
Provable Safe Reinforcement Learning with Binary Feedback

Andrew Bennett
Cornell University

Dipendra Misra
Microsoft Research

Nathan Kallus
Cornell University

Abstract

Safety is a crucial necessity in many applications of reinforcement learning (RL), whether robotic, automotive, or medical. Many existing approaches to safe RL rely on receiving numeric safety feedback, but in many cases this feedback can only take binary values; that is, whether an action in a given state is safe or unsafe. This is particularly true when feedback comes from human experts. We therefore consider the problem of provable safe RL when given access to an offline oracle providing binary feedback on the safety of state, action pairs. We provide a novel meta algorithm, SABRE, which can be applied to any MDP setting given access to a blackbox PAC RL algorithm for that setting. SABRE applies concepts from active learning to reinforcement learning to provably control the number of queries to the safety oracle. SABRE works by iteratively exploring the state space to find regions where the agent is currently uncertain about safety. Our main theoretical results shows that, under appropriate technical assumptions, SABRE never takes unsafe actions during training, and is guaranteed to return a near-optimal safe policy with high probability. We provide a discussion of how our meta-algorithm may be applied to various settings studied in both theoretical and empirical frameworks.

1 INTRODUCTION

Reinforcement learning (RL) is an important paradigm that can be used to solve important dynamic decision-making problems in a diverse set of fields, such as robotics, transportation, healthcare, and user assistance. In recent years there has been a significant increase in interest in this problem, with many proposed solutions. However, in many

such applications there are important safety considerations that are difficult to address with existing techniques.

Let us consider the running example of a cleaning robot, whose task is to learn how to vacuum the floor of a house. The primary goal of the robot, of course, is to learn to vacuum as efficiently as possible, which may be measured by the amount cleaned in a given time. However, we would also like to impose certain safety constraints on the robot's actions; for example, the robot should not roll off of a staircase where it could damage itself, it should not roll over electrical cords, or it should not vacuum up the owner's possessions. In this example, there are several desirable properties we would like a safety-aware learning algorithm to have, including:

1. The agent should avoid taking any unsafe actions, *even during training*
2. Since it is hard to concretely define a safety function from the robot's sensory observations *a priori*, we would like the agent to *learn* a safety function given feedback of observed states
3. Since the notion of safety is human-defined, and we would like the safety feedback to be manually provided by humans (*e.g.* the owner), we would want the agent to ask for *as little feedback as possible*
4. We would like to use *binary* feedback (*i.e.* is an action in a given state safe or unsafe) rather than numeric feedback, as this is more natural for humans to provide
5. Since the agent may need to act in real time without direct intervention, they should only ask for feedback *offline* in between episodes

Moving away from our specific example, the above five properties would be ideal for safe RL in many applications where safety is naturally human-defined. Unfortunately, there are no existing safe RL methods that can provably satisfy all of these properties. In particular, existing safe RL methods all fail to satisfy these properties for at least one of the following reasons: (1) they assume a safety function is fully known (*e.g.* Chow et al., 2018; Simão et al., 2021); (2) they learn safety using numeric feedback (*e.g.*

Wachi and Sui, 2020; Amani et al., 2021); or (3) they require a human-in-the-loop who can intervene and monitors the safety of every action taken in real time (Saunders et al., 2018). In addition, none of these existing methods address the issue of only asking for minimal feedback, which is an important consideration which places the problem at the interface of RL and active learning.

In this paper, we present a novel safe RL framework encompassing the above concerns, involving an offline safety-labeling oracle that can provide binary feedback on the safety of state, action pairs in between episodes. We provide a meta algorithm, SABRE, which can be applied to *any* MDP class within this safety framework, given a blackbox RL algorithm for the MDP class. This algorithm utilizes ideas from disagreement-based active learning, by iteratively exploring the state space to find all regions where there is disagreement on possible safety. Under some appropriate technical assumptions, including that the blackbox RL algorithm can provably optimize any given reward, SABRE will satisfy all of the above five properties, and will return an approximately optimal safe policy with high probability. Importantly, we provide high-probability bounds on the number of samples and calls to the labeling oracle needed, and show that they are both polynomial, and that the latter is lower order than the former. Finally, we provide some discussion of how this meta-algorithm approach may be applied in various settings of both theoretical and practical interest.

Math Notation For any natural number $n \geq \mathbb{N}$, we let $[n]$ denote the set $\{1, 2, \dots, n\}$. For any countable set X , we let $\mathcal{P}(X)$ denote the space of all probability distributions over X . Given sets X and Y , we let $X \rightarrow Y$ denote the set of all functions from X to Y . Lastly, $\text{sign}(y)$ denotes the sign function which takes a value of 1 if $y \geq 0$, or -1 if $y < 0$.

2 PROBLEM SETUP

Markov Decision Process We consider learning in finite-horizon Markov decision processes (MDPs). A *reward-free* MDP is characterized by a tuple $(S; A; T; \gamma; H)$, where S is a given (potentially infinitely large) state space, A is a finite action space with $|A| = A$, $T : S \times A \rightarrow \mathcal{P}(S)$ is a transition operator, $\gamma \in [0, 1)$ is an initial state distribution, and $H \geq \mathbb{N}$ is the horizon. An MDP (with reward) is a tuple $(S; A; T; R; \gamma; H)$, where $R : S \times A \rightarrow [0, 1]$ is the reward function. We assume that the transition operator and reward function are time homogeneous; *i.e.*, T and R do not depend explicitly on the time index.¹ We will let $(M; R)$ refer to the MDP we are learning in, where M is

¹This is w.l.o.g. since we can always include the current time index in the observation definition.

the corresponding reward-free MDP.

The agent interacts with an MDP over a series of rounds. In each round the agent successively takes actions according to some *policy*, which is a mapping from states to actions, in order to generate an *episode*. Specifically, in the n 'th round for each $n \geq \mathbb{N}$, the agent selects some policy $\pi_n \in \Pi(S \rightarrow A)$, which it then uses in order to generate an episode $\tau_n = (s_1^n; a_1^n; r_1^n; \dots; s_H^n; a_H^n; r_H^n; s_{H+1}^n)$, where $s_1^n \sim \gamma$, and for each $h \geq [H]$ we have $a_h^n \sim \pi_n(s_h^n)$, $r_h^n = R(s_h^n; a_h^n)$, and $s_{h+1}^n \sim T(s_h^n; a_h^n)$.

For any given policy π , we let $\mathbb{E}[\cdot]$ and $\mathbb{P}(\cdot)$ denote the expectation and probability operators over trajectories generated using π in the MDP $(M; R)$. Also, for any policy π , we define the *value function* $V : S \rightarrow \mathbb{R}$ according to $V(s) = \mathbb{E}[\sum_{h=1}^H R(s_h; a_h) \mid s_1 = s]$, and we similarly define the value of the policy according to $V(\pi) = \mathbb{E}_s[\mathbb{E}_\pi[V(s)]]$.

Our main metric of success for reinforcement learning is *suboptimality* (SubOpt). For any policy class $\Pi(S \rightarrow A)$ and any policy $\hat{\pi} \in \Pi$, we define $\text{SubOpt}(\hat{\pi}; \Pi) = \sup_{\pi \in \Pi} V(\pi) - V(\hat{\pi})$. Then, one of our primary goals is to learn a policy $\hat{\pi}$ that has low suboptimality with respect to a given set of policies that it must choose from, in a relatively small number of episodes.

Over the past few decades, many reinforcement learning algorithms have been developed that can provably obtain low suboptimality using a small number of episodes, for various classes of MDPs. Our focus is on adapting such RL algorithms, in order to additionally take into account safety considerations. For these reasons, we take a meta-learning approach, and assume access to a blackbox reinforcement learning algorithm Alg . We formalize this as follows:

Assumption 1. *We have access to a blackbox RL algorithm Alg for a set of reward-free MDPs \mathcal{M} that contains M . Specifically, given any reward function $R^0 : S \times A \rightarrow [0, 1]$ and policy class Π as input, along with any $\epsilon \in (0, 1)$, $\text{Alg}(R^0; \Pi; \epsilon)$ returns a policy $\hat{\pi} \in \Pi$ such that $\text{SubOpt}(\hat{\pi}; \Pi) \leq \epsilon$ with probability at least $1 - \delta$. Furthermore, it is guaranteed to do so after at most $n_{\text{Alg}}(\epsilon; \delta)$ episodes, for some $n_{\text{Alg}}(\epsilon; \delta) = \text{Poly}(\epsilon^{-1} \log(1/\delta))$, while only following policies in Π .*

We refer to $n_{\text{Alg}}(\epsilon; \delta)$ as the *sample complexity* of Alg , and note that it implicitly depends on the time horizon H . We also note that the requirement that Alg only follows policies in Π is important, as in practice we will call Alg using classes of policies that are guaranteed to be safe.

Safety Considerations In addition to the general MDP setup presented above, we would also like to take safety considerations into account. Specifically, we will use a binary notion of safety as follows: we assume there exists a function $f^? : S \rightarrow \{0, 1\}$ such that $f^?(s; a) = 1$ means that taking action a in state s is safe, and otherwise ac-

Figure 1: Illustrating the challenge with binary safety feedback. Left: we observe binary safe/unsafe labels (our setting) and we know the safety function is a halfspace; the safety of state-action pairs in the region of disagreement (shaded gold region) is uncertain, as there exist halfspaces consistent with both the observed data and these being either safe or unsafe. Right: we observe noisy numeric safety values (positive being safe) and we know their mean is linear; the region of disagreement (using 95% confidence) is far smaller because we can extrapolate to yet-unseen state-action pairs.

tion a in states is unsafe. Importantly, we assume that the that is, $f^? \in F$. Furthermore, there exists a policy $\pi_{safe} \in F$ that is known to always take safe actions. safety function $f^?$ is unknown, and we require the agent to only take actions a in states s such that $f^?(s; a) = 1$, even during training. That is, the agent must learn the safety relation, while simultaneously maintaining safety.

To explain the importance of this assumption, let us define the set of safe policies corresponding to each F , by

The assumption that the safety values are binary is in contrast to most of the literature, which mostly assumes that these values are numeric, takes values in \mathbb{R} rather than $\{0, 1\}$.

$$F_{safe}(f) = \{f \in F : f(s; (s)) = 1 \text{ a.s. } \forall s \in S\}$$

That is $F_{safe}(f)$ is the set of all policies that would be safe if f was the true safety function, along with the known safe policy. We similarly define the shorthand $F_{safe} = F_{safe}(f^?)$ for the actual set of safe policies. Then, the assumption that F is correctly specified allows us to reason about which policies are safe, since if $\pi \in F_{safe}(f)$ for all possible $f \in F$ that are consistent with our observations so far, we are guaranteed that $\pi \in F_{safe}$. Furthermore, the existence of a known safe policy π_{safe} allows the agent to avoid getting stuck with no known safe action to perform. Note that the definition of F_{safe} could be problem dependent. For example, in the previous cleaning robot example, the policy for almost all states except near the $f^?(s; a) = 0$ margin, whereas with the binary observations we are unsure of safety for all states in between the two observed clusters. Note that this is true even though the numeric observations have noise, while the binary ones don't. In other words, when we receive binary safety feedback, we cannot easily extrapolate safety beyond the observed states.

Obtaining Safety Feedback A key motivation behind using binary safety feedback rather than continuous is for applications where safety is human-defined. In such applications, obtaining safety labels for $(s; a)$ pairs may be expensive or otherwise burdensome. This motivates an active learning style approach to the problem, where we have to explicitly ask for labels of $(s; a)$ pairs, with an additional goal of minimizing the number of times we do so.

In order to make the task of learning given these safety constraints feasible, we model $f^?$ using some class \mathcal{F} of candidate safety functions. This allows us to ensure that $f^?$ is learnable, by using a sufficiently well-behaved class \mathcal{F} . In particular, in our theory we will focus on classes with finite Vapnik-Chervonenkis (VC) dimension (Dudley, 1987). Then, in order to make it feasible to guarantee safety during training, we make the following core assumption on the setting.

Assumption 2. The safety class \mathcal{F} is correctly specified;

Concretely, our model for the labeling process is as follows. We assume access to a labeling oracle, which can be given $(s; a)$ pairs and returns the corresponding values of $f^?(s; a)$. We assume that this oracle can be queried an unlimited number of times in between episodes, and that

it can only be queried with states that have been previously observed. The reason for the latter restriction is that in many applications the state corresponds to some kind of observation of the environment, such as an image, and therefore the total space of states is not a-priori known.

Problem Setup Summary. Given a class of MDPs \mathcal{M} , a policy class Π , a blackbox RL algorithm Alg , a safety function class \mathcal{F} , a desired sub-optimality ϵ , and failure probability δ , we want to propose an online learning algorithm that, for any unknown $M \in \mathcal{M}$ and $f \in \mathcal{F}$, interacts with the MDP over N rounds and returns a policy π such that, with probability at least $1 - \delta$, we have:

1. $\text{SubOpt}(\pi; \epsilon)$ (π is approximately optimal)
2. $\pi \in \Pi_{\text{safe}}$ (the returned policy is safe)
3. $f(s_h^n; a_h^n) = 1$ for all $h \in [H]$ and $n \in [N]$ (the agent never takes unsafe actions during training)

In addition, we would like to establish sample-complexity bounds on the number of episodes needed to ensure the above, in terms of $\epsilon, \delta, \log(1/\delta), H$, and $n_{\text{Alg}}(\cdot)$. Finally, we would like to establish corresponding high-probability bounds on the total number of calls to the labelling oracle, which are lower order than the total sample complexity.

