\mu}(\cdot) \in BV(\tilde{\chi})$  and any constant  $\lambda \in [0, 1]$ , their convex combination satisfies

$$\lambda\mu(\cdot) + (1-\lambda)\tilde{\mu}(\cdot) \in BV(\lambda\chi + (1-\lambda)\tilde{\chi}).$$

This shows that assuming 3.3 increases dramatically the level of generality: indeed, any Gaussian Mixture is included in the BV class, which makes this class extremely flexible. We now investigate how a small change in the algorithm can lead to a generalization of Theorem 3.4 to the case where the model (described in Section 2) is misspecified.

### 3.1 Extension: misspecified models

For the sake of this section, we relax assumption 2.1 in the following way:

**Assumption 3.6.** (Sparse misspecified feature map decomposition) For some feature map and  $\phi_h: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$  matrix  $W_h^* \in \mathbb{R}^{p_S \times d}$  such that each row of  $W_h^*$  has at most  $d_0$  non-zero entries, with  $p_S \ll d_0 \ll d$ , let  $\xi_h: \mathcal{S} \times \mathcal{A} \rightarrow [0, \infty)$  be the function

$$\xi_h(\cdot) := \|W_h^* \phi_h(\cdot) - f_h(\cdot)\|_{L^\infty}.$$

We assume that, for any  $h = 1, \dots, H$ ,  $\xi_h(\cdot)$  is uniformly bounded by a constant  $\xi_\infty$ .

We can see that Assumption 2.1 corresponds to Assumption 3.6 for  $\xi_\infty = 0$ . As the next result shows, a modification of Algorithm 1 allows us to generalize its regret bound to this, more challenging, setting.

**Theorem 3.7.** Assume 3.6 and 3.3. Let  $\chi := \max\{\chi_h : h = 1, \dots, H-1\}$ . Let us modify the construction of the confidence sets in Equation 5 by using as confidence radius  $\beta'(\delta) = \beta(\delta) + K\xi_\infty^2 \log(1/\delta)$ . Then, with probability at least  $1 - \delta$ , algorithm 1 achieves the following regret bound

$$R_K \leq \tilde{\mathcal{O}}\left(H^2 \chi p_S \sqrt{d_0 d K} + H^2 p_S \sqrt{d K} \xi_\infty\right).$$

The new regret bound contains a first term  $H^2 \chi p_S \sqrt{d_0 d K}$  which is identical to the previous regret bound. The dependence on the misspecification is all captured by the other term,  $H^2 p_S \sqrt{d K} \xi_\infty$ , where some things need to be commented on. First, the linear growth in  $K$  is expected: for big  $\xi_\infty$ , the regret grows linearly, as the approximation power

of our feature map is very limited. What is more concerning is the multiplicative factor  $\sqrt{d}$ , which creates a trade-off: bigger  $d$  may reduce the misspecification, but if the reduction is not sufficient, the term  $\sqrt{d}\xi_\infty$  gets bigger, actually. This potentially pernicious phenomenon affects many other algorithms for Bandits/Reinforcement Learning with linear function approximation, and cannot be avoided in general (Lattimore et al., 2020).

## 4 ALGORITHM: PRACTICAL VARIANT

The algorithmic solution presented in Section 3 relies on the usage of SEQSEW (Gerchinovitz, 2013) as a regressor to build the confidence balls over the transition parameters. However, a computationally feasible version of SEQSEW is still an open problem (see Lattimore and Szepesvári (2020) 23.5 Note 1). To overcome this limitation, we introduce Optimistic Simulation with Confidence-bAll Regression (OSCAR), a practical algorithm inspired by the theoretical principles of DPCBE, but adapted to be computationally efficient.

**From Theory to Practice.** OSCAR maintains the key idea of optimistic planning over confidence sets but replaces SEQSEW with the more computationally efficient SINDY as a regressor. Based on SINDY estimates, confidence balls are then constructed using the same formula as in the theoretical algorithm. The use of confidence balls is the main difference from our approach and the other adaptations of SINDY that were proposed in the RL literature. Although performing explicit optimistic planning within these sets would be computationally prohibitive, sampling models from them provides an ensemble of reasonable dynamics models. This approximate form of optimistic planning within the confidence ball encourages OSCAR to explore regions where uncertainty is highest.

### Algorithm Structure.

We now explain the structure of Algorithm 2 in more detail. Initially, a random policy is exploited to explore the environment and collect useful data for the initialization of the SINDY model. For each state-action pair collected from the environment, the next state prediction is used to update the confidence ball according to Lemma 3.1. The center of the confidence ball  $\widehat{W}$  can be computed online, initializing the matrix  $\Lambda$  and a vector  $P$  to zero value. For each new tuple state, action, and next state, they are updated with the following rules:

$$\Lambda_{h+1} \leftarrow \Lambda_h + \phi_h(s, a)\phi_h(s, a)^T \quad (9)$$

$$P_{h+1} = P_h + s' \phi_h(s, a) \quad (10)$$

So that each row of the matrix  $\widehat{W}$  can be computed as  $\widehat{W}_{h,\cdot} = (\Lambda_{h+1} + \alpha \cdot I)^{-1}P$ . with  $\alpha$  being an hyperparameter for regularization of the SINDY model and  $I$  the

---

### Algorithm 2 OPTIMISTIC SIMULATION WITH CONFIDENCE-BALL REGRESSION (OSCAR)

---

**Require:** Regularization parameter  $\alpha$ , Reward function  $r$ , Basis Functions  $\phi$ , Batch size  $b$ , Horizon  $H$ , Number of off-policy iterations  $N_e$ , Number of on-policy iterations  $N_o$ , Number of model interactions  $N_m$ , Policy  $\pi(\cdot|s_w)$ , Ball scale factor  $\lambda$

- 1: Initialize dataset  $D_e \leftarrow \emptyset$
- 2: Collect trajectories in  $D_e$  using Random Policy for  $N_e$  episodes
- 3:  $SINDy(D_e, \phi)$  ▷ Fit SINDy Model
- 4: **for** each  $(s, a)$  in  $D_e$  **do** ▷ Set Confidence Ball
- 5:     Predict  $s' = SINDy(s, a)$
- 6:     Update  $\mathcal{C}$  with  $s'$  according to equation (5)
- 7: **end for**
- 8: **for**  $n = N_e + 1, \dots, N_o$  **do** ▷ Start On Policy Interaction
- 9:     **for**  $h = 1, 2, \dots, H$  **do**
- 10:         Observe state  $s_h$
- 11:         Sample action  $a_h \sim \pi(\cdot|s_h^{\widehat{W}})$
- 12:         Predict  $s_{h+1} = SINDy(s_h, a_h)$
- 13:         Update  $\mathcal{C}$  with  $s_{h+1}$  according to equation (5)
- 14:         Observe next state  $s_{h+1}$
- 15:         Update  $D_e \leftarrow D_e \cup (s_h, a_h, s_{h+1})$
- 16:     **end for**
- 17:     Randomly sample a batch of  $b$  states  $S = \{s\} \sim D_e$  ▷ Start Model Interaction and Policy Update
- 18:     **for**  $i = 1, 2, \dots, N_m$  **do**
- 19:         Uniformly sample  $b$  scale factors  $\beta \sim U[0, \lambda]$
- 20:         Sample  $b$  models  $M \sim N(\widehat{W}, \beta(\Lambda + \alpha \cdot I)^{-1})$
- 21:         Uniformly sample  $b$  scale factors  $\beta' \sim U[0, \lambda]$
- 22:         Sample  $b$  models  $M' \sim N(\widehat{W}, \beta'(\Lambda + \alpha \cdot I)^{-1})$
- 23:         Initialize  $A \leftarrow \emptyset, S' \leftarrow \emptyset, K \leftarrow \emptyset$
- 24:         **for** each  $(s, m, m')$  in  $S \times M \times M'$  **do**
- 25:             Sample action  $a \sim \pi(\cdot|s^{m'})$
- 26:             Observe  $s' = m \cdot \phi(s, a)$
- 27:             Collect  $k = r(s, a)$
- 28:             Update  $A \leftarrow A \cup a, S' \leftarrow S' \cup s', K \leftarrow K \cup k$
- 29:         **end for**
- 30:         Perform a SAC update using batch  $(S \times M, A, K, S' \times M)$
- 31:     **end for**
- 32:      $SINDy(D_e, \phi)$  ▷ Update SINDy Model
- 33: **end for**

---

identity matrix. After this setup phase, the main algorithm follows a Dyna-style learning approach, alternating between interactions with the real and surrogate environment, where the real environment is used only for data gathering.

