
Reward-Free Exploration beyond Finite-Horizon

Jean Tarbouriech^{1 2} Matteo Pirodda¹ Michal Valko² Alessandro Lazaric¹

Abstract

We consider the reward-free exploration framework recently introduced by Jin et al. (2020), where an RL agent interacts with an unknown environment without any explicit reward function to maximize. The objective is to collect enough information during the exploration phase, so that a near-optimal policy can be immediately computed once a specific reward function is provided. In this paper, we move from the finite-horizon setting studied by Jin et al. (2020) to the more general setting of goal-conditioned RL, often referred to as stochastic shortest path (SSP). We first discuss the challenges specific to SSPs and then we study two scenarios: **1)** reward-free goal-free exploration in communicating MDPs, and **2)** reward-free goal-free incremental exploration in non-communicating MDPs where the agent is provided with an action to reset to an initial state. In both cases, we devise novel exploration algorithms and derive sample-complexity bounds.

1. Introduction

In problems where the reward function is sparse or even absent, a reinforcement learning (RL) agent needs to explore the environment driven by objectives other than reward maximization. Recent unsupervised exploration deep RL algorithms successfully tackled complex problems such as Montezuma’s Revenge (e.g., Ecoffet et al., 2020) or real-world robotic manipulation tasks (e.g., Pong et al., 2020) solely driven by the objective of *discovering* and *controlling* the environment. Nonetheless, the problem still lacks of a rigorous formalization and algorithms do not have solid theoretical guarantees. A first step in that direction is the reward-free exploration framework introduced by Jin et al. (2020) in finite-horizon Markov decision processes (MDPs). Jin et al. (2020) define an exploration phase where the agent interacts with an unknown environment and collects information about its dynamics. Then in a planning phase, the agent is provided with a reward function and it must return

a near-optimal policy without any further learning. The performance of the agent is evaluated by the number of samples collected during the exploration phase.

While the finite-horizon setting is very popular in theoretical RL, it is rarely representative of the type of problems considered in popular benchmarks and real applications in RL. In this paper, we rather focus on the strictly more general and more practical stochastic shortest path (SSP) setting (Bertsekas, 2012) (often referred to as goal-conditioned RL), where the objective is to compute a policy that minimizes the cost accumulated before reaching a specific goal state. We first reformulate the reward-free exploration setting by defining the objective of learning an accurate enough model of the environment so that a near-optimal policy can be computed for *any* SSP problem (i.e., for any initial state, any goal state, and any cost function). We illustrate how this problem may be considerably more difficult than in the finite-horizon setting. We then study two different scenarios (i.e., goal-free cost-free exploration in communicating MDPs and goal-free cost-free incremental exploration in non-communicating MDPs with restart), summarize the sample complexity results that we obtain and contrast them with the guarantees in the finite-horizon case.

2. Preliminaries

A Markov decision process (MDP) is defined as $M := \langle \mathcal{S}, \mathcal{A}, p, c \rangle$, where \mathcal{S} is the state space with $S := |\mathcal{S}|$ states and \mathcal{A} is the action space with $A := |\mathcal{A}|$ actions. Taking action a in state s incurs a cost¹ of $c(s, a) \in [0, 1]$ and the next state $s' \in \mathcal{S}$ is selected with probability $p(s'|s, a)$. We denote by $\Gamma := \max_{s,a} \|p(\cdot|s, a)\|_0$ the largest support of the transition model. In the SSP case, for a designated goal state \bar{s} , the objective is to compute a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ minimizing the cumulative cost before reaching \bar{s} . Formally, we define the (possibly unbounded) value function

$$V_\pi(\underline{s} \rightarrow \bar{s}) := \mathbb{E} \left[\sum_{t=1}^{\tau_\pi(\underline{s} \rightarrow \bar{s})} c(s_t, \pi(s_t)) \mid s_1 = \underline{s} \right],$$

where $\tau_\pi(\underline{s} \rightarrow \bar{s}) := \inf\{t \geq 0 : s_{t+1} = \bar{s} \mid s_1 = \underline{s}, \pi\}$ is the (random) number of steps needed to reach \bar{s} from \underline{s} when executing policy π . An optimal policy (if it exists) is

¹Facebook AI Research, Paris ²Inria Lille - Nord Europe.