3 LEARNING OF SAFETY VIA ACTIVE LEARNING IN REINFORCEMENT LEARNING

First, we establish some basic definitions and a result that will form the basis for safety learning. For any safety-labeled dataset \mathcal{D} consisting of $(s; a; f(s; a))$ tuples, we define the corresponding version space of safety functions consistent with \mathcal{D} by

$$V(\mathcal{D}) = \{f \in \mathcal{F} : f(s; a) = c \mid (s; a; c) \in \mathcal{D}\}$$

Similarly, for any given $a \in \mathcal{A}$, we can define the corresponding region of disagreement of states where safety of that action is not known given \mathcal{D} by

$$RD^a(\mathcal{D}) = \{s : \exists f, g \in V(\mathcal{D}) \text{ s.t. } f(s; a) \neq g(s; a)\}$$

Note that these are both standard definitions in disagreement-based active learning (Hanneke, 2014).

Next, define the set of policies known to surely be safe given \mathcal{D} as follows:

$$\Pi(\mathcal{D}) = \bigcap_{f \in V(\mathcal{D})} \Pi_{\text{safe}}(f)$$

Note that given Assumption 2, $\Pi(\mathcal{D})$ is always ensured to be non-empty, as it will at least contain π_{safe} .

Given these definitions, we have the following lemma.

Lemma 1. For any safety labeled dataset \mathcal{D} , let

$$U(\mathcal{D}) = \sup_{\pi \in \Pi_{\text{safe}}} P(\exists h \in [H]; \exists a \in \mathcal{A} : s_h \in RD^a(\mathcal{D}))$$

Then, for any $\epsilon \in (0, 1)$, we have

$$\text{SubOpt}(\pi; \epsilon) \leq \text{SubOpt}(\pi; \epsilon(\mathcal{D})) + HU(\mathcal{D})$$

3.1 Challenges with a Naive Approach

This lemma in fact suggests a rather simple, naive approach, based on greedily trying to follow the blackbox RL algorithm, and switching to π_{safe} and querying the safety oracle whenever we reach a state where safety is not known. Before we outline our novel solution that addresses the challenges outlined in the previous section, we consider the limitations of this naive approach.

Given Lemma 1 and Assumption 1, we can easily ensure sub-optimality of at most $\epsilon + HU(\mathcal{D})$ for any target using a naive approach like above, where \mathcal{D} is the labeled dataset incidentally collected following this approach. However, since this approach does nothing explicit to try to learn the safety and shrink the region of disagreement, it is difficult to provide any guarantees on how fast $U(\mathcal{D})$ shrinks.

A second issue is that this approach provides no control on the number of calls to the labelling oracle. Existing analyses from disagreement-based active learning provide bounds on the number of samples needed to shrink regions of disagreement under fixed distributions, but the agent may roll out with a different distribution every episode.

3.2 The SABRE Algorithm

Motivated by Lemma 1, as well as the limitations of the above naive baseline approach, we now present our novel algorithm in Algorithm 1. This algorithm is superficially similar to the baseline approach, in the sense that it builds a labeled dataset \mathcal{D} over a series of rounds, and then estimates an optimal policy following $\text{Alg}(R; \Pi(\mathcal{D}); \epsilon)$. However, the difference is in how the safety-labeled dataset \mathcal{D} is constructed. Instead of successively trying to optimize the environmental reward, and labeling the states we happen observe in the process, our algorithm uses a strategy for constructing \mathcal{D} that explicitly targets $\Pi(\mathcal{D})$.

Our algorithm uses two loops to construct the safety-labeled dataset. The outer loop runs over n epochs. In each epoch it starts with the set Π_n of policies known to be safe at the start of the epoch, and holds this set over the entire epoch. Within each epoch, the inner loop performs B iterations, alternating between the following two steps: (1) (approximately) optimize a policy within n for hitting the region of disagreement with the current \mathcal{D} , and (2) roll out with this policy for episodes to collect additional labeled data to expand \mathcal{D} with. The reason for having these

Algorithm 1 SAFe Binary-feedback REinforcement learning (SABRE)

Input: Number of epochs N , number of iterations per epoch B , number of rollouts per batch m , accuracy parameters $\epsilon_{\text{explore}}$ and ϵ_R , and probability parameters $\epsilon_{\text{explore}}$ and ϵ_R , initial safety-labeled dataset $D_0^{(B)}$ (optional)
 Output: Final policy π

- 1: for $n = 1; 2; \dots; N$ do
- 2: $\pi_n \leftarrow \text{Alg}(\mathbb{R}_n^{(B)}; \epsilon_{\text{explore}}; \epsilon_R)$ // de ne a safe policy class for the current safety labeled dataset $D_n^{(B)}$
- 3: $D_n^{(0)} \leftarrow D_n^{(B)}$
- 4: for $i = 1; 2; \dots; B$ do
- 5: $R_n^{(i)} \leftarrow \text{Alg}(\mathbb{R}_n^{(i)}; \epsilon_{\text{explore}}; \epsilon_R)$ // safety reward that incentivizes visiting region of disagreement
- 6: $\pi_n^{(i)} \leftarrow \text{Alg}(\mathbb{R}_n^{(i)}; \epsilon_{\text{explore}}; \epsilon_R)$ // call blackbox RL method with safety exploration reward
- 7: Roll out with $\pi_n^{(i)}$ for m episodes, and collect all observed states $S_n^{(i)}$
- 8: Expand $D_n^{(i-1)}$ to $D_n^{(i)}$ by labelling all $s \in S_n^{(i)}$ and $a \in A$ such that $R_n^{(i)}(s, a) \geq 0$
- 9: return $\text{Alg}(\mathbb{R}; \epsilon_{\text{explore}}; \epsilon_R)$ // return a safe policy by optimizing the environment reward

two separate loops is that it allows us to derive our formal guarantees to be safe. This holds with certainty, not just with high probability, since by definition it is possible that an improved analysis in future work contains only safe policies, for any obtainable D . Given this, allow the algorithm to be simplified to a single loop. We will focus on establishing approximate suboptimality, along with high-probability bounds on the sample and labeling complexities.

Comparing with the naive approach discussed above, this algorithm follows a strategy for expanding D based on actually optimizing hitting the region of disagreement for Before we give our result, we must give some additional data collection, which is better tailored to explicitly reducing technical assumptions. First, we require that the MDP $(\mathcal{S}, \mathcal{A}, P, U, \gamma)$ is ergodic. Furthermore, the data that is used for expanding D has reasonably low complexity, as follows.

D within each epoch is collected by rolling out with $\pi_n^{(i)}$ policies, which allows stronger control on the number of times the safety labeling oracle will be called.

Finally, we note that this algorithm is fairly generic and abstract, and it is not immediately clear how to apply and scale it to practical scenarios. We provide a detailed discussion of this issue in Appendix A. First, we discuss how the constraint of the blackbox algorithm to policies in (D) may be easily implemented for most kinds of general RL algorithms, as long as we can tractably check whether $(s; a)$ is known to be safe or not given D . Then, we discuss how we may check the possible safety of $(s; a)$ given D for various safety classes of the form $F = \{f : \text{sign}(g(s; a)) \geq 0\}$ for some base class $G \subseteq \mathcal{S} \times \mathcal{A}$. In particular, when G is a linear class safety checking reduces to solving linear programs, and when G is a kernel class safety checking reduces to solving quadratic programs. Finally, we discuss some heuristic approaches for reducing the amount of computation or labeling required for such implementations, or for dealing with more general safety classes based on e.g. generic machine learning methods.

4 THEORETICAL ANALYSIS

We now provide a theoretical analysis of our proposed algorithm. First, by design SABRE only ever takes safe actions during training, and the final returned policy is al-

Assumption 3. There exists some positive integer d such that, for any given set of policies $\{\pi_1, \dots, \pi_d\}$ safe there exists a set of policies $\{\pi_1, \dots, \pi_d\}$ satisfying

$$\frac{1}{H} \sum_{h=1}^H P(s_h \in \mathcal{S}) \leq \frac{1}{H} \sum_{h=1}^H P_i(s_h \in \mathcal{S}) ;$$

We call the smallest d for which this assumption holds the policy-cover dimension of the MDP. We give examples of d for some common MDP classes in Section 4.1.

In addition, we need need to assume a bound on the disagreement coefficient of the distribution induced by any policy, which is a standard joint complexity measure of both distribution and hypothesis class, and is used to derive upper and lower bounds for active learning in binary classification (Hanneke, 2014). Formally, for any $h \in [H]$ and $f \in \mathcal{F}$, let $E_h(f; f^0)$ denote the event that $f(s_h; a) \neq f^0(s_h; a)$ for some $a \in A$. That is, $E_h(f; f^0)$ denotes the event that f and f^0 do not fully agree on the safety of s_h . Then, for any $h \in [H]$, and $r \geq 0$, we define the pseudometric $b_{h,r}$ on \mathcal{F} and the corresponding ball $B_{h,r}(f)$ about f by

$$b_{h,r}(f; f^0) = P(E_h(f; f^0))$$

$$B_{h,r}(f) = \{f' \in \mathcal{F} : b_{h,r}(f; f') \leq r\}$$

Then, for any $h \in [H]$, and $r_0 \geq 0$, we define the

disagreement coefficient $\beta_{h;}(r_0)$ by

$$\beta_{h;}(r_0) = \sup_{r > r_0} \frac{P(\exists f; f \in \mathcal{B}_{h;}(r) : E_h(f; f^0))}{r}$$

Note that clearly by construction $\beta_{h;}(r_0)$ is non-increasing in r_0 . Then, we require the following technical assumption.

Assumption 4. There exists some $\alpha < 1$ such that $\beta_{h;}(0) \leq \alpha$ for all $h \in [H]$ and $\beta_{h;}(0) \leq 2$.

Note that for some problems we may have $\beta_{h;}(0) = 1$.

In this case, we can still obtain a similar PAC bound given a

complex technical condition in terms of the rate of growth of $\beta_{h;}(r)$ as $r \rightarrow 0$, which we present in the appendix. However, we focus here on using Assumption 4 instead since it is much simpler and already covers many important settings, such as linear classifiers with bounded density (Hanneke, 2014). We also refer readers to Hanneke (2014) for a detailed discussion of known results on disagreement coefficients.

Given these assumptions, we are ready to present our main theoretical result.

Theorem 2. Let d_{VC} denote the VC dimension of \mathcal{F} , and let some $\epsilon \in (0; 1)$ be given. Suppose we run the SABRE algorithm with $N = H$, $B = n_B(\epsilon; H; d)$, $m = n_m(\epsilon; H; d; d; d_{VC})$, $\epsilon_{\text{explore}} = \frac{1}{8}H^{-2}$, $R = \frac{1}{2}$, $\epsilon_{\text{explore}} = \frac{1}{4NB}$, and $R = \frac{1}{2}$, for some

$$n_B = O(\log(\epsilon^{-1}) \log(H) d)$$

$$n_m = O(\epsilon^{-1} \log(\epsilon^{-1}) H^3 \log(H)^2 d^2 d_{VC})$$

Then, under Assumptions 1 to 4, the returned policy satisfies $\text{SubOp}(\epsilon; \text{safe})$ with probability at least $1 - \epsilon$.

Ignoring log terms, it follows that the total sample complexity (number of episodes) of our algorithm is at most

$$n_{\text{sample}} = O(H^4 \epsilon^{-1} \log(\epsilon^{-1}) d^2 d_{VC} + H d n_{\text{Alg}}(H^{-2}; \epsilon))$$

Finally, under an additional event with probability at least $1 - \epsilon$, the number of calls to the labeling oracle with the above settings is bounded by

$$n_{\text{label}} = O(A \log(\epsilon^{-1}) \log(\epsilon^{-1}) H^2 d^2 d_{VC})$$

Importantly, Theorem 2 tells us that although we require a sample complexity that depends on H by at least $O(H^4)$ (possibly greater depending on the complexity of the blackbox algorithm), the corresponding dependence on the sampling oracle is much smaller at $O(\log(\epsilon^{-1}) H^2)$. This is very desirable compared with naive baselines as the sample complexity. We also note that finite classes \mathcal{F} have VC dimension at most $\log_2(|\mathcal{F}|)$, so in the case that \mathcal{F} is finite we can replace d_{VC} by $\log(|\mathcal{F}|)$.

Proof Overview Here we provide a brief overview of the proof of Theorem 2. Full proof details, along with a generalized result using a relaxed version of Assumption 4, as discussed above, are provided in the supplement.

First, by Lemma 1 and the guarantees of the blackbox algorithm, if we can ensure $U(D_N^{(B)}) \leq \frac{1}{2}H^{-1}$ with probability at least $1 - \frac{\epsilon}{2}$, then we have $\text{SubOp}(\epsilon; \text{safe})$ with probability at least $1 - \frac{\epsilon}{2}$ by a union bound. Therefore, the proof focuses on establishing this bound.