The original algorithm involves a nested optimization problem over actions and models (line 12). To address this in practice, we define an extended state  $s_w$  which includes both the observed environment state and the dynamical system parameters, enabling evaluation across all models within the confidence ball. To evaluate the *state-model-action* space, we adopt a Dyna-style approach with SAC (Haarnoja et al., 2018), considering models randomly sampled from the confidence ball at each step. Model sampling is performed through *Thompson Sampling* with posterior reshaping, drawing from a multivariate normal distribution  $N(W^*, \beta(\Lambda + \alpha \cdot I)^{-1})$  with  $\beta \sim U[0, \lambda]$  and with  $\lambda$  being a hyperparameter controlling the maximum scale of the

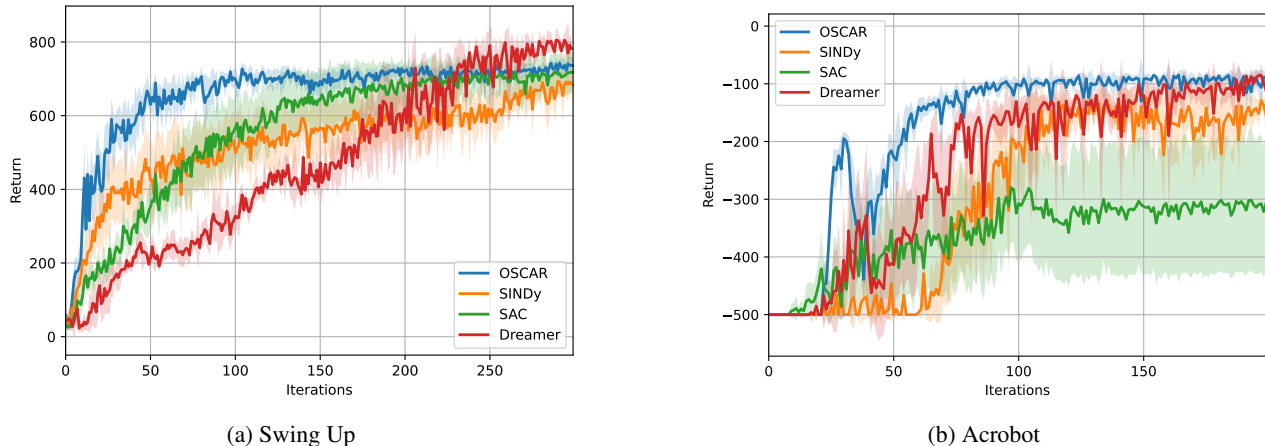


Figure 1: Learning Curves for Swing up (a) and Acrobot (b)

covariance matrix  $\Lambda$  to sample inside the confidence ball uniformly. To avoid storing model parameters in the replay buffer, SAC updates are performed online by simultaneously collecting batches of trajectories. For this reason, the neural networks of SAC are extended, including layer normalization as done in PQN (Gallici et al., 2025) to stabilize learning. Model exploration is done by randomly sampling a vector of models from the confidence ball at each step, with each sampled state evaluated under its corresponding model. Action exploration is guided not only by SAC but also by enforcing actions that are optimal in similar models, i.e., taking the optimal actions for the sampled states under different models randomly sampled from the confidence ball. Finally, we collect the reward, update the policy based on the SAC algorithm, and update the SINDY model accordingly. This approach can also be parallelized to reduce computing time, considering the set of random models as a vectorized environment.

#### 4.1 Experiments

In this section, we present the results obtained by adopting our approach over different RL benchmarks. We conducted experiments using three environments of increasing complexity from the Gymnasium suite (Towers et al., 2024), Acrobot and Continuous Mountain Car, and the dm\_control Swing Up (Tassa et al., 2018). We evaluate the algorithm by studying the effects of the confidence ball with respect to classical Soft Actor Critic, model-based Soft Actor Critic with SINDY and the Dreamer v3 algorithm (Hafner et al., 2023). In each environment, the confidence ball dimension is obtained scaling  $(\Lambda + \alpha \cdot I)^{-1}$  for a fixed hyperparameter  $\lambda$ . All the results presented are averaged over 10 different seeds. Appendix D provides additional experimental details and results. The code to reproduce the experiments can be found at <https://github.com/enegusme/OSCAR>.

**Swing-Up.** The Cartpole Swing-Up task involves swinging up and balancing a pole attached to a cart that moves horizontally. The state includes the cart position  $x$ , the cosine and sine of the pole angle  $\theta$ , the cart velocity  $\dot{x}$ , and the pole angular velocity  $\dot{\theta}$ . The reward function incentivizes the agent to keep the pole up, keeping the cart in a central position. The environment considers episodes of 1000 time steps. Initially, a single episode of data is collected for SINDY and confidence ball initialization. Subsequently, for each episode of data collection in the real environment, 4 episodes are taken in the surrogate one for a total of  $N_m = 4 \cdot 1000$ . The batch size  $b = 256$  and the scaling factor  $\lambda = 0.01$ . We can see in Figure 1a how, among the four approaches presented, OSCAR represents the most data-efficient solution, managing to converge to an almost optimal policy with significantly fewer iterations.

**Acrobot.** The Acrobot environment consists of two linked pendulum-like arms, with the second joint actuated. The goal is to swing the free end upward to reach a target height starting from the downward position. The state is represented by the cosine and sine of the two angles  $\theta_1, \theta_2$ , considering that the angle of the second joint is relative to the first one, and their angular velocities  $\dot{\theta}_1, \dot{\theta}_2$ . The agent always receives a negative reward unless the target position is reached. Each episode lasts 500 time steps. Initially, a single episode of data is collected for SINDY and confidence ball initialization. Then, for each trajectory collected in the real environment, we collect one episode in the surrogate environment. The batch size is  $b = 256$  and the scaling factor is  $\lambda = 1.0$ . Figure 1b shows the learning curves for all three approaches. OSCAR clearly outperforms the baselines, converging faster and reaching the optimal policy.

**Mountain Car.** The Mountain Car environment objective is to bring a car to the top of the hill starting from the bottom of a sinusoidal valley. The state is represented by the position of the cart  $x$  and its velocity  $\dot{x}$ . The agent receives