¹One can translate between costs and rewards by simply taking negation.

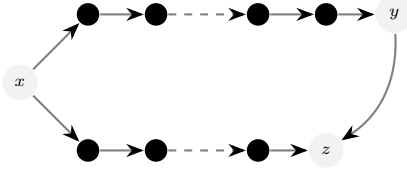


Figure 1. The agent starts at state x and reaches z in H steps with probability $1/2$, and y in $H + 1$ steps with probability $1/2$. From state y the agent deterministically transitions to state z in 1 step.

denoted by $\pi^* \in \arg \min_{\pi} V_{\pi}(\underline{s} \rightarrow \bar{s})$. For more details on the SSP problem we refer to e.g., Bertsekas (2012, Sect. 3).

Jin et al. (2020) introduced the reward-free framework in the finite-horizon case, which is a special case of the SSP problem where each episode terminates after exactly H steps. The agent receives as input an accuracy level $\varepsilon > 0$, a confidence level $\delta \in (0, 1)$, the state and action spaces, and the horizon H , while no knowledge is provided about the transition model p . The learning process is decomposed into two phases. ① *Exploration phase*: The agent first collects trajectories from the MDP without a pre-specified reward function and returns an estimate of the transition model \hat{p} . ② *Planning phase*: The agent receives an arbitrary reward function and is tasked with computing an ε -optimal policy with probability at least $1 - \delta$, without any additional interaction with the environment. The objective is to minimize the duration of the exploration phase needed to simultaneously enforce any requested planning guarantee.

Jin et al. (2020) study the reward-free exploration problem for any arbitrary MDP, where there may exist states that are difficult or impossible to reach. The core mechanism in their analysis is to partition the states depending on their ease of being reached within H steps. Specifically, they distinguish between *significant* states, that can be sufficiently visited and whose transition probability can thus be accurately estimated, and *insignificant* states that are too difficult to reach within H steps, but therefore have negligible contribution to any reward optimization.

Interestingly, in the goal-conditioned setting this distinction may no longer be meaningful. By way of illustration, consider any fixed horizon H and the toy environment in Fig. 1. Suppose that the objective is to quickly reach state z (i.e., the goal state is z , the starting state is x and all costs are equal to 1). Even though state y is *insignificant* within H steps (in the finite-horizon sense of Jin et al., 2020, for any positive “significance level”), it is actually crucial in solving the objective, as z can be reached deterministically in 1 step from y . Extrapolating this scenario, in the goal-conditioned setting, we may have an effective horizon of $H = +\infty$ for some goals, which implies that the transition model p must be accurately estimated across the state-action space to ensure that a near-optimal policy can be computed.

3. Goal-Free Cost-Free Exploration in Communicating MDPs

In order to guarantee that the environment can be estimated uniformly well, we introduce the following assumption.

Assumption 1 (In Sect. 3). *The MDP M is communicating, with finite and unknown diameter*

$$D = \max_{\underline{s}, \bar{s}} D_{\underline{s}, \bar{s}} = \max_{\underline{s}, \bar{s}} \min_{\pi} \mathbb{E}[\tau_{\pi}(\underline{s} \rightarrow \bar{s})] < +\infty.$$

We stress that the challenges that emerge in such setting are orthogonal to the ones in (Jin et al., 2020): a *constraint on the environment is added* (all states must now be reachable), allowing the *removal of the constraint on performance* (which is not limited to H steps anymore) and thus enabling to tackle the more general class of goal-oriented problems.