Now, define

$$G(\epsilon; D) = \frac{1}{H} \sum_{h=1}^H P(\exists a : s_h \geq R D^a(D))$$

and $\epsilon = \frac{1}{4}H^{-3}$. Given Assumption 4, we show that our choice of m ensures that, after rolling out with any fixed policy for m iterations to expand D to D^0 , we will have $G(\epsilon; D^0) \leq \frac{1}{2} \max(\epsilon; G(\epsilon; D))$ with high probability.

Furthermore, given Assumption 3 and any fixed policy π such that $G(\pi; D) \leq \frac{1}{2}$ (which is ensured by our choice of $\epsilon_{\text{explore}}$ since the expected sum of rewards $R_n^{(i)}$ under π is given by $H G(\pi; D_n^{(i)})$) and expand D to D^0 such that $G(\pi; D^0) \leq \frac{1}{2} \max(\epsilon; G(\pi; D))$ at least B times, then at end of the process we have $\sup_{\pi} G(\pi; D)$ with certainty. That is, putting the above together, the inner loop of our algorithm ensures with high probability that $\sup_{\pi} G(\pi; D_n^{(B)}) \leq \epsilon$.

Finally, we show that after expanding D to D^0 to ensure that $\sup_{\pi} G(\pi; D^0) \leq \frac{1}{2}$ at least H times, we are ensured with certainty that at the end of this process we have $U(D) \leq \frac{1}{2}H^{-1}$. This follows, intuitively, because after repeating the above process H times, we will always know safety values for the H actions of any safe policy.

4.1 Examples

Finally we provide some examples that instantiate our theory to some particular classes of MDPs. For each MDP class \mathcal{M} considered below we will provide a bound on, an example blackbox algorithm Alg for this class, its sample complexity n_{Alg} , and the corresponding sample complexity of SABRE. We provide all of these bounds in Table 1, and provide details on how the bounds are derived in the appendix.

Tabular MDPs Tabular MDPs is a simple MDP class where S is finite, and no other structure is assumed. Letting $j \in S$, it is easy to show that $d_{VC} \leq |S|$. In this example, we consider the UCB-VI algorithm (Azar et al., 2017), which is based on value iteration with an exploration reward bonus, and is known to be minimax optimal for this class. Note that Azar et al. (2017) only provided a regret bound; we provide details of how we converted this to a corresponding sample-complexity bound in the appendix.

M	d bound	Alg	n_{Alg}	Corresponding SABRE sample complexity
Tabular MDP	S	UCB-VI	$\mathcal{O}(H^3 S A^2 \log(\frac{1}{\delta}))^4$	$\mathcal{O}(H^8 S^2 A^2 \log(\frac{1}{\delta})^4 d^2 d_{\text{VC}})$
Block MDP	S	HOMER	$\mathcal{O}(H^4 S^8 A^4 \log(\frac{1}{\delta}))^2$	$\mathcal{O}(H^9 S^9 A^4 \log(\frac{1}{\delta})^2 d^2 d_{\text{VC}})$
Low NNR MDP	d_{NNR}	Rep-UCB	$\mathcal{O}(H^5 d^4 A^2 \log(\frac{1}{\delta}))^2$	$\mathcal{O}(H^{10} d^5 A^2 \log(\frac{1}{\delta})^2 d^2 d_{\text{VC}})$

Table 1: Summary of example instantiations of our theory for different MDP classes and corresponding blackbox algorithms Alg. We consider UCB-VI (Azar et al., 2017), HOMER (Misra et al., 2020), and Rep-UCB (Uehara et al., 2021). In each case we give a bound n_{Alg} , and the corresponding sample complexity of SABRE.

Block MDP Block MDP is an MDP class where \mathcal{S} can be general, but the transition dynamics are defined by an latent tabular MDP. Specifically, the observed states are sampled iid given the latent discrete state, and it is assumed that the observed state distributions for each latent state do not overlap.² Let S denote the number of discrete latent states. Then, it can be shown that we again have S . Here we consider the HOMER algorithm of Misra et al. (2020), Wachi and Sui, 2020; Cheng et al., 2019). However, these which works by systematically exploring the latent state space, while simultaneously learning a decoder to map observed to latent states.

Low NNR MDP Finally, we consider the class of MDPs with low non-negative rank (NNR). We define the NNR of an MDP as the smallest integer d_{NNR} such that the transition operator can be written as

$$T(s^0_j | s; a) = \sum_{z=1}^{d_{\text{NNR}}} P(z | j; s; a) P(s^0_j | z)$$

for some latent variable $z \in [d_{\text{NNR}}]$. It can easily be shown that this class generalizes both Block MDP and Tabular MDP, and that we always have $d_{\text{NNR}} \leq S$. For this example we consider the Rep-UCB algorithm (Uehara et al., 2021), which is an algorithm with an optimism-based exploration bonus, which simultaneously learns the low-rank decomposition while exploring.

5 RELATED WORK

There has been a vast body of work on safe RL in recent years; for an overview see for example Garcia and Fernández (2015), Gu et al. (2022), or Brunke et al. (2022). Past work has considered many different approaches, including heuristic deep reinforcement learning methods (Thomas et al., 2021; Luo and Ma, 2021), providing PAC bounds on both sample complexity and number of safety violations (Ding et al., 2021; Hasanzade Zonuzi et al., 2021), safety in the context of transfer learning (Srinivasan et al., 2020), safe of ine RL (Amani and Yang, 2022), safe multi-agent RL (Lu et al., 2021), or providing benchmark environments for safe RL (Ray et al., 2019).

²This ensures that the latent states are fully recoverable from the observed states, and therefore the process is an MDP.

Within the vast body of work on safe RL, there is a sub-area of particular relevance that focuses on guaranteeing safety during all of training. One significant line of work there focuses on the setting where the safety function is unknown and numeric safety feedback is given, and work by leveraging smoothness assumptions on the safety function (Sui et al., 2015; Turchetta et al., 2016; Wachi et al., 2018; Wachi and Sui, 2020; Cheng et al., 2019). However, these papers require numeric feedback and deterministic state transitions, and are limited to continuous control-like settings. A related approach is taken by Berkenkamp et al. (2017), who also consider a similar continuous control-like setting, but define safety based on staying within the region of attraction of the system. A different line focuses on settings where the safety function is known, with diverse approaches including safe versions of value and policy iteration (Chow et al., 2018), reducing to blackbox optimal control problems given differentiability of trajectories (Jin et al., 2021), or using MDP abstractions (Alshiekh et al., 2018; Sirão et al., 2021). However, these works are not relevant when safety must be learned. Other works in this sub-area include methods for low-rank MDPs with linear safety functions and numeric feedback (Amani et al., 2021), settings with cost-based safety functions given a total safety budget (Huang et al., 2022), or using a human-in-the-loop who can take control in real time (Saunders et al., 2018; Turchetta et al., 2020; Peng et al., 2022). However, these approaches all fail to meet our requirement of being able to provably ensure safety during training given of ine binary feedback. Alternatively, Roderick et al. (2021) using “analogy” relationships to expand known sets of safe state, action pairs, but this is limited to settings where such a relationship is available. We also note that one of the existing papers in this sub-area consider the problem of actively acquiring safety feedback, or minimizing the amount of safety feedback required.

Another relevant sub-eld within safe RL considers the very general constrained Markov decision (CMDP) setting (Altman, 1999), where safety is defined by constraints on the total accrued cost $\sum_{h=1}^H c(s_h; a_h)$ given one or more cost function c . These constraints may be bounds on the expected total cost, high probability bounds on the total cost, or other risk bounds on the distribution of the total cost. Note also that cost function in general is unrelated to

Figure 2: Left: Mean return of SABRE on the block MDP task. The red line shows the optimal return. Middle: Mean cumulative number of calls to the labeling oracle over episodes. Right: Mean cumulative number of unsafe actions over episodes. For all plots the error bars correspond to the standard deviation over the 5 replications.

the reward function. Importantly, this is a very general setting than RL settings where we would like to learn classifiers that can model binary safety feedback on state, action (or safety) that are accurate under the state distributions pairs like we consider, and can even model high-probability induced by a wide range of policies.

safety constraints on states alone rather than state, action pairs (Wagener et al., 2021). There is a vast literature that tackles this more general setting, with approaches based on e.g. applying local policy optimization to the system's Lagrangian (Stooke et al., 2020; Chow et al., 2019; Ray et al., 2019), or safe versions local policy search (Zhang et al., 2020). Of particular interest, Cowen-Rivers et al. (2022) provides a method inspired by active learning, but unlike us they use active learning to generically encourage exploration within the region currently known to be safe, whereas we directly query labels for learning the safety relation. However, although these methods generally try to ensure safety during training, they do not actually provide formal guarantees on this, and indeed in their presented results these methods make some safety violations. In addition, again, none of these works consider the problem of actively acquiring safety feedback, or minimizing the amount of safety feedback required.

Finally, another relevant set of literature is on that the agent follows by taking unsafe actions in a disagreement-based active learning in the realizable setting (Hanneke, 2014). In particular, both our algorithm and analysis adopt ideas from the CAL algorithm of Cohn et al. (1994), which is a simple approach that labels whenever we observe an input in the region of disagreement. This algorithm was first formally analyzed by Balcan et al. (2006), for agnostic active learning settings, and Hanneke (2007) showed that its sample- and label-complexity can be bounded in terms of the disagreement coefficient. In addition, Hanneke (2011) establishes similar results for the realizable setting. Other works related to this consider, for example, bounding disagreement coefficients for particular settings (for a detailed overview see Hanneke (2014) and references therein), or providing lower and upper bounds for the problem of realizable active learning (Dasgupta, 2005; Balcan et al., 2007). However, these works all consider standard binary classification settings where data is sampled from a fixed distribution, rather

6 PROOF OF CONCEPT EXPERIMENTS

Finally, we provide some brief “proof of concept” experiments to help highlight the correctness of our theory. The focus here is to demonstrate that SABRE can achieve near-optimal return in a reasonable number of episodes in a non-trivial scenario, while never taking unsafe actions. Code for fully reproducing our experiments is available at <https://github.com/CausalML/SABRE>.

Environment We consider a particular Block MDP scenario, which was introduced in Section 4.1. This scenario has $A = 4$ actions, a time horizon $H = 5$, and $S = 4H + 1$ discrete latent states. The agent receives an observed state along with safety features $(s; a)$ for every action a . The underlying latent state space contains four different paths the agent may take: **unsafe path** where the agent follows by taking unsafe actions in a high reward path where the agent receives high rewards but learns the safety features more slowly; **low reward path** where the agent receives low rewards but learns safety features more quickly; and **safe path** which the agent reaches by following `safe` which is absorbing and gives no reward. The idea of the scenario is that, to safely optimize reward quickly, the agent should first follow the low-reward path to learn the safety function quickly, rather than greedily try to follow the high-reward path where the safety function is learned more slowly. We provide full details of this environment in the Appendix G.

Blackbox Algorithm We use Proximal Policy Optimization as our blackbox RL algorithm Schulman et al. (2017). PPO is a popular empirical RL method, and while not provably efficient, it is quite effective in practice for problems that do not require strategic exploration. We describe the specific details of our PPO implementation in Appendix G.

Safety Class We let the safety class be given by $\mathcal{S} = \{s \mid \exists a, w \in \mathcal{A} : \text{sign}(w \cdot (s; a)) : kwk_1 \leq 1\}$. We note that checking whether $s \in \mathcal{S}$ for any $(s; a)$, as required by SABRE, can be reduced to solving linear programs. We provide details of this in Appendix G.

Results We ran SABRE, implemented with the PPO blackbox algorithm as discussed above on our Block MDP environment, over a total of 7000 episodes. In addition, for comparison we ran the PPO algorithm, directly optimizing reward without any safety considerations. We note that although the total number of episodes is the same for each algorithm, SABRE spends the first 1500 episodes exploring safety while ignoring return, the next 1000 episodes optimizing return, and the next 4500 episodes continuing to roll out with the estimated optimal safe policy. In comparison, the unsafe PPO baseline spends all 7000 episodes optimizing return. We repeated this over 5 random replications, and we present the results in Figure 2, with error bars corresponding to the variance over these replications.

The left figure shows the episodic return against episodes for both SABRE and PPO. In the case of SABRE, as expected, the agent initially receives negative reward as it explores the safety relation along the low reward path, and then at the end once safety is approximately known it computes a near-optimal policy. On the other hand, PPO takes a more direct path, although they both end up with a similarly optimal policy.

The middle figure shows the mean number of calls to the labeling oracle over the same set of episodes. The total number corresponds to only approximately 12% of the total number of state, action pairs encountered. This is very encouraging, since although the block MDP environment has a discrete latent state space, the actual state space is infinite, so there is no absolute maximum number of possible calls to the safety oracle.

Next, on the right we plot the mean number of safety violations of the two methods. As expected, the unsafe PPO baseline takes many unsafe actions, while SABRE never takes any unsafe actions, as guaranteed by our theory.