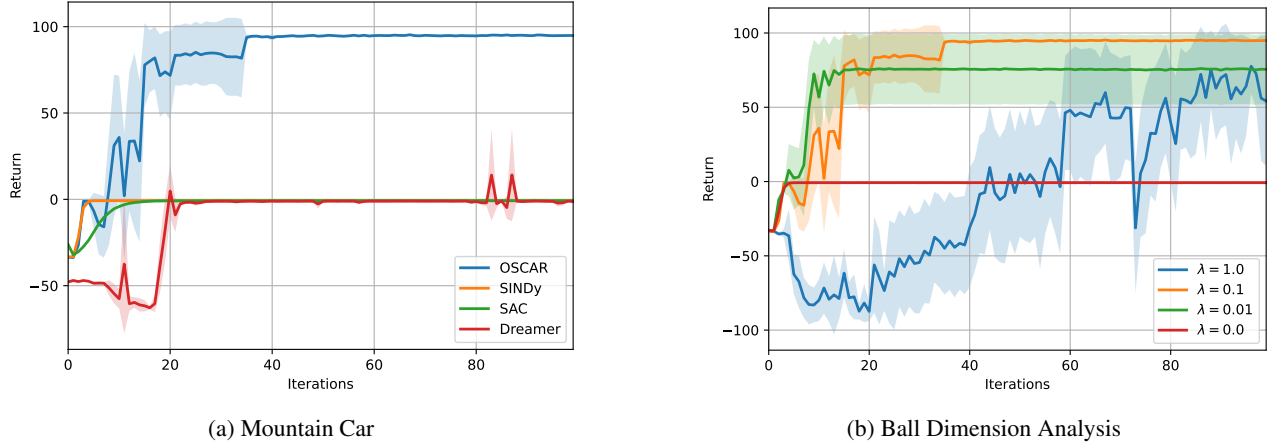


Figure 2: Learning Curves for Mountain Car

positive feedback only if the car reaches the goal position, making exploration essential to solve the task. The environment considers episodes of 1000 time steps and terminates if the goal is reached. Initially, a single episode of data is collected for SINDY and confidence ball initialization. Subsequently, for each episode of data collection in the real environment, 5 episodes are taken in the surrogate for a total of  $Nm = 5 \cdot 1000$ . The batch size  $b = 512$  and the scaling factor  $\lambda = 0.1$ . We show in Figure 2a the learning curves between the three proposed methods. We notice that OSCAR is the only approach that manages to complete the task. In environments where the exploration is crucial, the optimistic behavior of OSCAR allows for a more efficient exploration, exploiting the model ensemble. Instead, approaches like SAC, model-based SAC with SINDY or even Dreamer end up in a suboptimal behavior due to a lack of exploration. Furthermore, in Figure 2b we present an analysis over the  $\lambda$  parameter, which scales the radius of the confidence ball used. A larger  $\lambda$  corresponds to a more optimistic selection strategy, as the ball includes a wider range of models. However, we can notice that, while  $\lambda = 1.0$  incentivizes a broader exploration, it slows down the learning of the agent due to querying over uncertain models. On the other hand,  $\lambda = 0.01$  leads to a conservative behavior that limits exploration, resulting in suboptimal performance. Eventually, we notice that an intermediate,  $\lambda = 0.1$  value manages to balance the trade-off between exploration and model reliability, achieving the optimal performance.

## 5 RELATED WORK

The topics of model-based RL, as well as the optimal control of nonlinear dynamic systems, have been studied in countless papers. Therefore, we mention only the most important ones here, leaving an extensive study for the Appendix A.

**Model-based RL.** Most RL works assume transitions of the form  $s_{h+1} \sim P_h(\cdot | s_h, a_h)$ , but a large literature at the RL-control interface studies additive-noise models  $s_{h+1} = f_h(s_h, a_h) + \eta_h$ . Classical approaches almost always assume Gaussian noise, both in linear adaptive control and in nonlinear settings with Brownian or Gaussian process disturbances. More recent model-based RL methods such as PILCO (Deisenroth and Rasmussen, 2011) have reintroduced Gaussian additive noise for uncertainty propagation, but they do not leverage sparsity in the dynamics as we do (we refer to Appendix A for an extensive review of related approaches).

**SINDY.** Closest to our work are approaches combining SINDY with RL. For instance, (Arora et al., 2022) employ a Dyna-style algorithm using SINDY with Soft Actor-Critic, while (Zolman et al., 2024) integrate DRL with SINDY-guided dictionary learning. Differently from these methods, we provide a theoretical analysis with regret guarantees and design a practical algorithm that leverages confidence sets for efficient exploration, achieving improved performance.

**RL theory with linear function approximation.** From a theoretical perspective, our paper is the first to provide a regret bound under sparsity assumptions in the nonlinear setting. While regret bounds for LQR are classical, extensions to general nonlinear feature maps are relatively recent (Kakade et al., 2020). Compared to this line, we introduce three novelties: sparsity, noise with densities of bounded variation (beyond Gaussians), and model misspecification. The latter connects to a broader literature on misspecification amplification in RL with linear function approximation (Lattimore et al., 2020; Zanette et al., 2020; Weisz et al., 2023; Wu et al., 2024).

**Neuro-Dynamic programming.** A rich line of work in Adaptive and Neurodynamic Programming has studied optimal control for nonlinear systems using neural actor-critic

methods with stability guarantees (Vamvoudakis and Lewis, 2010; Lewis and Vamvoudakis, 2010; Liu et al., 2014; Bhasin et al., 2013; Jiang and Jiang, 2014). These approaches allow for general disturbances and thus share our robustness goal, but differ in focus and assumptions. They primarily establish ISS, UUB, or bounded suboptimality, often requiring persistence of excitation, affine structure, or Lipschitz error bounds. By contrast, our analysis provides regret guarantees under sparsity and non-Gaussian noise, without structural or excitation assumptions.

## 6 CONCLUSION

We have presented a unified theoretical and practical framework for RL in environments with sparse nonlinear dynamics. Our key theoretical contribution is Algorithm 1 (DPCBE), which is endowed with a regret bound that scales with sparsity, a more general noise condition and misspecification, significantly extending existing guarantees for model-based RL. To bridge theory and application, we introduced Algorithm 2 (OSCAR), a practical algorithm that uses SINDY to overcome the computational burden of DPCBE. Our experiments demonstrate that this approach not only maintains the theoretical spirit of optimism-driven exploration but also performs competitively or superiorly across classic continuous control benchmarks, especially in tasks where exploration is key.

## References

- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- Benjamin Recht. A tour of reinforcement learning: The view from continuous control. *Annual Review of Control, Robotics, and Autonomous Systems*, 2(1):253–279, 2019.
- Alex Simpkins. System identification: Theory for the user, (Ijung, I.; 1999)[on the shelf]. *IEEE Robotics & Automation Magazine*, 19(2):95–96, 2012.
- Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. *Advances in neural information processing systems*, 31, 2018.
- Steven L Brunton, Joshua L Proctor, and J Nathan Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the national academy of sciences*, 113(15):3932–3937, 2016.
- Joshua L Proctor, Steven L Brunton, and J Nathan Kutz. Dynamic mode decomposition with control. *SIAM Journal on Applied Dynamical Systems*, 15(1):142–161, 2016.
- Kathleen Champion, Bethany Lusch, J Nathan Kutz, and Steven L Brunton. Data-driven discovery of coordinates and governing equations. *Proceedings of the National Academy of Sciences*, 116(45):22445–22451, 2019.
- Davide Maran, Alberto Maria Metelli, Matteo Papini, and Marcello Restelli. No-regret reinforcement learning in smooth mdps. In *Forty-first International Conference on Machine Learning*, 2024a.
- Davide Maran, Alberto Maria Metelli, Matteo Papini, and Marcello Restelli. Projection by convolution: Optimal sample complexity for reinforcement learning in continuous-space mdps. *arXiv preprint arXiv:2405.06363*, 2024b.
- U. Fasel, J. N. Kutz, B. W. Brunton, and S. L. Brunton. Ensemble-sindy: Robust sparse model discovery in the low-data, high-noise limit, with active learning and control. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 478(2260), April 2022. ISSN 1471-2946. doi: 10.1098/rspa.2021.0904. URL <http://dx.doi.org/10.1098/rspa.2021.0904>.
- Peter Auer, Thomas Jaksch, and Ronald Ortner. Near-optimal regret bounds for reinforcement learning. *Advances in neural information processing systems*, 21, 2008.
- Alexander Strehl and Michael Littman. Online linear regression and its application to model-based reinforcement learning. *Advances in Neural Information Processing Systems*, 20, 2007.
- Sébastien Gerchinovitz. Sparsity regret bounds for individual sequences in online linear regression. *The Journal of Machine Learning Research*, 14(1):729–769, 2013.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on learning theory*, pages 2137–2143. PMLR, 2020.
- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24, 2011.
- Sham Kakade, Akshay Krishnamurthy, Kendall Lowrey, Motoya Ohnishi, and Wen Sun. Information theoretic regret bounds for online nonlinear control. *Advances in Neural Information Processing Systems*, 33:15312–15325, 2020.
- Yasin Abbasi-Yadkori, David Pal, and Csaba Szepesvari. Online-to-confidence-set conversions and application to sparse stochastic bandits. In *Artificial Intelligence and Statistics*, pages 1–9. PMLR, 2012.
- Tor Lattimore, Csaba Szepesvari, and Gellert Weisz. Learning with good feature representations in bandits and in rl with a generative model. In *International conference on machine learning*, pages 5662–5670. PMLR, 2020.

- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. Pmlr, 2018.
- Matteo Gallici, Mattie Fellows, Benjamin Ellis, Bartomeu Pou, Ivan Masmitja, Jakob Nicolaus Foerster, and Mario Martin. Simplifying deep temporal difference learning, 2025. URL <https://arxiv.org/abs/2407.04811>.
- Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U. Balis, Gianluca De Cola, Tristan Deleu, Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, Rodrigo Perez-Vicente, Andrea Pierré, Sander Schulhoff, Jun Jet Tai, Hannah Tan, and Omar G. Younis. Gymnasium: A standard interface for reinforcement learning environments, 2024. URL <https://arxiv.org/abs/2407.17032>.
- Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew LeFrancq, Timothy Lillicrap, and Martin Riedmiller. Deepmind control suite, 2018. URL <https://arxiv.org/abs/1801.00690>.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- Marc Deisenroth and Carl E Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on machine learning (ICML-11)*, pages 465–472, 2011.
- Rushiv Arora, Bruno Castro da Silva, and Eliot Moss. Model-based reinforcement learning with sindy. *arXiv preprint arXiv:2208.14501*, 2022.
- Nicholas Zolman, Urban Fasel, J Nathan Kutz, and Steven L Brunton. Sindy-rl: Interpretable and efficient model-based reinforcement learning. *arXiv preprint arXiv:2403.09110*, 2024.
- Andrea Zanette, Alessandro Lazaric, Mykel Kochenderfer, and Emma Brunskill. Learning near optimal policies with low inherent bellman error. In *International Conference on Machine Learning*, pages 10978–10989. PMLR, 2020.
- Gellért Weisz, András György, and Csaba Szepesvári. Online rl in linearly  $q^\pi$ -realizable mdps is as easy as in linear mdps if you learn what to ignore. *Advances in Neural Information Processing Systems*, 36:59172–59205, 2023.
- Runzhe Wu, Ayush Sekhari, Akshay Krishnamurthy, and Wen Sun. Computationally efficient rl under linear bellman completeness for deterministic dynamics. *arXiv preprint arXiv:2406.11810*, 2024.
- Kyriakos G Vamvoudakis and Frank L Lewis. Online actor-critic algorithm to solve the continuous-time infinite horizon optimal control problem. *Automatica*, 46(5):878–888, 2010.
- Frank L Lewis and Kyriakos G Vamvoudakis. Reinforcement learning for partially observable dynamic processes: Adaptive dynamic programming using measured output data. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 41(1):14–25, 2010.
- Derong Liu, Ding Wang, Fei-Yue Wang, Hongliang Li, and Xiong Yang. Neural-network-based online hjb solution for optimal robust guaranteed cost control of continuous-time uncertain nonlinear systems. *IEEE transactions on cybernetics*, 44(12):2834–2847, 2014.
- Shubhendu Bhasin, Rushikesh Kamalapurkar, Marcus Johnson, Kyriakos G Vamvoudakis, Frank L Lewis, and Warren E Dixon. A novel actor-critic-identifier architecture for approximate optimal control of uncertain nonlinear systems. *Automatica*, 49(1):82–92, 2013.
- Yu Jiang and Zhong-Ping Jiang. Robust adaptive dynamic programming and feedback stabilization of nonlinear systems. *IEEE Transactions on Neural Networks and Learning Systems*, 25(5):882–893, 2014.
- Steven J Bradtke, B Erik Ydstie, and Andrew G Barto. Adaptive linear quadratic control using policy iteration. In *Proceedings of 1994 American Control Conference-ACC'94*, volume 3, pages 3475–3479. IEEE, 1994.
- Yasin Abbasi-Yadkori and Csaba Szepesvári. Regret bounds for the adaptive control of linear quadratic systems. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 1–26. JMLR Workshop and Conference Proceedings, 2011.
- Kenji Doya. Reinforcement learning in continuous time and space. *Neural computation*, 12(1):219–245, 2000.
- Alex Simpkins, Raymond De Callafon, and Emanuel Todorov. Optimal trade-off between exploration and exploitation. In *2008 American control conference*, pages 33–38. IEEE, 2008.
- Christopher G Atkeson and Juan Carlos Santamaria. A comparison of direct and model-based reinforcement learning. In *Proceedings of international conference on robotics and automation*, volume 4, pages 3557–3564. IEEE, 1997.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. *arXiv preprint arXiv:2010.02193*, 2020.
- Luigi Ambrosio, Nicola Fusco, and Diego Pallara. *Functions of bounded variation and free discontinuity problems*. Oxford university press, 2000.

Davide Maran, Alberto Maria Metelli, Matteo Papini, and Marcello Restelli. Local linearity: the key for no-regret reinforcement learning in continuous mdps. *arXiv preprint arXiv:2410.24071*, 2024c.

Jialin Dong and Lin Yang. Does sparsity help in learning misspecified linear bandits? In *International Conference on Machine Learning*, pages 8317–8333. PMLR, 2023.

Brian M. de Silva, Kathleen Champion, Markus Quade, Jean-Christophe Loiseau, J. Nathan Kutz, and Steven L. Brunton. Pysindy: A python package for the sparse identification of nonlinear dynamics from data, 2020. URL <https://arxiv.org/abs/2004.08424>.

Alan Kaptanoglu, Brian de Silva, Urban Fasel, Kadierdan Kaheman, Andy Goldschmidt, Jared Callahan, Charles Delahunt, Zachary Nicolaou, Kathleen Champion, Jean-Christophe Loiseau, J. Kutz, and Steven Brunton. Pysindy: A comprehensive python package for robust sparse system identification. *Journal of Open Source Software*, 7(69):3994, January 2022. ISSN 2475-9066. doi: 10.21105/joss.03994. URL <http://dx.doi.org/10.21105/joss.03994>.

## Checklist

1. For all models and algorithms presented, check if you include:
  - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. Yes
  - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. Yes
  - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. Yes
2. For any theoretical claim, check if you include:
  - (a) Statements of the full set of assumptions of all theoretical results. Yes
  - (b) Complete proofs of all theoretical results. Yes
  - (c) Clear explanations of any assumptions. Yes
3. For all figures and tables that present empirical results, check if you include:
  - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). Yes
  - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). Yes
  - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). Yes
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
  - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). Yes
  - (a) Citations of the creator If your work uses existing assets. Yes
  - (b) The license information of the assets, if applicable. Not Applicable
  - (c) New assets either in the supplemental material or as a URL, if applicable. Not Applicable
  - (d) Information about consent from data providers/curators. Yes
  - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. Not Applicable
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
  - (a) The full text of instructions given to participants and screenshots. Not Applicable
  - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. Not Applicable
  - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. Not Applicable

---

## Supplementary Materials

---

### A Related Work

**Model-based RL.** While the majority of RL papers focus on a transition function of the form  $s_{h+1} \sim P_h(\cdot|s_h, a_h)$ , there is still a huge literature that, lying at the intersection between RL and control theory, considers systems where  $s_{h+1} = f_h(s_h, a_h) + \eta_h$ , as in this paper. Most of these works assume that the additive noise term  $\eta_h$  is Gaussian, which we do not. This assumption is classical in the control community, appearing in early works on adaptive control and reinforcement learning for linear systems with Gaussian disturbances (e.g., (Bradtke et al., 1994), (Abbasi-Yadkori and Szepesvári, 2011)). In the nonlinear setting, (Doya, 2000) and (Simpkins et al., 2008) model continuous-time systems subject to Brownian or Gaussian process noise, and derive actor-critic or Bayesian optimal policies accordingly.