Without loss of generality, we consider throughout that the maximum c_{\max} of the cost functions that we intend to consider in the planning phase is equal to 1. On the other hand, the minimum value c_{\min} has a more subtle impact on the type of performance guarantees we can obtain. In particular, for any cost function c and any pair of initial and goal states \underline{s} and \bar{s} , we introduce a slack parameter $\theta \in [1, +\infty]$ and we say that a policy $\hat{\pi}$ is (ε, θ) -optimal if²

$$V^{\hat{\pi}}(\underline{s} \rightarrow \bar{s}) \leq \min_{\pi: \mathbb{E}[\tau_{\pi}(\underline{s} \rightarrow \bar{s})] \leq \theta D_{\underline{s}, \bar{s}}} V^{\pi}(\underline{s} \rightarrow \bar{s}) + \varepsilon.$$

In the following theorem, we show that depending on the minimum cost c_{\min} in the cost functions of interest and the slack θ , we can solve the goal-free cost-free exploration problem with a bounded sample complexity.

Theorem 1. *Consider any unknown environment satisfying Asm. 1 and the goal-free cost-free exploration problem characterized by an accuracy level $0 < \varepsilon \leq 1$, a confidence level $\delta \in (0, 1)$, a minimum cost $c_{\min} \in [0, 1]$ and a slack parameter $\theta \in [1, +\infty]$. There exists an algorithm \mathfrak{A} whose exploration phase (i.e., number of time steps) is bounded with probability at least $1 - \delta$ by*

$$\tilde{O}\left(\frac{D^4 \Gamma S A}{\omega \varepsilon^2} + \frac{D^3 S^2 A}{\omega \varepsilon} + \frac{D^3 \Gamma S A}{\omega^2}\right),$$

where: $\omega := \max\left\{c_{\min}, \frac{\varepsilon}{\theta D}\right\}.$

Note that we can have either $c_{\min} = 0$ or $\theta = +\infty$, but not both simultaneously, to guarantee that $\omega > 0$. Following this exploration phase, the algorithm \mathfrak{A} can compute in the planning phase, for any pair of starting and goal state $(\underline{s}, \bar{s}) \in \mathcal{S}^2$, and for any cost function c in $[c_{\min}, 1]$, a policy $\hat{\pi}$ (depending on $c, \underline{s}, \bar{s}$) that is (ε, θ) -optimal.

Algorithmic principle (see App. A). We first use a sample complexity analysis to solve SSP problems with a generative model (Tarbouriech et al., 2020a) and define the number of samples that are needed in each state-action pair to

²This reduces to standard ε -optimality for $\theta \rightarrow \infty$.

compute an estimated model that is accurate enough that a near-optimal policy can be computed for any cost function. Then we leverage the online learning algorithm GOSPRL (Tarbouriech et al., 2020b) that explicitly collects the desired amount of samples in any communicating environment. Interestingly, such algorithm is simply defining a sequence of SSP problems, where the goal state is any state for which the required number of samples is not achieved yet.

4. Goal-Free Cost-Free Incremental Exploration

In this section, we seek to provide cost-free guarantees in MDPs with possibly very large state space and diameter (e.g., for non-communicating MDPs where $D = \infty$). In order to make such setting feasible, we need to restrict the type of SSP problems we would like to solve during the planning phase. We propose an alternative approach that builds on the setting of incremental autonomous exploration introduced by Lim & Auer (2012).

Assumption 2 (In Sect. 4). *The MDP M has a finite, possibly large state space \mathcal{S} for which an upper bound S on its cardinality is known, i.e., $|\mathcal{S}| \leq S$.³ It contains a designated initial state $s_0 \in \mathcal{S}$. Since the learner may get stuck in a state without being able to return to s_0 , we assume that the action space contains a RESET action s.t. $p(s_0|s, \text{RESET}) = 1$ for any $s \in \mathcal{S}$.*

We make explicit the states where a policy π takes action RESET in the following definition.

Definition 1. *For $\mathcal{S}' \subseteq \mathcal{S}$ a policy π is restricted on \mathcal{S}' if $\pi(s) = \text{RESET}$ for any $s \notin \mathcal{S}'$. We denote by $\Pi(\mathcal{S}')$ the set of policies restricted on \mathcal{S}' .*

We denote by $\Gamma_{\mathcal{S}'} := \max_{s \in \mathcal{S}', a} \|\{p(s'|s, a)\}_{s' \in \mathcal{S}'}\|_0$ the largest support of the model p restricted to states in $\mathcal{S}' \subseteq \mathcal{S}$.