Finally, we provide a more detailed presentation of results in Appendix G. There, we provide similar results for a range of different values on the number of iterations used for safety exploration versus return optimization. In addition, we provide additional results for an implementation of the naive baseline approach described in Section 3.1 which greedily tries to optimize return using PPO while only taking actions known to be safe, and labeling observed states where safety is unknown at the same frequency as for SABRE. In summary, we find that both SABRE and this baseline approach perform very well in practice. However, SABRE is better able to learn the safety function accurately in fewer rounds of data collection for labeling, and conse-

quently obtain a near-optimal policy faster. This is consistent with our theory, since unlike the baseline it can learn to follow the low reward path where it learns safety functions more quickly.

In summary, SABRE is able to achieve optimal return similar to a standard PPO baseline, while taking no unsafe actions, and only making a tiny number of oracle calls. This supports the theoretical findings described above.

CONCLUSION

We presented a novel algorithm, SABRE, which addresses the problem of safe RL, where the safety function must be learned via binary feedback that is actively acquired in between episodes. Under appropriate assumptions, SABRE is guaranteed to return a policy that only takes safe actions, and to never take unsafe actions during training. In addition, given access to a PAC blackbox RL algorithm for optimizing arbitrary reward functions in the underlying MDP class, it is guaranteed with high probability an approximately-optimal safe policy with polynomial sample complexity. Furthermore, it only requires labels for a relatively minimal number of state, action pairs, with a label complexity that scales as $n^2 \log(n)$ as opposed to a sample complexity that scales as n^4 .

There are multiple avenues for future research. First, improved concepts from realizable active learning could be used to further reduce the label complexity of the algorithm, by using a more sophisticated strategy rather than always labeling states in the region of disagreement (Dasgupta, 2005). Second, it is possible that an improved theoretical analysis could find settings under which the simpler baseline approach described in Section 3.1 has some provable guarantees. Indeed, it is apparent from the additional experimental results in Appendix G that this can be a strong algorithm, so it would be interesting to determine both whether it can obtain similar formal guarantees to SABRE, and also whether this baseline approach has any failure cases. Similarly, an improved theoretical analysis could allow SABRE to be simplified or improved; for example, it may be sufficient to have a single loop where the policy class is updated every iteration, or it may be possible to simultaneously perform reward optimization and safety learning. Future work may also consider how to relax the realizability assumption, and provide some kind of guarantees when function approximation is used to model the environment. This is important, since a major limitation of our work is that the safety guarantees only hold as long as we satisfy Assumption 2. In addition, it may consider how to implement SABRE with more complex choices of than the linear classes we considered in our experiments; for example, it may consider classes defined by kernels or neural nets. Finally, it would be exciting to see applications of our theory to practical safety problems.

- lations. *Advances in Neural Information Processing Systems* 34:25621–25632, 2021.
- D. Misra, M. Henaff, A. Krishnamurthy, and J. Langford. Kinematic state abstraction and provably efficient rich-observation reinforcement learning. *International conference on machine learning*, pages 6961–6971. PMLR, 2020.
- Z. Peng, Q. Li, C. Liu, and B. Zhou. Safe driving via expert guided policy optimization. *Conference on Robot Learning* pages 1554–1563. PMLR, 2022.
- A. Ray, J. Achiam, and D. Amodei. Benchmarking safe exploration in deep reinforcement learning. *arXiv preprint arXiv:1910.01708* 2019.
- P. Rebeschini. Covering numbers bounds for rademacher chaining, November 2020.
- M. Roderick, V. Nagarajan, and Z. Kolter. Provably safe pac-mdp exploration using analogies. *International Conference on Artificial Intelligence and Statistics*, pages 1216–1224. PMLR, 2021.
- W. Saunders, G. Sastry, A. Stuhler, and O. Evans. Trial without error: Towards safe reinforcement learning via human intervention. *Proceedings of the 17th International Conference on Autonomous Agents and Multi-Agent Systems*, pages 2067–2069, 2018.
- J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* 2017.
- T. D. Simão, N. Jansen, and M. T. Spaan. Always safe: Reinforcement learning without safety constraint violations during training. *Proceedings of the 20th International Conference on Autonomous Agents and Multi-Agent Systems* pages 1226–1235, 2021.
- K. Srinivasan, B. Eysenbach, S. Ha, J. Tan, and C. Finn. Learning to be safe: Deep rl with a safety critic, 2020.
- A. Stooke, J. Achiam, and P. Abbeel. Responsive safety in reinforcement learning by pid lagrangian methods. *International Conference on Machine Learning*, pages 9133–9143. PMLR, 2020.
- Y. Sui, A. Gotovos, J. Burdick, and A. Krause. Safe exploration for optimization with gaussian processes. *International Conference on Machine Learning*, pages 997–1005. PMLR, 2015.
- G. Thomas, Y. Luo, and T. Ma. Safe reinforcement learning by imagining the near future. *Advances in Neural Information Processing Systems* 34:13859–13869, 2021.
- M. Turchetta, F. Berkenkamp, and A. Krause. Safe exploration in finite markov decision processes with gaussian processes. *Advances in Neural Information Processing Systems* 29:4312–4320, 2016.
- M. Turchetta, A. Kolobov, S. Shah, A. Krause, and A. Agarwal. Safe reinforcement learning via curriculum induction. *Advances in Neural Information Processing Systems* 33:12151–12162, 2020.
- M. Uehara, X. Zhang, and W. Sun. Representation learning for online and of line rl in low-rank mdp. *arXiv preprint arXiv:2110.04652* 2021.
- A. Wachi and Y. Sui. Safe reinforcement learning in constrained markov decision processes. *International Conference on Machine Learning*, pages 9797–9806. PMLR, 2020.
- A. Wachi, Y. Sui, Y. Yue, and M. Ono. Safe exploration and optimization of constrained mdp using gaussian processes. *Thirty-Second AAAI Conference on Artificial Intelligence* 2018.
- M. C. Wagener, B. Boots, and C.-A. Cheng. Safe reinforcement learning using advantage-based intervention. *International Conference on Machine Learning*, pages 10630–10640. PMLR, 2021.
- Y. Zhang, Q. Vuong, and K. Ross. First order constrained optimization in policy space. *Advances in Neural Information Processing Systems* 33:15338–15349, 2020.

A DISCUSSION OF PRACTICAL IMPLEMENTATIONS OF SABRE

A.1 Implementation of Blackbox RL Algorithm with Policy Constraints

One concern an astute reader may have about our requirement of the blackbox algorithm in Assumption 1 is that it not only needs to be able to optimize any reward function, but it needs to be able to do so over any given policy while only taking actions in \mathcal{A}^0 . While this may seem like a major requirement and limitation, we explain here why this is not so, and provide a general recipe for dealing with these constraints, as long as we can check whether any given action is known to be safe in any given state given any safety-labeled dataset.

First, note that the policy sets \mathcal{P}^0 that we require Alg to be actually able to work with are all of the form $\mathcal{P}(D)$. These all look like \mathcal{P} , except with some constraints added on which actions are allowed in any given state based on whether or not that state, action pair is known to be safe or not given D . Now, for any given D , let us consider the $\text{MDP}(D)$, which is defined identically to the actual MDP , except whenever the agent would take an action a in states that is not known to be safe according to D the MDP transitions following $T(j; s; \text{safe}(s))$ rather than $T(j; s; a)$. Next, note that if we add a layer of abstraction between the agent and the MDP which checks any action the agent takes, and converts it to $\text{safe}(s)$ if it is not known to be safe, then interacting with $\text{MDP}(D)$ via this abstraction layer is equivalent to interacting with $M(D)$.

Now, let π^* be an ϵ -optimal policy over \mathcal{P} for $M(D)$, and $\tilde{\pi}^?$ be the policy given by following π^* , and converting its action to $\text{safe}(s)$ whenever the action is not known to be safe given D . By construction, following $\tilde{\pi}^?$ in $M(D)$ is equivalent to following π^* in M , and also by construction the optimal return V^* over $\mathcal{P}(D)$ is the same as the optimal return V^* over \mathcal{P} . Therefore, it is clear that we have $\text{SubOpt}(\tilde{\pi}^?; \mathcal{P}(D)) \leq \epsilon$. Furthermore, as long as the class is not defined in some absurd or pathological way, we should have $\mathcal{P} \in \mathcal{M}$, in which case we would also have $\mathcal{P}(D) \in \mathcal{M}$. That is, $\tilde{\pi}^?$ would be an ϵ -optimal policy over $\mathcal{P}(D)$ for M . Note that the above condition that $\mathcal{P} \in \mathcal{M}$ is trivial for example if \mathcal{M} consists of all policies, or all deterministic policies.

Given the above reasoning, as long as the $\text{MDP}(D)$ is still within the class \mathcal{M} , we can implement the blackbox algorithm by estimating an ϵ -optimal policy for $M(D)$ over $\mathcal{P}(D)$, using the kind of abstraction layer discussed above. Since $M(D) \in \mathcal{M}$, the blackbox PAC guarantees of Assumption 1 will still apply to this estimation. The returned policy will then correspond to an optimal policy over $\mathcal{P}(D)$, given by converting potentially unsafe actions to $\text{safe}(s)$ as required. That being said, this is not necessarily the best approach for implementing since optimizing over the duplicated $\text{safe}(s)$ actions may lead to both computational and statistical inefficiencies. More efficient implementations will be dependent on M and the choice of blackbox algorithm (we discuss a particular case for a better implementation of PPO with arbitrary $\mathcal{P}(D)$ in our experimental details).

Finally, we note that the above argument relied on the assumption $M(D) \in \mathcal{M}$ was still in the class \mathcal{M} , so the PAC guarantees of the algorithm would still apply to $M(D)$. This is definitely the case for all of the MDP classes considered in Section 4.1. For example, in the case of low non-negative rank MDPs, the $\text{MDP}(D)$ will have non-negative rank less than or equal to that of M . To see this, note that for any $(j; s; a)$ that is known to be safe, we still have the decomposition $T(j; s; a) = \sum_{z=1}^d P(z; j; s; a) P^0(j; z)$, for some P and P^0 and for any $(j; s; a)$ that is not known to be safe we have $T(j; s; a) = \sum_{z=1}^d P(z; j; s; \text{safe}(s)) P^0(j; z)$, which is a decomposition of the same rank. Similar logic easily extends to Tabular and Block MDPs.

A.2 Implementation of Safety Checking

Next, we discuss how we can determine the possible safety of a given pair given an arbitrary safety labeled dataset D . Here we will focus on safety classes that take the form $F = \{f(s; a) \geq \gamma \mid \text{sign}(g(s; a)) : g \in G\}$ for some base numeric class $G \subseteq \mathbb{R}$. Note that this is without loss of generality, since for any $G \subseteq \mathbb{R}$ we have $F = \{f(s; a) \geq \gamma \mid \text{sign}(g(s; a)) : g \in G\}$. In general given F specified in such a way, we can reduce the problem of checking the possible safety of $(s; a)$ given $D = \{(s_1; a_1; c_1), \dots, (s_n; a_n; c_n)\}$ to computing the maximum and minimum possible values of $g(s; a)$ for $g \in G$ such that $g(s_i; a_i) c_i \geq 0$ for all $i \in [n]$. If both minimum and maximum are non-negative then the action is known to be safe, if both are negative then the action is known to be unsafe, and otherwise we have $s \in \text{RD}^a(D)$. We discuss below some specific cases where this is a tractable optimization problem.

First, we consider the case that G is a linear class. Specifically, suppose that $G = \{w^T (s; a) \mid w \in \mathbb{R}^d\}$ for some fixed embedding function $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$, and some $d \geq 2$. In this case, the above safety checking problem reduces to computing the maximum and minimum possible values of $w^T \phi(s; a)$ such that $w^T \phi(s_i; a_i) c_i \geq 0$ for all $i \in [n]$, and

$\sum_{j=1}^d w_j = 1$ for all $j \in [d]$. We note that this is a linear program with $2d$ constraints and variables. Therefore, with linear G , checking safety reduces to linear programming. This is very appealing, since linear programming problems are very quick and tractable to solve, which suggests that such linear classes with an appropriate embedding function may be very practical.

For some scenarios, the above approach based on a linear safety class with a fixed embedding function may not be reasonable. One approach in such settings may be to treat a safety features embedding functions such that we can apply this method. For example, we may try to learn to be able to linearly classify observed training data using deep learning. We could then treat these learnt embeddings as the true ones, and hope that the sign of $\hat{g}(s; a) = \sum_{i=1}^d w_i g_i(s; a)$ contains $f^?$. Alternatively, we could be more robust by incorporating the error in these embeddings within our safety class; for example, we could consider $\hat{g}(s; a) = \sum_{i=1}^d w_i (g_i(s; a) + \epsilon_i)$ with $\epsilon_i \in [-\kappa, \kappa]$ for some maximum and norm on the error. In the latter case, if we modeled the error by a reproducing kernel Hilbert space (RKHS), then checking the safety of an $(s; a)$ could be reduced to solving some quadratically constrained linear programs (QCLPs). To see this, note that by the representer theorem, we can optimize over the form $g(s; a) = \sum_{i=1}^d \alpha_i K((s; a); (s_i; a_i))$ without loss of generality, in which case $\|g\|_K^2 = \sum_{i,j} \alpha_i \alpha_j K((s_i; a_i); (s_j; a_j))$. Although quadratic programs are more expensive to solve than linear programs, this still may be a reasonable approach. Alternatively, if we used a generic class for modeling the error, such as neural networks, we may be able to solve the corresponding min/max optimization problems using Lagrangian-based approaches.