From a more practical point of view, although models of this type have been known for a long time (Atkeson and Santamaria, 1997), recent model-based RL approaches "re-discovered" this assumption to capture uncertainty in learned dynamics. For instance, PILCO (Deisenroth and Rasmussen, 2011) explicitly posit  $s' \sim \mathcal{N}(\phi(s, a), \Sigma)$  and exploit the additive structure to propagate uncertainty during planning. Differently, more recent approaches focused more on learning complex latent environment models, as for instance the Dreamer family of algorithms (Hafner et al., 2019, 2020, 2023) which adopts recurrent state-space models to capture non additive stochasticity in high-dimensional settings. These approaches, however, do not typically impose or exploit sparsity in the transition model.

**SINDY.** On the point of view of the application, the works that are closer to ours are the ones that try to leverage SINDY algorithm on top of RL methods. The usage of SINDY has been paired with classical model-based paradigm. In (Arora et al., 2022), the authors propose a Dyna-style learning algorithm adopting SINDY as a method to estimate the surrogate environment paired with Soft Actor-Critic algorithm. A more complex approach was introduced in (Zolman et al., 2024), which involves the usage of Deep Reinforcement Learning (DRL) and dictionary learning guided by SINDY-based sparse identification. However, unlike the cited methodologies, our work tackle the setting also from a theoretical perspective providing regret guarantees, paired with competitive practical results. Moreover, our practical algorithm uses confidence sets over models to guide exploration, resulting in more efficient exploration and improved performance compared to the cited approaches.

**RL theory with linear function approximation.** From the point of view of theory, we are the first paper providing a regret bound in our setting under the sparsity assumption. While regret bounds for the LQR are already well-known, papers generalizing these results to an arbitrary non-linear feature map of the state-action pair (like in our case (2)) are relatively novel (Kakade et al., 2020). With respect to the latter paper, we generalize over three sides: we use a sparsity assumption, we extend the family of the noise allowed for the model, from Gaussians to any density of bounded variation, and we allow for model misspecification. Related to the first point, the work of (Abbasi-Yadkori et al., 2012) is worth mentioning, as it solves the problem of sparse linear bandits. The bounded variation assumption on the density of the noise is, to the best of our knowledge, novel in the RL setting, while this class of functions plays a central role in variational calculus and partial differential equations (Ambrosio et al., 2000). Lastly, introducing the misspecification into a linear function approximation scheme for RL is known to introduce one problem that we briefly mention in Subsection 3.1: the influence of  $\xi_\infty$  in the regret bound is multiplied by  $\sqrt{d}$ . As (Lattimore et al., 2020) argues, this problem arises from the fact that, in bandit/RL settings, the agent needs uniform control over the error of the model on the state-action space, and is not solvable in general. In fact, several recent seminal works in RL with linear function approximation (Zanette et al., 2020) (Weisz et al., 2023) (Wu et al., 2024) provide regret bounds that suffer from misspecification amplification. To deal with this problem was the main goal of multiple papers, like (Maran et al., 2024c), which removes the amplification factor in case of a locally linear feature map, and (Dong and Yang, 2023), which improves the  $\sqrt{d}$  amplification in case of sparsity, but just in case of bandits, not in RL. None of these works is able to tackle the problem in our setting, so that, at the current state-of-the-art, misspecification amplification remains one big open problem in RL with linear function approximation.

**Neuro-Dynamic programming.** A large body of research in Adaptive Dynamic Programming (ADP) and Neurodynamic Programming (NDP) has investigated optimal control for nonlinear systems, often through neural-network-based actor–critic architectures with provable stability and performance guarantees. Seminal contributions in this field include works by (Vamvoudakis and Lewis, 2010), (Lewis and Vamvoudakis, 2010), (Liu et al., 2014), (Bhasin et al., 2013) and (Jiang and Jiang, 2014) among others.

These methods are of particular interest from our perspective since they account for general stochastic disturbances, and therefore do not rely on Gaussian noise assumptions, like in our case. In this sense, they share the goal of robust learning under uncertainty, similar to our work. The difference stays mostly in the goal: these papers focus primarily on Input-to-State Stability (ISS), uniform ultimate boundedness (UUB), and bounded suboptimality of the learned policy. Moreover, the assumption made on the exploration are different from the ones of our work, and they do not assume sparsity.

In particular, (Lewis and Vamvoudakis, 2010) proposed one of the first online actor–critic algorithms for continuous-time nonlinear systems, proving uniform ultimate boundedness and near-optimality under bounded approximation errors. Their analysis, however, crucially relies on persistence of excitation or artificially injected exploration signals, whereas our regret guarantees hold without requiring such assumptions.

(Liu et al., 2014) (Bhasin et al., 2013) developed an adaptive dynamic programming method for nonlinear systems with guaranteed cost control, establishing stability and bounded suboptimality under the assumption of known affine structure and bounded approximation errors. In contrast, our work requires no structural affinity assumption and provides regret bounds, which are a stronger type of guarantee.

(Jiang and Jiang, 2014) introduced robust adaptive dynamic programming for nonlinear systems, proving input-to-state stability and robust near-optimality under bounded disturbances and Lipschitz approximation errors. Unlike their framework, our analysis does not rely on predefined Lyapunov functions or Lipschitz error bounds, as it leverages on the noise assumption to get rid of this smoothness requirement.

## B Notation

In this section, we leave, for the reader’s convenience, a table of the notation used in this paper.

$\mathcal{S}$	State space $\subset \mathbb{R}^{p_S}$
$\mathcal{A}$	Action space $\subset \mathbb{R}^{p_A}$
$H$	Time horizon
$f_h$	Transition function $\mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ at stage $h \in \{1, \dots, H\}$
$\eta_h$	Noise at step $h$
$\mu_h$	Density function of the noise at step $h$
$r_h$	Reward function $\mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$
$K$	Number of episodes
$\pi_h^k$	Policy played at stage $h$ of episode $k$
$s_h^k$	state visited at stage $h$ of episode $k$
$a_h^k$	action played at stage $h$ of episode $k$
$\Lambda_h^k$	design matrix at stage $h$ of episode $k$
$V_h^*$	Optimal value function
$R_K$	Regret
$W_h^*$	System matrix (see equation 2)
$\phi_h$	Feature map (see equation 2)
$d$	Dimension of the feature maps $\phi_h$
$d_0$	Sparsity parameter (see assumption 2.1)

## C Proofs

**Theorem 3.2.** Fix an error probability  $\delta > 0$ ,  $\iota \in \{1, \dots, p_S\}$  and  $h \in \{1, \dots, H\}$ . Let  $\hat{y}_{h,\iota}^k$  be the output of SEQSEW for the sequences  $\phi_h(s_h^k, a_h^k), s_{h+1,\iota}^k$ , that are indexed by  $k$ . Define the confidence balls

$$\mathcal{C}_{h,\iota}^k := \left\{ \theta \in \mathbb{R}^d : \sum_{\tau=1}^k (\hat{y}_{h,\iota}^\tau - \phi_h(s_h^\tau, a_h^\tau)^\top \theta)^2 \leq \beta(\delta) + \|W_{h,\iota}^*\|_2^2 \right\} \quad (5)$$

where  $W_{h,\iota}^*$  is the  $\iota$ -th row of matrix  $W_h^*$ , and  $\beta(\delta) = 1 + 2\rho_{\max} + 32 \log\left(\frac{\sqrt{8} + \sqrt{1 + \rho_{\max}}}{\delta}\right) = \tilde{\mathcal{O}}(\rho_{\max} + \log(1/\delta))$  (where  $\rho_{\max}$  comes from equation 4). Then,

$$\mathbb{P}(\exists k : W_{h,\iota}^* \notin \mathcal{C}_{h,\iota}^k) \leq \delta.$$