In (Lim & Auer, 2012), given an input parameter $L \geq 1$ and accuracy $\varepsilon > 0$, the objective of the agent is to identify the set of *incrementally L -controllable states* $\mathcal{S}_L^{\rightarrow}$ (Def. 2), as well as a set of goal-conditioned policies to reach each state in $\mathcal{S}_L^{\rightarrow}$ from s_0 in at most $L + \varepsilon$ steps on average.⁴

Definition 2 (Incrementally controllable states $\mathcal{S}_L^{\rightarrow}$). *Let \prec be some partial order on \mathcal{S} . The set $\mathcal{S}_L^{\rightarrow}$ of states controllable in L steps w.r.t. \prec is defined inductively as follows. The initial state s_0 belongs to $\mathcal{S}_L^{\rightarrow}$ by definition and if there exists a policy π restricted on $\{s' \in \mathcal{S}_L^{\rightarrow} : s' \prec s\}$ with*

³Lim & Auer (2012) consider a countable, possibly infinite state space; however this leads to a technical issue in the analysis of UCBEXPLORE (acknowledged by the authors via personal communication), which requires considering finite state spaces.

⁴Lim & Auer (2012) showed that discovering all states in $\mathcal{S}_L := \{s \in \mathcal{S} : \min_{\pi \in \Pi} \mathbb{E}[\tau_{\pi}(s_0 \rightarrow s)] \leq L\}$ may require a number of exploration steps that is *exponential* in L or $|\mathcal{S}_L|$, hence the definition of incrementally controllable states.

$\mathbb{E}[\tau_{\pi}(s_0 \rightarrow s)] \leq L$, then $s \in \mathcal{S}_L^{\rightarrow}$. The set $\mathcal{S}_L^{\rightarrow}$ of incrementally L -controllable states is defined as $\mathcal{S}_L^{\rightarrow} := \bigcup_{\prec} \mathcal{S}_L^{\rightarrow}$, where the union is over all possible partial orders.

Finally, we introduce $S_L := |\mathcal{S}_L^{\rightarrow}|$ and $\Gamma_L := \Gamma_{\mathcal{S}_L^{\rightarrow}}$.

We extend the formalism of (Lim & Auer, 2012) and define a more challenging cost-free objective. At the end of the exploration phase, an algorithm should be able to compute a near-optimal policy restricted on $\mathcal{S}_L^{\rightarrow}$ for any SSP problem with initial state s_0 , any goal state $s \in \mathcal{S}_L^{\rightarrow}$, and any cost function. As the state space \mathcal{S} may be very large, the set $\mathcal{S}_L^{\rightarrow}$ effectively captures our area of interest, with L being the radius of interest provided as input. Note that $\mathcal{S}_L^{\rightarrow}$ is unknown in advance and is hard to estimate online.

Similar to Thm. 1 we provide an (ε, θ) -optimality guarantee for the planning phase with the additional condition that we consider policies restricted to the initially unknown set $\mathcal{S}_L^{\rightarrow}$.

Theorem 2. *Consider any unknown environment satisfying Asm. 2 and the goal-free cost-free incremental exploration problem characterized by an accuracy level $0 < \varepsilon \leq 1$, a confidence level $\delta \in (0, 1)$, a minimum cost $c_{\min} \in [0, 1]$ and a slack parameter $\theta \in [1, +\infty]$. There exists an algorithm \mathfrak{A} whose exploration phase (i.e., number of time steps) is bounded with probability at least $1 - \delta$ by*

$$\tilde{O}\left(\frac{L^5 \Gamma_{L+\varepsilon} S_{L+\varepsilon} A}{\omega^2 \varepsilon^2} + \frac{L^3 S_{L+\varepsilon}^2 A}{\omega \varepsilon}\right),$$

$$\text{where: } \omega := \max\left\{c_{\min}, \frac{\varepsilon}{\theta L}\right\}.$$