An alternative to the above linear class would be to consider kernel-based classes within a norm-constrained ball of a reproducing kernel Hilbert space, which are much more expressive. If we use a universal kernel (such as the Gaussian kernel), then this approach has a slight complication that all all possible labellings are feasible, so all points will always be in the region of disagreement. Note that this is still an issue if we have a norm-constrained since $\hat{g}(s; a) = \sum_{i=1}^d \alpha_i g_i(s; a)$ with $\sum_{i=1}^d \alpha_i^2 \leq M$ for every $M; M \geq 0$; this latter fact follows because $\text{sign}(c g(s; a)) = \text{sign}(g(s; a))$ for every $c > 0$. However, we could get around this by enforcing a minimum margin size in the separation of the safe and unsafe points; that is, we could enforce that $|\hat{g}(s; a)| \geq \gamma$ for some γ^2 such that $\text{sign}(\gamma^2(\cdot)) = f^2(\cdot)$ and $\|g\|_K \leq M$. Then, we can determine the possible safety of $(s; a)$ by computing the minimum and maximum possible values α_i such that $|\hat{g}(s; a)| \geq \gamma$ for all safety-labelled triplets $(s_i; a_i; c_i)$, and $\|g\|_K \leq M$. Then, following the representer theorem, as described above, these min and max problems reduce to QCLPs.

A.3 Additional Heuristics

Finally, we discuss some additional heuristics that may be useful for practical implementations of SABRE.

In the case that we use a class such that safety checking reduces to solving some convex optimization problems, we may apply some heuristics to reduce the number of such problems that need to be solved, or to simplify these problems. As an example, we could justify that some new point $(s; a)$ is known to be safe/unsafe by finding feasible dual solutions to the max/min problems with the same sign; then, for example, if we cache previous dual solutions, then we could check the sign/feasibility of these solutions first with the new $(s; a)$ to possibly avoid re-solving. In addition, we may improve scalability by identifying redundant labeled triplets $(s_i; a_i; c_i)$ that do not change the set of feasible solutions, and removing these; for example, it is easy to argue that $(s_i; a_i)$ is not in the region of disagreement given all other labeled triplets, then $(s_i; a_i; c_i)$ is redundant. This may be helpful in reducing the number of constraints needed for solving the required convex optimization problems. Also, when we are required to label a batch $(s; a)$ at each stage of the inner loop of our algorithm, we may apply heuristics to decide what order to label these in, in order to maximize the chance that the safety of later pairs becomes known during the labeling so fewer labels are needed.

In addition, we could apply heuristic methods for using classes for which safety checking does not easily reduce to solving some tractable convex optimization problems. For example, if we fit the function based on the labeled data using some generic class and method that permits uniform confidence intervals, then we could decide on the possible values of $f^?(s; a)$ based on the confidence interval of $\hat{g}(s; a)$. As another example, we could consider using ensemble methods to estimate $f^?$, and decide on the possible values of $\hat{g}(s; a)$ based on the range of predictions within the ensemble. We caution that such heuristic methods could represent a departure from our theory, and therefore our safety guarantees may no longer hold. However, they may be practical and perform well, and the safety guarantees may still mostly hold in practice if the heuristics are reasonably accurate.

B PROOF OF LEMMA 1

Proof. First, let some arbitrary $\pi \in \arg\max_{\pi \in \Pi_{\text{safe}}} V(\pi)$ be given, and let us define a policy $\tilde{\pi}(D)$ as follows:

- If $s_h \in \text{RD}^a(D)$ for any $a \in \mathcal{A}$, then $\tilde{\pi}(D)(\cdot | s_h) = \pi(\cdot | s_h)$
- Otherwise, $\tilde{\pi}(D)(\cdot | s_h) = \pi_{\text{safe}}(\cdot | s_h)$

Note that by construction we have $\tilde{\pi}(D) \in \Pi_{\text{safe}}$, since it only ever takes actions that are known to be safe. In addition, let E denote the event where $s_h \in \text{RD}^a(D)$ for any $a \in \mathcal{A}$ or $h \in [H]$. Note that under event E , it follows $\tilde{\pi}(D)$ takes actions identically to π . Therefore, we have

$$\text{SubOpt}(\tilde{\pi}(D); \Pi_{\text{safe}}) = \mathbb{P}(E) \text{SubOpt}(\pi; \Pi_{\text{safe}}) + \mathbb{P}(E^c) \text{SubOpt}(\pi_{\text{safe}}; \Pi_{\text{safe}});$$

since the difference in the sum of rewards received between π and $\tilde{\pi}(D)$ can never be greater than ϵ . Furthermore, by the definition of $U(D)$ we have $\text{SubOpt}(\pi_{\text{safe}}; \Pi_{\text{safe}}) = U(D)$ for every $D \in \mathcal{D}$ so therefore

$$\text{SubOpt}(\tilde{\pi}(D); \Pi_{\text{safe}}) = \mathbb{P}(E) \text{SubOpt}(\pi; \Pi_{\text{safe}}) + \mathbb{P}(E^c) U(D);$$

Then, given this, we can bound

$$\begin{aligned} \text{SubOpt}(\pi; \Pi_{\text{safe}}) &= V(\pi) - V(\pi_{\text{safe}}) \\ &= V(\tilde{\pi}(D)) - V(\pi_{\text{safe}}) + \text{SubOpt}(\tilde{\pi}(D); \Pi_{\text{safe}}) \\ &\quad - \text{SubOpt}(\pi; \Pi_{\text{safe}}) + \text{SubOpt}(\tilde{\pi}(D); \Pi_{\text{safe}}) \\ &= \text{SubOpt}(\pi; \Pi_{\text{safe}}) + \text{SubOpt}(\tilde{\pi}(D); \Pi_{\text{safe}}) - \text{SubOpt}(\pi; \Pi_{\text{safe}}) \\ &\quad + \text{SubOpt}(\tilde{\pi}(D); \Pi_{\text{safe}}) - \text{SubOpt}(\pi_{\text{safe}}; \Pi_{\text{safe}}) \\ &= \text{SubOpt}(\pi; \Pi_{\text{safe}}) + \text{SubOpt}(\tilde{\pi}(D); \Pi_{\text{safe}}) - U(D); \end{aligned}$$

where the second equality follows by adding and subtracting $V(\pi_{\text{safe}})$ and plugging in the suboptimality definition, the first inequality follows since $\tilde{\pi}(D) \in \Pi_{\text{safe}}$ by construction, and the second follows from the previous bounds. This is our required bound, so we can conclude. \square

C PROOF OF THEOREM 2

Before we present the main proof, we present some additional lemmas from which the proof can be built. In these lemmas, we will define

$$\text{RD}(D) = \bigcup_{a \in \mathcal{A}} \text{RD}^a(D);$$

for every safety-labeled dataset D . Also, let

$$G(\pi; D) = \frac{1}{H} \sum_{h=1}^H \mathbb{P}(s_h \in \text{RD}(D));$$

for any policy π and safety-labeled dataset D .

First, the following lemma provides a guarantee on the behavior of the inner loop of SABRE, as long as sufficiently large.

Lemma 3. Let Assumption 3 be given. Suppose that for some given $n \in \mathbb{N}$ in the execution of Algorithm 1, we have the events

$$G(\pi_n^{(i)}; D_n^{(i-1)}) \geq \frac{1}{2} \sup_{\pi \in \Pi_{\text{safe}}} G(\pi; D_n^{(i-1)}) - \frac{1}{2} \epsilon \quad \forall i \in [B]; \quad (1)$$

and

$$G(\pi_n^{(i)}; D_n^{(i)}) \geq \frac{1}{2} \max_{\pi \in \Pi_{\text{safe}}} G(\pi; D_n^{(i)}) - \frac{1}{2} \epsilon \quad \forall i \in [B]; \quad (2)$$

for some $\epsilon \in (0, 1)$. In addition, suppose that $\epsilon \leq \frac{1}{8d \log_2(\frac{1}{\epsilon})}$. Then, under these conditions, we are ensured that

$$\sup_{\pi \in \Pi_{\text{safe}}} G(\pi; D_n^{(B)}) \geq \frac{1}{2} \max_{\pi \in \Pi_{\text{safe}}} G(\pi; D_n^{(0)}) - \frac{1}{2} \epsilon;$$

Next, the following lemma analyzes the outer loop of SABRE, under the assumptions and events detailed in Lemma 3.

Lemma 4. Suppose that the conditions of Lemma 3 holds for every $n \in [N]$, with $\frac{1}{2}H^{-2}$ for some $\delta \in (0, \frac{1}{2})$. Suppose also, that $\delta \leq H$. Then, the final dataset $D_N^{(B)}$ satisfies

$$U(D_N^{(B)}) \leq \delta$$

Next, we note that in order to apply the above lemmas, we need to be able to establish the conditions of Lemma 3. For (1), we can rely on the guarantees of the blackbox algorithm. However, for (2), we can instead appeal to an additional lemma. For this lemma, we introduce the following assumption, which is a weaker version of Assumption 4.

Assumption 5. There exists some continuous, non-increasing function $\phi: \mathbb{R}^+ \rightarrow \mathbb{R}^+$ such that

$$\phi(r) \leq h; \quad \phi(r) \geq 8r > 0; \quad h \leq H; \quad \phi \text{ safe:}$$

Furthermore, this function satisfies $\phi(r) \leq 1$, and $\lim_{r \rightarrow 0} \phi(r) = 0$

Furthermore, we will define the function $\phi_{\min}: \mathbb{R}^+ \rightarrow \mathbb{R}^+$, according to

$$\phi_{\min}(x) := \inf_{r > 0: \phi(r) \leq x} \phi(r)$$

We note that given Assumption 5 $\phi_{\min}(x) > 0$ must be non-empty for every $x > 0$, so this is a well-defined function. Also, we note that in the case that $\phi(0) < 1$, we clearly have $\phi_{\min}(x) \leq \phi(0)$ for every $x > 0$, since ϕ is non-increasing.

Lemma 5. Let Assumption 4 be given. Consider some arbitrary $n \in [N]$ and $\delta \in [B]$ in the running of Algorithm 1. Then, as long as $\delta \leq \phi_{\min}(\delta; \min; d_{VC}; \delta; H)$, for some

$$\phi_{\min}(\delta; \min; d_{VC}; \delta; H) = O\left(\frac{1}{\phi_{\min}(\delta)} \log(1/\delta) \log(H)\right)$$

we have

$$G(\wedge_n^{(i)}; D_n^{(i)}) \leq \frac{1}{2} \max_n \left\{ G(\wedge_n^{(i)}; D_n^{(i-1)}) \right\}$$

with probability at least $1 - \delta$.

Given the above lemmas, we can provide the main proof.

Proof of Theorem 2 First, we note that applying Lemma 5 with $\delta = \frac{1}{4}H^{-3}$ and $\delta = \frac{1}{4}H^{-3}$, the corresponding choice of m given by this lemma ensures that Lemma 5 holds at each iteration of the algorithm with probability at least $1 - \delta$. Then, applying a union bound over the $n \in [N]$ events, we have that (2) holds for every $n \in [N]$ with probability at least $1 - 4\delta$.

Next, given the guarantee of Alg, and the setting of δ_{explore} and noting that $HG(\delta; D_n^{(i-1)})$ corresponds to the expected sum of rewards policy under $\mathcal{R}_n^{(i)}$, we are ensured with probability at least $1 - \delta_{\text{explore}}$

$$HG(\wedge_n^{(i)}; D_n^{(i-1)}) \leq \sup_n HG(\delta; D_n^{(i-1)}) \leq \frac{1}{8}H^{-2}$$

$$\implies G(\wedge_n^{(i)}; D_n^{(i-1)}) \leq \sup_n G(\delta; D_n^{(i-1)}) \leq \frac{1}{2}$$

at any given $n \in [B]$ and $n \in [N]$. Then, by a union bound over all $n \in [B]$ calls to Alg during the exploratory loops of the algorithm, and noting the value of δ_{explore} we have that (1) holds for every $n \in [N]$ and $\delta = \frac{1}{4}H^{-3}$ with probability at least $1 - 4\delta$. Then, as long as δ is at least $\frac{1}{4}H^{-3}$, the conditions of Lemma 3 hold for every $n \in [N]$, so by Lemma 4 we have

$$U(D_N^{(B)}) \leq \frac{1}{2}H^{-1}$$

under the above high probability events. Then, by Lemma 1, we have

$$\begin{aligned} \text{SubOp}(\hat{\pi}; \text{safe}) & \leq \text{SubOp}(\hat{\pi}; (D_N^{(B)})) + \text{HU}(D_N^{(B)}) \\ & \leq \text{SubOp}(\hat{\pi}; (D_N^{(B)})) + \frac{1}{2} \end{aligned}$$

Furthermore, following the high-probability guarantee of the `nal calAlg`, and the settings of r and R , we have

$$\text{SubOp}(\hat{\pi}; (D_N^{(B)})) \leq \frac{1}{2};$$

with probability at least $1 - \frac{1}{2} = \frac{1}{2}$. Then, combining the above three high probability events, and the previous two bounds, a union bound gives us

$$\text{SubOp}(\hat{\pi}; \text{safe}) \leq \frac{1}{2};$$

with probability at least $\frac{1}{2}$, as required.