*Proof.* In this proof, we fix  $h, \iota$  and call

$$\mathbf{x}_k = \phi_h(s_h^k, a_h^k) \quad \hat{y}_k = \hat{y}_{h,\iota}^k.$$

Under our transition (2)

$$y_k = \mathbf{x}_k^\top \theta^* + \eta_k \quad \theta^* := W_{h,\iota}^*.$$

At this point applying Theorem 23.4 from (Lattimore and Szepesvári, 2020) gives, with probability  $1 - \delta$ , at the same time for every  $k = 1, \dots, K$

$$\theta^* \in \mathcal{C}^k := \left\{ \theta \in \mathbb{R}^d : \|\theta\|_2^2 + \sum_{\tau=1}^k (\hat{y}_\tau - \mathbf{x}_\tau^\top \theta)^2 \leq \beta(\delta) + \|\theta^*\|_2^2 \right\},$$

with our definition of  $\beta(\delta)$ . Erasing  $\|\theta\|_2^2$ , which is non-negative, makes only the set bigger. This ends the proof.  $\square$

**Lemma 3.1.** For every  $h, k, \iota$ , let  $\Lambda_h^k := I + \sum_{\tau=1}^k \phi_h(s_h^\tau, a_h^\tau) \phi_h(s_h^\tau, a_h^\tau)^\top$ , where  $I$  stands for the identity matrix of size  $d$ . Moreover, let  $\widehat{W}_{h,\iota}^k := (\Lambda_h^k)^{-1} \sum_{\tau=1}^k \hat{y}_{h,\iota}^\tau \phi_h(s_h^\tau, a_h^\tau)$ . Then,

$$\mathcal{C}_{h,\iota}^k \subseteq \left\{ \theta \in \mathbb{R}^d : \|\theta - \widehat{W}_{h,\iota}^k\|_{\Lambda_h^k}^2 \leq \beta(\delta) + \|W_{h,\iota}^*\|_2^2 \right\}.$$

*Proof.* Since, also in this case,  $h, \iota$  are fixed, we keep calling

$$\mathbf{x}_k = \phi_h(s_h^k, a_h^k) \quad \hat{y}_k = \hat{y}_{h,\iota}^k,$$

and

$$y_k = \mathbf{x}_k^\top \theta^* + \eta_k \quad \theta^* = W_{h,\iota}^* \quad \hat{\theta}_k = \widehat{W}_{h,\iota}^k \quad \Lambda_k = \Lambda_{h,\iota}^k.$$

In this notation, and using theorem 3.2, we have

$$\mathcal{C}^k = \left\{ \theta \in \mathbb{R}^d : \|\theta\|_2^2 + \sum_{\tau=1}^k (\hat{y}_\tau - \mathbf{x}_\tau^\top \theta)^2 \leq \beta(\delta) + \|\theta^*\|_2^2 \right\},$$

and we want to show it is contained in

$$\tilde{\mathcal{C}}^k = \left\{ \theta \in \mathbb{R}^d : \|\hat{\theta}_k - \theta\|_{\Lambda_k}^2 \leq \beta(\delta) + \|\theta^*\|_2^2 \right\}.$$

This comes from the following identity (see Exercise 23.5 from Lattimore and Szepesvári (2020)):

$$\|\hat{\theta}_k - \theta\|_{\Lambda_k}^2 + \|\hat{\theta}_k\|_2^2 + \sum_{\tau=1}^k (\hat{y}_\tau - \mathbf{x}_\tau^\top \hat{\theta}_k)^2 = \|\theta\|_2^2 + \sum_{\tau=1}^k (\hat{y}_\tau - \mathbf{x}_\tau^\top \theta)^2,$$

as we are adding two terms that are nonnegative to the LHS. □

### C.1 Regret bound

The regret bound builds on the following fundamental properties of bounded variation noises.

**Proposition C.1.** *Let  $\mu(\cdot) \in BV(\chi)$ . Then, for any  $x_0 \in \mathbb{R}^p$ ,*

$$TV(\mu(\cdot), \mu(\cdot + x_0)) \leq \chi \|x_0\|.$$

where  $TV$  denotes the total variation distance.

We prove the proposition for probability distribution defined on the whole space  $\mathbb{R}^p$ . When clipping them in a compact set, the theorem still holds as clipping does not increase the TV distance.

*Proof.* By definition,

$$TV(\mu(\cdot), \mu(\cdot + x_0)) = \sup_{f \in \mathcal{C}(\mathbb{R}^p), \|f\|_{L^\infty} \leq 1} \int_{\mathbb{R}^p} (f(y) - f(y - x_0)) \mu(y) dy.$$

Now, fix a function  $f$ . By Lagrange's theorem,

$$\int_{\mathbb{R}^p} (f(y) - f(y - x_0)) \mu(y) dy \leq \|x_0\| \sup_{x \in \mathbb{R}^p} \left\| \nabla_x \int_{\mathbb{R}^p} f(y - x) \mu(y) dy \right\|.$$

At this point, we have

$$\begin{aligned} \|x_0\| \sup_{x \in \mathbb{R}^p} \left\| \nabla_x \int_{\mathbb{R}^p} f(y - x) \mu(y) dy \right\| &= \|x_0\| \sup_{x \in \mathbb{R}^p} \left\| \nabla_x \int_{\mathbb{R}^p} f(y) \mu(y + x) dy \right\| \\ &= \|x_0\| \sup_{x \in \mathbb{R}^p} \left\| \int_{\mathbb{R}^p} f(y) \nabla_x \mu(y + x) dy \right\| \\ &\leq \|x_0\| \sup_{x \in \mathbb{R}^p} \|\nabla \mu(\cdot)\|_{\mathcal{M}} \leq \|x_0\| \chi. \end{aligned}$$

The second passage comes from the fact that, using the weak gradient, we can bring it under the integral sign under the condition that the integral is finite, which holds since  $f$  is bounded  $\mu$  is a probability density function. □

**Theorem 3.4.** *Assume 2.1 and 3.3. Let  $\chi := \max\{\chi_h : h = 1, \dots, H - 1\}$ . Then, with probability at least  $1 - \delta$ , Algorithm 1 achieves the following regret bound*

$$R_K \leq \tilde{\mathcal{O}}(H^2 \chi p_S \sqrt{d_0 d_K}).$$

*Proof.* We work under the event defined by theorem 3.2, which has probability at least  $1 - \delta$ . Under this even, the true matrix  $W_{h,\ell}^*$  of the transition at step  $h$  (component  $\ell$ ) belongs to  $\mathcal{C}_{h,\ell}^k$  at any episode. Since algorithm 1 chooses its policy by maximizing over the possible transitions in  $\mathcal{C}_{h,\ell}^k$ . We the following bound on the estimated value function at each episode

$$V_1^*(s) \leq V_1^k(s) \quad \forall s_1 \in \mathcal{S}.$$

Thus, the regret can be written in the following form

$$R_K = \sum_{k=1}^K V_1^*(s_1^k) - V_1^{\pi^k}(s_1^k) \quad (11)$$

$$\leq \sum_{k=1}^K V_1^k(s_1^k) - V_1^{\pi^k}(s_1^k) \quad (12)$$

$$\leq \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s_h^k, a_h^k} [(\widehat{P}_h^k(\cdot | s_h^k, a_h^k) - P_h(\cdot | s_h^k, a_h^k))(V_h^{\pi^k})], \quad (13)$$

where step 13 comes from the performance difference lemma and we have defined  $\widehat{P}_h^k(\cdot | s, a)/P_h(\cdot | s, a)$  to be the estimated/true distribution of  $s'$  following  $s, a$ . Indeed, by equation 2, these correspond to

$$\mu_h(\cdot - W_h^* \phi_h(s, a)) \quad \mu_h(\cdot - \widehat{W}_h^k \phi_h(s, a)).$$

Using the fact that  $V_h^{\pi^k}$  is bounded by  $H$ , we can proceed as follows.