Note that we can have either $c_{\min} = 0$ or $\theta = +\infty$, but not both simultaneously, to guarantee that $\omega > 0$. Following this exploration phase, the algorithm \mathfrak{A} has confidently identified a set $\mathcal{S}_L^{\rightarrow} \subseteq \mathcal{K} \subseteq \mathcal{S}_{L+\varepsilon}^{\rightarrow}$, and has collected enough information such that for any goal state $\bar{s} \in \mathcal{S}_L^{\rightarrow}$ and any cost function c in $[c_{\min}, 1]$, it can compute in the planning phase a policy $\hat{\pi}$ (depending on c, \bar{s}) that verifies

$$V^{\hat{\pi}}(s_0 \rightarrow \bar{s}) \leq \min_{\pi \in \Pi(\mathcal{S}_L^{\rightarrow}) : \mathbb{E}[\tau_{\pi}(s_0 \rightarrow \bar{s})] \leq \theta L} V^{\pi}(s_0 \rightarrow \bar{s}) + \varepsilon.$$

Algorithmic principle. Despite the difference in the setting, we leverage similar algorithmic principles as in Sect. 3. In this case, we define the sample requirements *limited* to the states that have currently been discovered and for which a shortest-path policy is available. Such policies are then used to collect new samples and a one-step random exploration is used to expand the set of controllable states until all states in $\mathcal{S}_L^{\rightarrow}$ have been identified. The resulting algorithm, called DISCO, is presented in (Tarbouriech et al., 2020c).⁵

⁵While in (Tarbouriech et al., 2020c) the presentation of the algorithm deals with the unit-cost case, the extension to handle general costs is straightforward as explained in the paper's Sect. 2.3.

Reference	RF-FH — (Jin et al., 2020)	RF-COMM — Sect. 3 of this paper (Thm. 1)		RF-INC — Sect. 4 of this paper (Thm. 2)	
Setting	Finite-horizon RL	Goal-conditioned RL (i.e., SSP)		Goal-conditioned RL (i.e., SSP)	
Feedback	Any rewards $r \in [0, 1]$	Any goal state, any costs $c \in [c_{\min}, 1]$ with $c_{\min} \geq 0$		Any goal state in S_L^{\rightarrow} , any costs $c \in [c_{\min}, 1]$ with $c_{\min} \geq 0$	
MDP	👉 Non-communicating and resetting after H steps	👈 Communicating with diameter D		👉 Non-communicating and reset action	
Optimality	👈 Restricted to H steps	👉 Arbitrary* length to goal		👉 Arbitrary* length to goal + 👈 Incremental Optimality	
State dep.	👈 Total state space S	👈 Total state space S		👉 State space of interest $S_L^{\rightarrow} \ll S$	
Sample comp.	$\tilde{O}\left(\frac{S^2 A \text{poly}(H)}{\varepsilon^2} + \frac{S^4 A \text{poly}(H)}{\varepsilon}\right)$	$\tilde{O}\left(\frac{S^2 A \text{poly}(D)}{\varepsilon^2}\right)$ for not too small $c_{\min} > 0$	$\tilde{O}\left(\frac{S^2 A \text{poly}(D)}{\varepsilon^3}\right)$ for very small $c_{\min} \simeq 0$	$\tilde{O}\left(\frac{S_L^2 A \text{poly}(L)}{\varepsilon^2}\right)$ for not too small $c_{\min} > 0$	$\tilde{O}\left(\frac{S_L^2 A \text{poly}(L)}{\varepsilon^4}\right)$ for very small $c_{\min} \simeq 0$

Table 1. High-level comparison between (Jin et al., 2020) and this paper. Asterisk* introduces the subtlety that, only in the case of $c_{\min} \simeq 0$ (the second sub-column), the length to goal targeted by the candidate policy is restricted if it is “too long”.

5. Discussion

We stress that **RF-FH** (Jin et al., 2020), **RF-COMM** (Sect. 3) and **RF-INC** (Sect. 4) tackle orthogonal settings, each posing different challenges. That notwithstanding, we believe that it is insightful to compare the three settings in terms of algorithmic approach and resulting bound (see Table 1).