Next, let us analyze the above choices `of` and `B`. Plugging in the settings for `n` and `B`, the above choices `of` and `B` satisfy

$$\begin{aligned} m & = O\left(\frac{1}{8} H^3 \log(H) \min(H^3 = 128)^2 d_{VC} \log\left(\frac{1}{4} (4NB)\right)\right) \\ & = O\left(\frac{1}{8} \log\log\left(\frac{1}{4} H^3 \log(H)^2 \log\log(H) \log(d)\right) \min(H^3 = 128)^2 d_{VC} \log\left(\frac{1}{4} H^3\right)\right) \\ & = O\left(\frac{1}{8} \log\left(\frac{1}{4} H^3 \log(H)^3 \log(d)\right) \min(H^3 = 128)^2 d_{VC} \log\left(\frac{1}{4} H^3\right)\right) \end{aligned}$$

and

$$B = O\left(\log\left(\frac{1}{4} H^3\right) \log(H) d\right);$$

Then, throwing away log-factors, and noting again that $H \leq H$, the total number of iterations is given by

$$\begin{aligned} n_{\text{sample}}(\hat{\pi}; H; d; d; d_{VC}) & = NB \cdot n_{\text{Alg}}\left(\frac{1}{8} H^2; \frac{1}{4} (4NB)\right) \\ & \quad + O\left(\frac{1}{8} H^3 \min(H^3 = 128)^2 d_{VC} \log\left(\frac{1}{4} H^3\right)\right) \\ & = O\left(H d \cdot n_{\text{Alg}}\left(H^2; \frac{1}{4} (4NB)\right) + H^3 \min(H^3 = 128)^2 d_{VC} \log\left(\frac{1}{4} H^3\right)\right); \end{aligned}$$

where in the second equality we apply the assumption that $n_{\text{Alg}}(\cdot; \cdot)$ is polynomial in $\frac{1}{8} H^2$ and $\log\left(\frac{1}{4} (4NB)\right)$, so therefore hiding log terms we have $n_{\text{Alg}}\left(\frac{1}{8} H^2; \frac{1}{4} (4NB)\right) = O\left(n_{\text{Alg}}\left(H^2; \frac{1}{4} (4NB)\right)\right)$.

Next, let us analyze the number of calls to the labeling oracle. Following the proof of Lemma 5, for any given labeling at most k states for each $h \in [H]$, for some

$$k = O\left(\min(H^3 = 128)^2 d_{VC} \log\left(\frac{1}{4} H^3\right) \log(H)\right);$$

we are ensured that, with probability at least $1 - \frac{1}{2}$, the probability of encountering additional $2k$ $\text{RD}(D)$ for each $h \in [H]$ is either reduced by a factor of 4, or that probability was already less than $\frac{1}{2}$. This implies that after at most

$$O\left(H \log\left(\frac{1}{4} H^3\right) \min(H^3 = 128)^2 d_{VC} \log\left(\frac{1}{4} H^3\right) A\right)$$

queries to the labeling oracle in a given iteration, we are ensured that the probability of encountering $2k$ $\text{RD}(D)$ is at most $O\left(\frac{1}{4}\right)$ for every $h \in [H]$, with probability at most $\frac{1}{2}$. Then, applying a union bound over the rollouts in a single iteration with $m = \frac{1}{2} m$, after the above maximum number of

$$O\left(H \log(m) \min(H^3 = 128)^2 d_{VC} \log\left(\frac{1}{4} H^3\right)^2 A\right)$$

queries to the labeling oracle, we will make no further queries in that iteration with probability at least $\frac{1}{2}$. Given this, replacing δ in this bound with $\frac{\delta}{2}$, and applying a union bound over all iterations, with probability at most $\frac{1}{2}$ we make no more than

$$\begin{aligned} & \mathcal{O}(NBH \log(m) \log(NB) \min(\delta, \frac{\delta}{2})^2 d_{VC} \log(\frac{1}{\delta})^2 A) \\ & = \mathcal{O}(H^2 \log(\frac{1}{\delta}) d \min(\delta, \frac{\delta}{2})^2 d_{VC} \log(\frac{1}{\delta})^2 A) \end{aligned}$$

total calls to the labeling oracle.

Finally, we note that for the simplified version of the results presented in the main paper, we assumed that Assumption 4, which is stronger than Assumption 5, and implies that we have $\delta(x) \geq \delta$ for all $x > 0$. Therefore, we can obtain the result in the main paper by replacing $\delta(x)$ with δ everywhere. □

D PROOFS OF ADDITIONAL LEMMAS

D.1 Proof of Lemma 3

Proof. First, for any $i \in \{0, 1, \dots, B\}$, let

$$\delta_i = \sup_{D_n^{(i)}} G(\cdot; D_n^{(i)}):$$

Under the events of the lemma statement, we will argue that

$$\delta_{8d+i} \leq \max_{i \in \{0, \dots, B\}} \left(\frac{1}{2} \delta_i \right); \quad (3)$$

for every $i \in \{0, \dots, B\}$. Given this result the overall lemma easily follows, by chaining this result, where $\delta = \frac{\delta}{2}$. Under the assumption that δ_i is non-increasing, this gives us $\delta_{8d+i} \leq \frac{1}{2} \delta_i$. Note that the last inequality follows since our choice of δ_i satisfies $\delta_i \geq \delta_{i+1}$.

Now, let $\tilde{\pi}_1, \dots, \tilde{\pi}_d \in \Pi_n$ be the policy cover of Π_n ensured by Assumption 3. Note that for any $D_n^{(i)}$ and $i \in \{0, \dots, B\}$, we have

$$\begin{aligned} \sum_{k=1}^d G(\tilde{\pi}_k; D_n^{(i-1)}) - G(\tilde{\pi}_k; D_n^{(i)}) &= \sum_{k=1}^d \frac{1}{H} P_{\tilde{\pi}_k}(s_h \geq RD(D_n^{(i-1)}) \wedge n RD(D_n^{(i)})) \\ &\quad - \frac{1}{H} P(s_h \geq RD(D_n^{(i-1)}) \wedge n RD(D_n^{(i)})) \\ &= G(\tilde{\pi}_k; D_n^{(i-1)}) - G(\tilde{\pi}_k; D_n^{(i)}); \end{aligned}$$

where the first and last equalities follow because $\mathbb{R}^d \setminus \mathbb{R}^d = \emptyset$.

Next, let some arbitrary $j \in \{0, \dots, B\}$ be given, and suppose that (3) does not hold. Given this, it must be the case that $\delta_{i+j} > \frac{1}{2} \delta_{i+j-1}$ for all $j \in \{2, \dots, B\}$, so for the corresponding iterations (2) can be strengthened to

$$G(\wedge_n^{(i+j)}; D_n^{(i+j)}) \leq \frac{1}{2} G(\wedge_n^{(i+j)}; D_n^{(i+j-1)}) \quad \forall j \in \{2, \dots, B\}; \quad (4)$$

Then, for every $j \in [8d]$ we have

$$\begin{aligned} \sum_{k=1}^j & G(\sim_k; D_n^{(i+j-1)}) - G(\sim_k; D_n^{(i+j)}) - G(\wedge_n^{(i+j)}; D_n^{(i+j-1)}) + G(\wedge_n^{(i+j)}; D_n^{(i+j)}) \\ & \geq \frac{1}{2} G(\wedge_n^{(i+j)}; D_n^{(i+j-1)}) \\ & \geq \frac{1}{2} \frac{1}{i+j-1} - \frac{1}{4} \\ & \geq \frac{1}{2} \frac{1}{i+8d} - \frac{1}{4} \\ & > \frac{1}{2} \max\left\{\frac{1}{n}, \frac{1}{2} \frac{1}{i}\right\} - \frac{1}{4} \\ & \geq \frac{1}{4} \max\left\{\frac{1}{n}, \frac{1}{2} \frac{1}{i}\right\} \\ & \geq \frac{1}{8} \frac{1}{i}; \end{aligned}$$

where the strict inequality step follows from the assumption that (3) doesn't hold. Then, summing and noting that this is a telescoping sum, we have

$$\begin{aligned} \frac{1}{8} (8d) \frac{1}{i} & < \sum_{k=1}^j G(\sim_k; D_n^{(i)}) - G(\sim_k; D_n^{(i+8d)}) \\ & < \sum_{k=1}^j G(\sim_k; D_n^{(i)}) \\ & < d \sup_{2 \leq n} G(\cdot; D_n^{(i)}) \\ & = d \frac{1}{i}; \end{aligned}$$

But this is impossible, because we cannot have $\frac{1}{8} < \frac{1}{i}$. Therefore, we have proved by contradiction that (3) holds. Also, since we argued the above for an arbitrary $B \in [8d]$, (3) must hold for every $s \in [H]$, so we are done. □

D.2 Proof of Lemma 4

Proof. We will prove that under these conditions, we are ensured that

$$\sup_{2 \leq n \leq \frac{n}{H}} P_{\text{safe}}(s_h \in [n] : s_h \in \text{RD}(D_n^{(B)})) \leq \frac{n}{H}; \tag{5}$$

for every $n \in [H]$. The required result then follows from this by plugging in $n = H$, since

$$U(D) = \sup_{2 \leq n \leq H} P_{\text{safe}}(s_h \in [H] : s_h \in \text{RD}(D));$$

and because $U(D_N^{(B)});_{\text{safe}} \leq U(D_H^{(B)})$, as we assumed $n \leq H$.

In the remainder of this proof, we establish (5) for every $n \in [H]$ by forward induction on n .

Base Case

First, for the base case, we note that after the first inner loop, by Lemma 3 we are ensured that

$$\sup_{2 \leq n \leq H} \frac{1}{n} \sum_{h=1}^j P(s_h \in \text{RD}(D_1^{(B)})) \leq \frac{1}{2H^2} + \frac{1}{H^2};$$

Therefore, by simple algebra, for every $\delta \in [0, 1]$ we have

$$P(s_1 \in \text{RD}(D_1^{(B)})) \geq \frac{1}{H} \delta$$

Now, for any policy π , the distribution π_1 is always the same (it is always equal to the initial state distribution), so therefore the previous inequality is equivalent to

$$\sup_{\pi \text{ safe}} P(s_1 \in \text{RD}(D_1^{(B)})) \geq \frac{1}{H} \delta$$

which establishes the base case.

Inductive Case

Suppose that (5) holds for some arbitrary $\delta \in [0, 1]$. We will argue that it then holds for $\delta/2$. Now, let

$$\delta/2 \text{ argmax}_{\pi \text{ safe}} P(s_h \in \text{RD}(D_n^{(B)}))$$

Also, let $\tilde{\pi}$ be some element of Π that is defined identically to π at every s_h such that $s_h \in \text{RD}(D_n^{(B)})$, and is defined arbitrarily otherwise. Also, let E denote the event that $s_h \in \text{RD}(D_n^{(B)})$ for every $h \in [n-1]$. Note that under every $\pi, \tilde{\pi}$ and $\tilde{\pi}$ take identical actions at the first $n-1$ time steps, so therefore the distributions of s_1, \dots, s_{n-1} conditional on E is the same under π and $\tilde{\pi}$. Therefore, we can bound

$$\begin{aligned} & P_{\tilde{\pi}}(s_h \in \text{RD}(D_n^{(B)})) \\ &= P_{\tilde{\pi}}(E) P_{\tilde{\pi}}(s_h \in \text{RD}(D_n) \mid E) \\ &\quad + P_{\tilde{\pi}}(\neg E) P_{\tilde{\pi}}(s_h \in \text{RD}(D_n) \mid \neg E) \\ &= P_{\tilde{\pi}}(E) P_{\tilde{\pi}}(s_h \in \text{RD}(D_n) \mid E) \\ &\quad + \frac{n-1}{H} P_{\tilde{\pi}}(s_h \in \text{RD}(D_n) \mid \neg E) \\ &= P_{\tilde{\pi}}(s_h \in \text{RD}(D_n) \mid E) + \frac{n-1}{H} \\ &= \frac{P_{\tilde{\pi}}(s_h \in \text{RD}(D_n))}{P_{\tilde{\pi}}(E)} + \frac{n-1}{H} \\ &= \frac{P_{\tilde{\pi}}(s_h \in \text{RD}(D_n))}{1 - \frac{n-1}{H}} + \frac{n-1}{H} \\ &= \frac{P_{\tilde{\pi}}(s_h \in \text{RD}(D_n))}{1 - \frac{n-1}{H}} + \frac{n-1}{H} \end{aligned}$$

Note that in these bounds we apply the inductive assumption, which gives us

$$\begin{aligned} P_{\tilde{\pi}}(s_h \in \text{RD}(D_n^{(B)})) &= P_{\tilde{\pi}}(s_h \in \text{RD}(D_n^{(B)})) \\ &\quad + \frac{n-1}{H} \delta/2 \\ &\geq \frac{n-1}{H} \delta/2 \end{aligned}$$

for every $\delta \in [0, 1]$. Also, we apply the assumption that $\frac{1}{2} \geq \frac{n-1}{H}$, which gives us $(1 - \frac{n-1}{H})^{-1} \leq 2$, the fact that $P_{\tilde{\pi}}(\neg E) = P_{\tilde{\pi}}(s_h \in \text{RD}(D_n) \mid \neg E)$, and the fact that $P(A \mid B) = P(A) = P(B)$ for generic events A and B .