$$\begin{aligned} \text{R.H.S.} &\leq H \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s_h^k, a_h^k} [\text{TV}(\mu_h(\cdot - W_h^* \phi_h(s_h^k, a_h^k)), \mu_h(\cdot - \widehat{W}_h^k \phi_h(s_h^k, a_h^k)))] \\ &\leq H \sum_{k=1}^K \sum_{h=1}^H \sum_{\ell=1}^{p_S} 2 \wedge \chi \mathbb{E}_{s_h^k, a_h^k} [\|W_{h,\ell}^* \phi_h(s_h^k, a_h^k) - \widehat{W}_{h,\ell}^k \phi_h(s_h^k, a_h^k)\|], \end{aligned}$$

where in the second passage we have used proposition C.1, dividing the norm into the  $p_S$  component of the state (indeed  $\|x\| \leq \|x\|_1 \leq |x_1| + |x_2| + \dots$ ). The presence of 2 takes into account that the total variation distance can never exceed 2. We then proceed as follows:

$$\text{R.H.S.} \leq H \sum_{k=1}^K \sum_{h=1}^H \sum_{\ell=1}^{p_S} 2 \wedge \chi \mathbb{E}_{s_h^k, a_h^k} [\|W_{h,\ell}^* - \widehat{W}_{h,\ell}^k\|_{\Lambda_{h,\ell}^k} \|\phi_h(s_h^k, a_h^k)\|_{(\Lambda_{h,\ell}^k)^{-1}}] \quad (14)$$

$$\leq H^2 \mathbb{E} \left[ \sum_{k=1}^K \sum_{\ell=1}^{p_S} 2 \wedge \chi \|W_{h,\ell}^* - \widehat{W}_{h,\ell}^k\|_{\Lambda_{h,\ell}^k} \|\phi_h(s_h^k, a_h^k)\|_{(\Lambda_{h,\ell}^k)^{-1}} \right] \quad (15)$$

$$= H^2 \mathbb{E} \left[ \sum_{\ell=1}^{p_S} \sum_{k=1}^K 2 \wedge \chi \|W_{h,\ell}^* - \widehat{W}_{h,\ell}^k\|_{\Lambda_{h,\ell}^k} \|\phi_h(s_h^k, a_h^k)\|_{(\Lambda_{h,\ell}^k)^{-1}} \right] \quad (16)$$

$$\leq H^2 \chi \mathbb{E} \left[ \sum_{\ell=1}^{p_S} \sqrt{K} \sqrt{\sum_{k=1}^K 4 \wedge \|W_{h,\ell}^* - \widehat{W}_{h,\ell}^k\|_{\Lambda_{h,\ell}^k}^2 \|\phi_h(s_h^k, a_h^k)\|_{(\Lambda_{h,\ell}^k)^{-1}}^2} \right] \quad (17)$$

$$\leq H^2 \chi \mathbb{E} \left[ \sum_{\ell=1}^{p_S} \sqrt{K} \sqrt{\sum_{k=1}^K 4(\beta_T(\delta) + m_2) \left(1 \wedge \|\phi_h(s_h^k, a_h^k)\|_{(\Lambda_{h,\ell}^k)^{-1}}^2\right)} \right] \quad (18)$$

$$= H^2 \chi \mathbb{E} \left[ \sum_{\ell=1}^{p_S} \sqrt{4K(\beta_K(\delta) + m_2)} \sqrt{\sum_{k=1}^K \left(1 \wedge \|\phi_h(s_h^k, a_h^k)\|_{(\Lambda_{h,\ell}^k)^{-1}}^2\right)} \right] \quad (19)$$

$$\lesssim H^2 \chi p_S \sqrt{(\beta_K(\delta) + m_2) K} \sqrt{d \log(K)}. \quad (20)$$

In the previous derivation, we have used the shortcut  $m_2 = \max_{h,\ell} \|W_{h,\ell}^*\|_2^2$ . Step 17 comes from the Cauchy-Schwartz inequality. Step 18 comes from our lemma 3.1. Step 20 comes from the elliptical potential lemma (Lemma 11 in (Abbasi-Yadkori and Szepesvári, 2011)). Replacing the definition of  $\beta_K(\delta)$  ends the proof.  $\square$

**Theorem 3.7.** Assume 3.6 and 3.3. Let  $\chi := \max\{\chi_h : h = 1, \dots, H-1\}$ . Let us modify the construction of the confidence sets in Equation 5 by using as confidence radius  $\beta'(\delta) = \beta(\delta) + K\xi_\infty^2 \log(1/\delta)$ . Then, with probability at least  $1 - \delta$ , algorithm 1 achieves the following regret bound

$$R_K \leq \tilde{\mathcal{O}}\left(H^2 \chi p_S \sqrt{d_0 d K} + H^2 p_S \sqrt{d K} \xi_\infty\right).$$

*Proof.* The proof of this result is very similar to the one of the previous theorem 3.4. First, we have to show that, letting  $\beta'(\delta) = \beta(\delta) + K\xi_\infty^2 \log(1/\delta)$ , the optimism is maintained, even with misspecification. This time, we cannot use directly the result by (Lattimore and Szepesvári, 2020) and we need to derive the result from theorem 3.1, assuming that the samples take the form

$$y_t = \mathbf{x}_t^\top \theta + \xi(\mathbf{x}_t) + \eta_t,$$

where  $\eta_t$  is an i.i.d. noise. Our goal is to bound

$$\mathcal{L}_T := \sum_{t=1}^T (x_t^\top \theta - \hat{y}_t)^2, \quad (21)$$

so that, for  $\theta = \theta^*$ , we can use this bound to build a confidence ball. Indeed, theorem 3.1, (which makes no assumptions of  $y_t$  being well-specified), gives, for each  $\theta \in \mathbb{R}^d$

$$\rho_T(\theta) := \sum_{t=1}^T (y_t - \hat{y}_t)^2 - (y_t - \mathbf{x}_t^\top \theta)^2 \leq \rho_{\max}.$$

By definition, calling  $\zeta_t := \eta_t + \xi(x_t)$ ,

$$\begin{aligned} \sum_{t=1}^T (x_t^\top \theta - \hat{y}_t)^2 &= \rho_T(\theta) + \sum_{t=1}^T (x_t^\top \theta - \hat{y}_t)^2 - (y_t - \hat{y}_t)^2 + (y_t - x_t^\top \theta)^2 \\ &= \rho_T(\theta) + \sum_{t=1}^T x_t^\top \theta^2 + \hat{y}_t^2 - 2x_t^\top \theta \hat{y}_t - y_t^2 - \hat{y}_t^2 + 2\hat{y}_t y_t + x_t^\top \theta^2 - 2y_t x_t^\top \theta + y_t^2 \\ &= \rho_T(\theta) + \sum_{t=1}^T 2x_t^\top \theta^2 - 2x_t^\top \theta \hat{y}_t + 2\hat{y}_t y_t - 2y_t x_t^\top \theta \\ &= \rho_T(\theta) + \sum_{t=1}^T 2x_t^\top \theta^2 - 2x_t^\top \theta \hat{y}_t + 2\hat{y}_t (x_t^\top \theta + \zeta_t) - 2(x_t^\top \theta + \zeta_t) x_t^\top \theta \\ &= \rho_T(\theta) + \sum_{t=1}^T 2\hat{y}_t \zeta_t - 2\zeta_t x_t^\top \theta \\ &= \rho_T(\theta) + \sum_{t=1}^T 2(\hat{y}_t - x_t^\top \theta) \zeta_t. \end{aligned}$$

We can rearrange the last result as follows:

$$\sum_{t=1}^T 2(\hat{y}_t - x_t^\top \theta) \zeta_t = \sum_{t=1}^T 2(\hat{y}_t - x_t^\top \theta) \xi(x_t) + \sum_{t=1}^T 2(\hat{y}_t - x_t^\top \theta) \eta_t.$$

These two terms can be bounded by means of  $\Psi_T$  (equation (21)), as follows: i) the random variable  $\sum_{t=1}^T 2(\hat{y}_t - x_t^\top \theta) \eta_t$  is  $\sigma$ -subgaussian with  $\sigma = \sqrt{\Psi_t}$ . Therefore, w.p.  $1 - \delta$ ,

$$\sum_{t=1}^T 2(\widehat{y}_t - x_t^\top \theta) \eta_t \leq 2\sqrt{\Psi_T \log(T/\delta)}.$$

ii) By Cauchy-Schwartz inequality:

$$\sum_{t=1}^T (\widehat{y}_t - x_t^\top \theta) \xi(x_t) \leq \sqrt{\Psi_T \sum_{t=1}^T \xi(x_t)^2} \leq \sqrt{\Psi_T T \xi_\infty^2}.$$