Similarities in the three algorithmic designs. All three approaches construct accurate estimates of the transitions. **RF-FH** (Jin et al., 2020) restrict their attention to “significant” states within H steps. As previously explained, such a reasoning cannot be directly extended to general SSP problems, as there is no more notion of fixed horizon, with some states possibly becoming non-negligible for value optimization at some random point before the goal state is reached. This is why **RF-COMM** enforces to visit uniformly enough the state-action space, which explains the need for the communicating assumption (Asm. 1). By focusing on incremental exploration, **RF-INC** can effectively restrict its attention to the (unknown) state space of interest S_L^{\rightarrow} , which removes the need for the communicating assumption. Finally, note that to collect the sought-after samples, Jin et al. (2020) deploy a finite-horizon algorithm for regret minimization, whereas our algorithms leverage SSP policies.

Comparison between RF-FH and RF-COMM. In the main order term w.r.t. ε , the dependencies in S^2 and A are equivalent, matching the lower bound derived in the finite-horizon case (Jin et al., 2020, Thm. 4.1). Moreover, the role of the horizon H in **RF-FH** is captured by the ratio D/c_{\min} in **RF-COMM** (when $c_{\min} > 0$). Note that this ratio is not a strict horizon (as the performance may last longer, as opposed to finite-horizon which always truncates it at H steps), and it is environment-dependent and thus crucially *unknown*, which introduces an additional layer of complexity to the

problem. **RF-COMM** (Thm. 1) is the first result tackling the reward-free framework beyond finite-horizon, for goal-conditioned RL in communicating MDPs. The resulting exploration bound scales polynomially with D , which is somewhat unavoidable. Finally, the bound of **RF-COMM** inherits a $\tilde{O}(\varepsilon^{-2})$ dependency (as in **RF-FH** of Jin et al., 2020) whenever c_{\min} is not too small (and can therefore be considered as a constant). Otherwise, **RF-COMM** can cope with very small (or even zero-valued) c_{\min} , yet the bound worsens to $\tilde{O}(\varepsilon^{-3})$, and the performance becomes *restricted* to policies with not too large expected goal-reaching time (via the slack parameter θ). This interesting behavior does not appear in the finite-horizon case (where the range of rewards has no influence on the rate in ε), and it captures the key role of the minimum cost played in the behavior of the optimal goal-reaching policy.

Specificity of RF-INC. The incremental focus of **RF-INC** enables to tackle goal-conditioned tasks while removing the communicating assumption of **RF-COMM**, where the dependency on the diameter D is replaced by the parameter L , which may be designed to be much smaller than D . In fact, while L defines the horizon of interest, resetting after every L steps (as in finite-horizon) would prevent the agent to identify incrementally L -controllable states and lead to poor performance. Another interesting element of comparison is the dependency on the size of the state space. While the **RF-FH** algorithm of (Jin et al., 2020) is robust w.r.t. states that can be reached with very low probability, it still displays a polynomial dependency on the global state space S . On the other hand, in virtue of its incremental focus, **RF-INC** (Thm. 2) depends polynomially on the number of $(L + \varepsilon)$ -controllable states and only *logarithmically* on S . This result is significant since not only $S_{L+\varepsilon}$ can be arbitrarily smaller than S , but also because the set $S_{L+\varepsilon}^{\rightarrow}$ itself is initially unknown to the learner.