Next, given the assumption that we satisfy the conditions of Lemma 3 with $\frac{1}{2}H^{-2}$, we have

$$\sup_{2 \leq n} \frac{1}{H} \sum_{h=1}^n \mathbb{P} (s_h \geq 2 \text{RD}(D_n^{(B)})) \leq \frac{1}{2H^2} ;$$

This then implies that

$$\sum_{h=1}^n \mathbb{P} (s_h \geq 2 \text{RD}(D_n^{(B)})) \leq \frac{n}{2H} ;$$

for every $2 \leq n$.

Then, noting that by construction $n \geq 2$, putting together the prior bounds gives us

$$\mathbb{P} (\exists h \geq 2 [n] : s_h \geq 2 \text{RD}(D_n^{(B)})) \leq \frac{1}{H} + \frac{n-1}{H} \leq \frac{n}{H} ;$$

Finally, noting the definition of τ , the above is equivalent to

$$\sup_{2 \leq n} \mathbb{P} (\exists h \geq 2 [n] : s_h \geq 2 \text{RD}(D_n^{(B)})) \leq \frac{n}{H} ;$$

which completes the inductive case. □

D.3 Proof of Lemma 5

Proof. Let D be an arbitrary label dataset, π be an arbitrary policy to run, and D^0 be the dataset that is obtained by running π for m episodes, and for all encountered states labeling and adding to D all $(s_h; a)$ such that $s_h \geq 2 \text{RD}^a(D)$. It is sufficient to argue that, for any choice of D and π such that

$$G(\pi; D) \geq \frac{1}{2} ; \tag{6}$$

we have with probability at least $\frac{1}{2}$ after performing the above process for m episodes that

$$G(\pi; D^0) \geq \frac{1}{2} G(\pi; D) ;$$

In the remainder of the proof, we will consider an arbitrarily chosen D that satisfy (6), and we will reason about how large m needs to be to ensure the above fact.

First, for every $h \in [H]$ and safety labeled dataset D , let

$$g_h(D) = \mathbb{P} (s_h \geq 2 \text{RD}(D)) ;$$

Note that according to this definition, we have $G(\pi; D) = \frac{1}{H} \sum_{h=1}^H g_h(D)$. Now, even though (6) ensures that we can bound $G(\pi; D)$ from below on average across $h \in [H]$, there might be some h for which $g_h(D)$ is arbitrarily small, or even zero. Therefore, it may be very inefficient or even infeasible to try to ensure that $G(\pi; D^0) \geq \frac{1}{2} G(\pi; D)$ simultaneously for every $h \in [H]$.

Instead, we will proceed by defining a target set $S = \{h \in [H] : g_h(D) \geq \frac{1}{8}\}$ of time indices for which we would like to

reduce $\mathbb{P}_h(D)$. Suppose we were able to ensure that $\mathbb{P}_h(D^0) \leq \frac{1}{4} \mathbb{P}_h(D)$ for all $h \geq H$. Then, we would have

$$\begin{aligned} \frac{1}{H} \sum_{h=1}^H \mathbb{P}_h(D^0) &\leq \frac{1}{4H} \sum_{h \geq 2H} \mathbb{P}_h(D) + \frac{1}{H} \sum_{h \geq H} \mathbb{P}_h(D) \\ &\leq \frac{1}{4H} \sum_{h=1}^H \mathbb{P}_h(D) + \frac{1}{8} \\ &\leq \frac{1}{4H} \sum_{h=1}^H \mathbb{P}_h(D) + \frac{1}{4H} \sum_{h=1}^H \mathbb{P}_h(D) \\ &= \frac{1}{2H} \sum_{h=1}^H \mathbb{P}_h(D); \end{aligned}$$

where in the final inequality we apply (6). Therefore, if we were able to prove that

$$\mathbb{P}_h(D^0) \leq \frac{1}{4} \mathbb{P}_h(D) \quad \text{w.p. at least } 1 - \frac{1}{8h^2} \quad (7)$$

then our required result would follow by a union bound over H . Therefore, in the remainder of the proof we will consider an arbitrary h such that $h \geq 2H$.

Next, let r_h be defined such that

$$r_h \mathbb{P}_h(r) = \frac{1}{4} \mathbb{P}_h(D):$$

Note that by Assumption 5, $\mathbb{P}_h(r)$ is continuous and non-decreasing with $\lim_{r \downarrow 0} \mathbb{P}_h(r) = 0$, and $\lim_{r \uparrow 1} \mathbb{P}_h(r) = 1$, so clearly such an r_h must exist. Furthermore, given the condition on $\mathbb{P}_h(D)$ we have $r_h \mathbb{P}_h(r) \leq \frac{1}{32}$, which allows us to bound

$$\mathbb{P}_h(r) \leq \min\left(\frac{1}{32}, 1\right):$$

Now, let $P_{h; \cdot}^U$ denote the distribution of s_h induced by \mathbb{P}_h , conditional on $s_h \in \text{RD}(D)$. Given the above definitions, for any given $f \in \mathcal{B}_{h; \cdot}(r_h)$, we can bound

$$\begin{aligned} P_{h; \cdot}^U(f) &= \frac{\mathbb{P}_h(f(s_h; a) \leq f^*(s_h; a))}{\mathbb{P}_h(D)} \\ &= \frac{r_h}{\mathbb{P}_h(D)} \\ &= \frac{1}{4} \mathbb{P}_h(r_h^0) \leq \frac{1}{4} \min\left(\frac{1}{32}, 1\right); \end{aligned} \quad (8)$$

Next, suppose we draw iid samples $s_h^{(1)}, \dots, s_h^{(k)}$ from $P_{h; \cdot}^U$, and let $F_h = \{f \in \mathcal{B}_{h; \cdot}(r_h) : f(s_h^{(i)}; a) \leq f^*(s_h^{(i)}; a) \text{ for some } i \in [k] \text{ and } a \in \mathcal{A}\}$. If this were the case, then by labeling the points we would eliminate all hypotheses outside $\mathcal{B}_{h; \cdot}(r_h)$. Therefore, letting D_k denote the dataset obtained by labeling all of these points and adding the rest, we would have

$$\begin{aligned} \mathbb{P}_h(D_k) &= \mathbb{P}(s_h \in \text{RD}(V(D_k))) \\ &= \mathbb{P}(s_h \in \text{RD}(B_h(r_h))) \\ &= \mathbb{P}_h(r_h) r_h \\ &= \frac{1}{4} \mathbb{P}_h(D); \end{aligned}$$

which is our required condition. Therefore, we can proceed by: (1) finding m such that the above occurs with probability at least $1 - \frac{1}{8m}$; and (2) finding k such that event $s_h \in \text{RD}(D)$ occurs at least k times with probability at least $1 - \frac{1}{8m}$. Then, applying a union bound we would be done. We will do these things one at a time.

Finding sufficiently large k

First, let us define the following stochastic process indexed by

$$G_k(f) = \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f) - E_{s_h}^U [c(s_h; f)];$$

where

$$c(s; f) = 1 - f \cdot a + 2A : f(s; a) = f \cdot ?(s; a)g;$$

Then, we have

$$\begin{aligned} \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f) &= \sup_{f \in \mathcal{F}_h} E_{s_h}^U [c(s_h; f)] + \sup_{f \in \mathcal{F}_h} G_k(f) \\ &\leq \frac{1}{4} \min(\epsilon, 32)^{-1} + \sup_{f \in \mathcal{F}_h} G_k(f); \end{aligned}$$

where in the second inequality follows by applying (8). Furthermore, we have

$$\sup_{f \in \mathcal{F}_h} G_k(f) = \sup_{f \in \mathcal{F}_h} G_k(f) - E[\sup_{f \in \mathcal{F}_h} G_k(f)] + E[\sup_{f \in \mathcal{F}_h} G_k(f)];$$

where $E[\sup_{f \in \mathcal{F}_h} G_k(f)]$ is the expected value of the random variable $\sup_{f \in \mathcal{F}_h} G_k(f)$ under the iid draws from $P_{s_h}^U$ (we will let this distribution be implicit in the rest of this sub-section of the proof). We will proceed by bounding these two terms one by one. Specifically, under any event where both terms are at most $\frac{1}{10} \min(\epsilon, 32)^{-1}$, we would have $\sup_{f \in \mathcal{F}_h} G_k(f) \leq \frac{1}{5} \min(\epsilon, 32)^{-1} < \frac{1}{4} \min(\epsilon, 32)^{-1}$. Therefore, we would have $\sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f) < 1$, which implies that for every $f \in \mathcal{F}_h$ we have $f(s_h^{(i)}; a) \leq f \cdot ?(s_h^{(i)}; a)$ for some $i \in [k]$ and $a \in \mathcal{A}$, which is our required condition.

For the first of these terms, let $s_h^{(1)}; \dots; s_h^{(k)}$ be an arbitrary sequence such that $s_h^{(i)} = s_h^{(j)}$ for all $i \in [j]$, for some $j \in [k]$. Then, we can bound

$$\begin{aligned} \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f) &= \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f) \\ &= \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f) - c(s_h^{(i)}; f) \\ &= \frac{1}{k} \sup_{f \in \mathcal{F}_h} c(s_h^{(j)}; f) - c(s_h^{(j)}; f) \\ &= \frac{1}{k}; \end{aligned}$$

Therefore, we can apply the bounded differences inequality to the first term, which gives us

$$\sup_{f \in \mathcal{F}_h} G_k(f) - E[\sup_{f \in \mathcal{F}_h} G_k(f)] \leq \frac{1}{10} \min(\epsilon, 32)^{-1};$$

with probability at least $\exp(-k \min(\epsilon, 32)^{-2}) = 50$. We can ensure that this event occurs with probability at least $1 - \epsilon(2H)$, by noting

$$\begin{aligned} \exp(-k \min(\epsilon, 32)^{-2}) &= 50 \\ \implies k &\geq 50 \min(\epsilon, 32)^2 \log(2H) : \end{aligned}$$

Next, for the second term above, we will follow a standard symmetrization argument. Let $\{s_h^{(i)}\}_{i=1}^k$ be iid Rademacher random variables (random variables such that $\text{Pr}(s_h^{(i)} = 1) = \text{Pr}(s_h^{(i)} = -1) = \frac{1}{2}$ for every $i \in [k]$). Also, let

$s_h^{(1)}, \dots, s_h^{(k)}$ be a collection shadow variables that are distributed independent from and identically to $s_h^{(1)}$; $s_h^{(k)}$. Then, we have

$$\begin{aligned}
 E[\sup_{f \in \mathcal{F}_h} G_k(f)] &= E_{s_h^{(1:k)}} \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f) \\
 &= E_{s_h^{(1:k)}} \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f) \\
 &= E_{s_h^{(1:k)}; s_h^{(1:k)}} \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f) \\
 &= E_{s_h^{(1:k)}; s_h^{(1:k)}} \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f) \\
 &= 2E_{s_h^{(1:k)}} \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f)
 \end{aligned}$$

where the first inequality follows from Jensen's, the subsequent equality follows by symmetry, and the final inequality follows since $\sup(a + b) \leq \sup a + \sup b$, and since by symmetry on the definitions of $s_h^{(i)}$, and $s_h^{(i)}$. Next, suppressing the subscript in the expectation for brevity, we note that

$$\begin{aligned}
 2E \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k c(s_h^{(i)}; f) &= 2E \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k f(s_h^{(i)}; a) \\
 &= 2E \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k f(s_h^{(i)}; a^{(i)}) \\
 &= E \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k f(s_h^{(i)}; a^{(i)}) \\
 &= E \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k f(s_h^{(i)}; a^{(i)})
 \end{aligned}$$

where the inequality step follows because

$$\sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k f(s_h^{(i)}; a) = \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k f(s_h^{(i)}; a^{(i)})$$

and also by applying Jensen's inequality, the second equality follows because $\sum_{i=1}^k a^{(i)} = 0$, and the final equality follows because $\sup f(s; a) = \sup f(s; a) - 1 = f(s; a) - 1$, and by symmetry, $f(s_h^{(i)}; a^{(i)})$ and $f(s_h^{(i)}; a^{(i)})$ are identically distributed for each $i \in [k]$.