Putting everything together, this shows the following bound on  $\Psi_T$  (eq. 21):

$$\Psi_T \leq \rho_{\max} + \sqrt{\Psi_T (2 \log(T/\delta) + T \xi_\infty^2)},$$

which is satisfied by  $\Psi_T = \beta'(\delta) = \beta(\delta) + T \xi_\infty^2 \log(1/\delta) = \widetilde{\mathcal{O}}(\rho_{\max} + T \xi_\infty^2 \log(1/\delta))$ . This proves that

$$\sum_{t=1}^T (x_t^\top \theta^* - \widehat{y}_t)^2 \leq \beta'(\delta).$$

Replacing back  $x_t = \phi_h(s_h^k, a_h^k)$ ,  $y_t = s_{h+1,\iota}^k$ ,  $\theta^* = W_{h,\iota}^*$ , and following the same passages of the proofs of theorem 3.2 and lemma 3.1, this leads to

$$\mathcal{C}_{h,\iota}^k \subseteq \left\{ \theta \in \mathbb{R}^d : \|\theta - \widehat{W}_{h,\iota}^k\|_{\Lambda_h^k}^2 \leq \beta'(\delta) + \|W_{h,\iota}^*\|_2^2 \right\},$$

with probability  $1 - \delta$ . This proves the optimism, and allows us to bound the regret as

$$R_K = \sum_{k=1}^K V_1^*(s_1^k) - V_1^{\pi^k}(s_1^k).$$

The sequent passages follow as in the proof of theorem 3.4, until getting

$$H \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s_h^k, a_h^k} [\text{TV}((\widehat{P}_h^k(\cdot | s_h^k, a_h^k) - P_h(\cdot | s_h^k, a_h^k)))] .$$

This term is bounded by

$$H \sum_{k=1}^K \sum_{h=1}^H \sum_{\iota=1}^{p_S} \left( 2 \wedge \chi \mathbb{E}_{s_h^k, a_h^k} [ \|W_{h,\iota}^* \phi_h(s_h^k, a_h^k) - \widehat{W}_{h,\iota}^k \phi_h(s_h^k, a_h^k)\| ] + 2\xi_\infty \right),$$

where the first part is bounded as in the previous theorem, giving  $H^2 \chi p_S \sqrt{(\beta'_K(\delta) + m_2)K} \sqrt{d \log(K)}$ , and the second one gives

$$H \sum_{k=1}^K \sum_{h=1}^H \sum_{\iota=1}^{p_S} 2\xi_\infty \leq 2H^2 K p_S \xi_\infty.$$

Since  $\beta'_K(\delta) = \widetilde{\mathcal{O}}(\rho_{\max} + K \xi_\infty^2)$ , this ends the proof.  $\square$

## D Experimental Settings and Additional Results

In this appendix, we discuss the experimental setting for the simulations provided in the main paper, and we provide further results.

### D.1 Experimental setting

**Algorithm Setting** For each test, the policy adopted is parametrized using a neural network with 2 hidden layers and 256 nodes per layer for both the actor and the critic using a ReLu as a non-linear layer. The model-based SINDY version adopts two different reply buffers, one filled with data collected from the real environment and one filled with data from the surrogate one. Differently, for our implementation, a single reply buffer is employed since data generated using the confidence ball is consumed directly from the network for a gradient step. We use PySINDY’s E-SINDY (de Silva et al., 2020; Kaptanoglu et al., 2022) implementation with an ensemble of 20 models with STLRidge and the coefficients of all the models are *ensembled*, taking the median of the coefficients. All the trained models use a threshold of  $1 \times 10^{-3}$  and RIDGE regularization of  $1 \times 10^{-3}$  for all the experiments.

Experiments have been conducted with a discount factor  $\gamma = 0.99$  and a learning rate of  $\eta_a = 3 \cdot 10^{-4}$  for the actor and  $\eta_c = 1 \cdot 10^{-3}$  for the critic except for the Acrobot environment where the learning rate of the critic is lowered to  $\eta_c = 3 \times 10^{-4}$ . Since the model error accumulates at each step the horizon of trajectories in the surrogate environment needs to be tuned in relation to the quality of the model. Following (Zolman et al., 2024, Appendix D), rollouts in the surrogate environment are truncated, and the surrogate environment is reset every time the predicted state exceeds certain thresholds. We assume that physical limitations of the environment are well known a priori, i.e., limits that could damage the physical assets, and we used those as boundaries, as shown in Table 1. For what concern Dreamer we used the v3 version (Hafner et al., 2023) adopting the same configuration the Authors adopted in their work, with a learning rate  $\eta = 4 \cdot 10^{-5}$  and  $\gamma = 0.997$ . For a fair comparison we limited the amount of interaction with the real environment to 256.

Environment	Threshold					
Swing-up	$ x  < 5$	$ \dot{x}  < 10$	$ \cos \theta  < 1.1$	$ \sin \theta  < 1.1$	$ \dot{\theta}  < 10$	
Mountain Car		$x < 0.45$	$x > -1.2$	$ \dot{x}  < 0.07$		
Acrobot	$ \cos \theta_1  < 1.1$	$ \sin \theta_1  < 1.1$	$ \cos \theta_2  < 1.1$	$ \sin \theta_2  < 1.1$	$ \dot{\theta}_1  < 4\pi$	$ \dot{\theta}_2  < 9\pi$

Table 1: Model Thresholds

**Computational Resources.** All the experiments are run on a server equipped as follows:

- CPU: Intel Xeon Gold 6238R (112 cores, 2.20 GHz);
- RAM: 256GB GB;

In particular,  $N = 300$  trajectories for each environments with  $H = 1000$  scored  $\approx 58.82$  iterations per second for SAC and  $\approx 12.71$  iterations per second with model-based SINDY, while  $\approx 8.92$  iterations per second for OSCAR. All the performance are run over a single CPU core.

**Additional Results** In Figure 3 we present an analysis of the performance of OSCAR w.r.t. the scale factor  $\lambda$  in the Swing-Up and Acrobot environment. In Figure 3a, we can notice how small values of  $\lambda$  manage to converge to the optimal behavior, while forcing a broader exploration ( $\lambda = 1$ ) leads us to a suboptimal solution. Since the environment provides dense reward that incentivizes keeping the pole upright and the cart centered, policies that focus their exploration around the optimal region (i.e., the upright position) receive more consistent learning signals. An excessive exploration can thus lead to policies that swing the pole but fail to stabilize it. On the other hand, in environments like Acrobot, with sparse reward signals, wider exploration helps find the optimal solution. We can, in fact, notice, in Figure 3b, how with higher values of  $\lambda$  we manage to converge faster to the optimal solution.

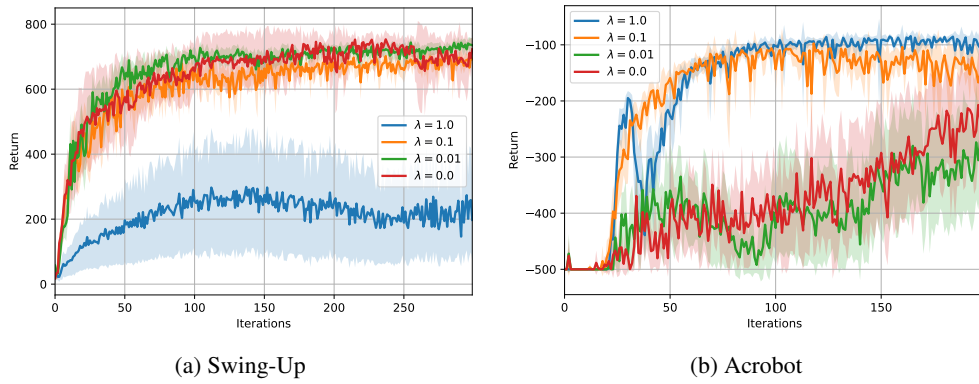


Figure 3: Ball Dimension Analysis