REFERENCES

- Bertsekas, D. *Dynamic programming and optimal control*, volume 2. 2012.
- Ecoffet, A., Huizinga, J., Lehman, J., Stanley, K. O., and Clune, J. First return then explore. *arXiv preprint arXiv:2004.12919*, 2020.
- Jin, C., Krishnamurthy, A., Simchowitz, M., and Yu, T. Reward-free exploration for reinforcement learning. In *International Conference on Machine Learning*, 2020.
- Lim, S. H. and Auer, P. Autonomous exploration for navigating in mdps. In *Conference on Learning Theory*, pp. 40–1, 2012.
- Pong, V. H., Dalal, M., Lin, S., Nair, A., Bahl, S., and Levine, S. Skew-fit: State-covering self-supervised reinforcement learning. In *International Conference on Machine Learning*, 2020.
- Tarbouriech, J., Pirotta, M., Valko, M., and Lazaric, A. On the sample complexity of stochastic shortest path with a generative model, 2020a. URL https://jtarbouriech.github.io/docs/ssp_genmodel.pdf.
- Tarbouriech, J., Pirotta, M., Valko, M., and Lazaric, A. A provably efficient sample collection strategy for reinforcement learning. *arXiv preprint arXiv:2007.06437*, 2020b.
- Tarbouriech, J., Pirotta, M., Valko, M., and Lazaric, A. Improved sample complexity for incremental autonomous exploration in mdps. *Advances in Neural Information Processing Systems*, 33, 2020c.

A. Cost-Free Goal-Free Exploration in Communicating MDPs (Sect. 3)

We leverage the GOSPRL algorithm of (Tarbouriech et al., 2020b), an algorithm that mimics the behavior of a generative model in communicating MDPs. Specifically, in any unknown communicating environment with diameter D and for any arbitrary (possibly time-varying) requirement of samples $b_t(s, a)$ (where the sequence is bounded from above by $\bar{b}(s, a)$), GOSPRL requires (with high probability) at most $\tilde{O}(BD + D^{3/2}S^2A)$ times steps to collect the sought-after samples for each state-action pair (s, a) , where $B \leq \sum_{s,a} \bar{b}(s, a)$.

We now show that instantiating GOSPRL for carefully selected sampling requirements $b_t(s, a)$ enables to obtain the guarantee of Thm. 1. To do so, we build on the sample complexity analysis of solving SSP problems with a generative model derived in (Tarbouriech et al., 2020a, Thm. 1). As such, we introduce the following sampling requirement function

$$\phi(X, y) := \alpha \cdot \left(\frac{X^3 \hat{\Gamma}}{y \varepsilon^2} \log \left(\frac{XSA}{y \varepsilon \delta} \right) + \frac{X^2 S}{y \varepsilon} \log \left(\frac{XSA}{y \varepsilon \delta} \right) + \frac{X^2 \hat{\Gamma}}{y^2} \log^2 \left(\frac{XSA}{y \delta} \right) \right), \quad (1)$$

where $\alpha > 0$ is a numerical constant and $\hat{\Gamma} := \max_{s,a} \|\hat{p}(\cdot|s, a)\|_0 \leq \Gamma$ is the largest support of \hat{p} .

This sampling requirement function for carefully selected values of X and y is used to guide the GOSPRL algorithm. Specifically, we set y to be equal to the minimum cost (in either the true or cost-perturbed model), i.e., $y := \omega^{-1}$. As for the value of X , let us perform the following distinction of cases.

① First let us assume that the learning agent has prior knowledge of the diameter D . Then we set $X = D$. From (Tarbouriech et al., 2020a), collecting at least $\phi(D, \omega^{-1})$ samples from each state-action pair enables to guarantee the ε -optimality cost-free planning guarantee of Thm. 1. The total time required to collect such samples is upper bounded by $DSA\phi(D, \omega^{-1})$, which directly yields the sample complexity guarantee stated in Thm. 1.

② Second we show that we can relax the assumption of knowing the diameter D without altering the sample complexity guarantee. To do so, we begin the algorithm by a procedure which computes a quantity \hat{D} such that $D \leq \hat{D} \leq D(1 + \varepsilon)$ with high probability. From (Tarbouriech et al., 2020b, App. H), this can be done in $\tilde{O}(D^3 S^2 A / \varepsilon^2)$ time steps by leveraging GOSPRL. We thus begin the algorithm by running such diameter-estimation subroutine. Crucially, we note that its sample complexity is *subsumed* in the total sample complexity of Thm. 1. Then we simply apply the reasoning in case ① by considering $X = \hat{D}$ in the allocation of Eq. 1 instead of $X = D$. Since \hat{D} is a sufficiently tight upper bound on D (i.e., $\hat{D} = O(D)$), we ultimately obtain the same sample complexity guarantee as in case ①.