Next, note that the final bound is the Rademacher complexity (under some distribution of state, action pairs). Since by assumption \mathcal{F} has VC dimension at most d_{VC} , so does \mathcal{F}_h , so by Rebeschini (2020, Theorem 5.6) we have

$$E \sup_{f \in \mathcal{F}_h} \frac{1}{k} \sum_{i=1}^k f(s_h^{(i)}; a) \leq C \sqrt{\frac{d_{VC}}{k}}$$

for some universal constant C , regardless of the distribution over state, action pairs. Therefore, we have

$$E \sup_{f \in \mathcal{F}_h} G_k(f) \leq C \sqrt{\frac{d_{VC}}{k}}$$

so we can ensure that $E[\sup_{f \in \mathcal{F}_h} G_k(f)] \leq \frac{1}{10} \min(\epsilon, 32)^{-1}$ as long as

$$k \geq 100C^2 \min(\epsilon, 32)^2 d_{VC}$$

Therefore, putting the above together, as long as

$$k \geq \max \left\{ 50 \min \left(\frac{1}{\epsilon} \right)^2 \log(2H), 100C^2 \min \left(\frac{1}{\epsilon} \right)^2 d_{VC} \right. \\ \left. 50 \min \left(\frac{1}{\epsilon} \right)^2 (2C^2 d_{VC} + \log(2H)) \right\};$$

then both bounds hold with probability at least $1 - \epsilon$, and therefore we have our required event that \mathcal{F}_h disagrees with π^* on at least one of the sampled states with this probability.

Finding sufficiently large m

Next, for our target value k of samples from P , the remaining question is how large m needs to be to ensure we have at least k samples of π^* where $\pi^* \in \mathcal{RD}(D)$, with probability at least $1 - \epsilon$? We know that for samples X_1, \dots, X_m , this event occurs with probability at least $1 - \epsilon$ by the construction of \mathcal{F}_h .

Now, suppose we set $k = (1 - \epsilon)^{-1} + t$ for some $t \geq 0$, and let X_1, \dots, X_m be a set of iid $\{0, 1\}$ -valued random variable such that $\text{Prob}(X_i = 1) = 1 - \epsilon$ for each $i \in [m]$. Also, let $X = \sum_{i=1}^m X_i$. Applying Bernstein's inequality, we have

$$\text{Prob}(X \leq k) = \text{Prob}(X \leq \mathbb{E}[X] - t) \leq \exp \left(- \frac{\frac{1}{2}t^2(1 - \epsilon)^2}{m(1 - \epsilon)(1 - \epsilon) + \frac{1}{3}t(1 - \epsilon)} \right);$$

Now, in order to ensure that this probability is smaller than ϵ , it is sufficient to solve

$$\exp \left(- \frac{\frac{1}{2}t^2(1 - \epsilon)^2}{m(1 - \epsilon)(1 - \epsilon) + \frac{1}{3}t(1 - \epsilon)} \right) = \epsilon \\ \Leftrightarrow \frac{1}{2}(1 - \epsilon)^2 t^2 - (4 - 3\epsilon)(1 - \epsilon) \log(2H) t - (1 - \epsilon) \log(2H) k = 0;$$

Now, letting t^+ be the greater solution of this quadratic equation,

$$t^+ = \frac{(4 - 3\epsilon)(1 - \epsilon) \log(2H)}{(1 - \epsilon)^2} + \sqrt{\frac{(4 - 3\epsilon)(1 - \epsilon) \log(2H)}{(1 - \epsilon)^2} + 2(1 - \epsilon) \log(2H) k} \\ = \frac{2(4 - 3\epsilon)(1 - \epsilon) \log(2H) + (1 - \epsilon) \sqrt{2(1 - \epsilon) \log(2H) k}}{3(1 - \epsilon)^2} + k;$$

where in the first inequality we apply the fact that $\frac{1}{a} + \frac{1}{b} \geq \frac{1}{a+b}$ for $a, b > 0$, and in the second inequality we note that $1 - \epsilon \geq \frac{1}{3}$ for any $\epsilon \in (0, 1)$, so $(4 - 3\epsilon)(1 - \epsilon) \geq 1$, and $2(1 - \epsilon) \log(2H) k \geq 0$.

Finally, we note that since the probability that $\pi^* \in \mathcal{RD}(D)$ is at least $1 - \epsilon$, the probability that we have fewer than k successes in m samples is no greater than $\text{Prob}(X \leq k)$, and therefore from the above bounds any choice of

$$m \geq (1 - \epsilon)^{-1} k + 3 \log(2H) \frac{k}{1 - \epsilon}$$

is sufficient to ensure that this probability is at most ϵ .

Putting Everything Together

Given the analysis, if we set

$$k = 50 \min \left(\frac{1}{\epsilon} \right)^2 (2C^2 d_{VC} + \log(2H)) \\ \text{and } m = \frac{k}{1 - \epsilon} + 3 \log(2H) \frac{k}{1 - \epsilon};$$

then our required result is ensured with probability at least .

We note that this setting satisfies the growth bound

$$m = \frac{1}{\epsilon} \min(\frac{1}{\epsilon}, 32)^2 d_{VC} \log(1/\epsilon) \log(H) ;$$

which is our required result. □

E DERIVATIONS OF POLICY COVER DIMENSION BOUNDS

Here we provide the derivations of our policy cover dimension bounds given in Section 4.1.

Tabular MDP Let some policy class Π be given, and assume the MDP has total states, and denote these by s_1, \dots, s_S . Then, for each $i \in [d]$, we define the policy $\pi_i \in \Pi$ according to

$$\pi_i = \operatorname{argmin}_{\pi \in \Pi} E_{\pi} \sum_{h=1}^H \mathbb{1}_{\{s_h = s_i\}} g ;$$

That is, π_i is the optimal policy in Π for the reward function given by reaching the i th state. Now, let some arbitrary subset $\mathcal{S} \subseteq S$, and arbitrary $\epsilon \in (0, 1]$ be given. Then, we have

$$\begin{aligned} \frac{1}{H} \sum_{h=1}^H P(s_h \in \mathcal{S}) &= \frac{1}{H} \sum_{h=1}^H \sum_{s \in \mathcal{S}} P(s_h = s) \\ &= \frac{1}{H} \sum_{h=1}^H \sum_{i=1}^S P(s_h = s_i) \\ &= \sum_{i=1}^S E_{\pi_i} \frac{1}{H} \sum_{h=1}^H \mathbb{1}_{\{s_h = s_i\}} g \\ &= \sum_{i=1}^S E_{\pi_i} \frac{1}{H} \sum_{h=1}^H \mathbb{1}_{\{s_h = s_i\}} g ; \end{aligned}$$

Since the above holds for arbitrary \mathcal{S} and ϵ , we have that $\{\pi_1, \dots, \pi_S\}$ is a policy cover for Π . That is, we always have a policy cover of size at most S , so $d_{\Pi} \leq S$.

Block MDP and Non-negative Rank MDP The reasoning here in both cases is similar to the tabular MDP case. In both cases, we can model the MDP as transitioning (from) to an intermediate discrete latent state following some distribution $P(j | s; a)$, and then sampling the next state following another distribution $P^0(j | z)$. Let d be the number of discrete latent states, which corresponds to $d_{\text{Block MDP}}$ and d_{NRR} for Non-negative Rank MDP in terms of our notation in Section 4.1. Also, for each $i \in [d]$, let z_i denote the discrete latent state that generates state i at time h , and let Z_1, \dots, Z_d denote the possible values of z . Then, following analogous reasoning as above, define

$$\pi_i = \operatorname{argmin}_{\pi \in \Pi} E_{\pi} \sum_{h=1}^H \mathbb{1}_{\{z_h = z_i\}} g ;$$

for some given policy class \mathcal{P} . Also, let some arbitrary measurable \mathcal{S} be given, along with some arbitrary policy $\pi \in \mathcal{P}$. Then, we have

$$\begin{aligned} \frac{1}{H} \sum_{h=1}^H P(s_h \in \mathcal{S}) &= \frac{1}{H} \sum_{h=1}^H \sum_{i=1}^d P(z_h = Z_i) P^0(\mathcal{S} | z = Z_i) \\ &= \sum_{i=1}^d P^0(\mathcal{S} | z = Z_i) E \left[\frac{1}{H} \sum_{h=1}^H 1_{z_h = Z_i} \right] \\ &= \sum_{i=1}^d P^0(\mathcal{S} | z = Z_i) E \left[\frac{1}{H} \sum_{h=1}^H 1_{z_h = Z_i} \right] \\ &= \frac{1}{H} \sum_{h=1}^H \sum_{i=1}^d P_i(z_h = Z_i) P^0(\mathcal{S} | z = Z_i) \\ &= \sum_{i=1}^d \frac{1}{H} \sum_{h=1}^H \sum_{j=1}^d P_i(z_h = Z_j) P^0(\mathcal{S} | z = Z_j) \\ &= \sum_{i=1}^d \frac{1}{H} \sum_{h=1}^H P_i(s_h \in \mathcal{S}); \end{aligned}$$

where the above bounds we apply the fact that the distribution for generating z does not depend on either the policy or h , and the second inequality follows because we are just introducing additional non-negative summands. That is, we always have a policy cover of size d so the policy cover dimension is at most d . So, for Block MDP we have $d = S$, and for Low Non-negative Rank MDP we have $d = d_{\text{NNR}}$.

F TABULAR MDP SAMPLE COMPLEXITY BOUND

The actual bound available in Azar et al. (2017) is a regret bound. Given N total episodes, as long as N is sufficiently large (specifically, if $N \geq S^3 A$), and $H \leq SA$ (which is reasonable in non-trivial settings) they bound the total regret by

$$\text{Regret}(N) \leq \sqrt{P H^3 S A N} \log(1/\delta)^2;$$

with probability at least $1 - \delta$. Now, suppose we run UCB-VI over N total episodes, for large N , and let $\hat{\pi}$ denote the average policy over all these episodes. Then, under the same high-probability event, the suboptimality of this average policy must be bounded by $\text{Regret}(\hat{\pi})/N$. That is,

$$\text{SubOpt}(\hat{\pi}; \pi^*) \leq \sqrt{\log(1/\delta)^2 \frac{H^3 S A}{N}};$$

Now, suppose we wish the sub-optimality to be bounded by some given $\epsilon \in (0, 1)$. Then solving for N , this corresponds to

$$\begin{aligned} \sqrt{\log(1/\delta)^2 \frac{H^3 S A}{N}} &\leq \epsilon \\ \implies N &\leq \frac{H^3 S A}{\epsilon^2} \log(1/\delta)^4; \end{aligned}$$

which is our presented sample complexity.

G EXPERIMENTAL DETAILS AND ADDITIONAL EXPERIMENTS

We show the promise of SABRE on a rich observation MDP setting called Block MDP Du et al. (2019); Misra et al. (2020). In this setting, the agent receives observations generated from a latent state. Further, no two states can generate the same observation. This is an expressive MDP setting which can capture non-trivial practical problems. Our goal is to discuss the implementation details for SABRE and show that it can achieve optimal return, while not taking unsafe actions, and minimizing the calls to safety oracle. Note that all details of our environment and methodology will also be available in our code release, which can be used to fully reproduce all of our results, which is currently withheld to preserve anonymity.

Environment. For a given horizon H , an instance of the environment has $4H + 1$ states and $K = 4$ actions. One of the action is a known safe action a_{safe} . These states are labeled as $S = \{s_{1,0}\} \cup \{s_{i,h} \mid i \in \{1,2,3,4\}, h \in [H]\}$. For a state $s_{i,h}$, the index i denotes its type and h denotes the level. We view states with index $i \in \{1,2\}$ as normal states, while states index $i = 3$ are considered safe states and those with $i = 4$ are unsafe states. The agent never observes the state but instead receives an observation $x = q(j|x)$ generated by an emission process q . The rich-observation setting implies that no two different states can generate the same observation.

The agent starts deterministically in $s_{1,0}$. All latent state transitions occur deterministically. After taking h actions, the agent is in one of these four states: $\{s_{j,h} \mid j \in \{1,2,3,4\}\}$. Each environment instance has two action sequences $(a_1; \dots; a_H)$ and $(a_1^o; \dots; a_H^o)$ with $a_h \in \{a_{\text{safe}}, a_h^o\}$ and $a_h^o \in \{a_{\text{safe}}, a_h^o\}$, for all $h \in [H]$. We view the a_h^o as the *unsafe action* and a_h as the *continue action* for time step h .

For any $i \in \{1,2\}$ and $h \in [H]$, transitions in the normal state $s_{i,h}$ operate as follows: taking the continue action a_h leads to the next normal state $s_{i,h+1}$ of same type, taking a_{safe} leads to the safe state $s_{3,h+1}$, taking the unsafe action a_h^o leads to the unsafe state $s_{4,h+1}$, while the remaining action leads to the normal state $s_{3-i,h}$ of the other type. All actions in the safe state $s_{3,h}$ leads to the next safe state $s_{3,h+1}$. Finally, in the unsafe state $s_{4,h}$, taking any action except a_{safe} leads to the next unsafe state $s_{4,h+1}$, while taking a_{safe} leads to the safe state $s_{3,h+1}$.

For a given transition $(s; a; s^o)$, the reward function $R(s; a; s^o)$ depends only on s^o . For $s^o = s_{i,h}$, the reward r is defined as follows: if $i = 1$ and $h < H$, then $r = 1-H$; for $i = 1$ and $h = H + 1$, $r = 2.0$; if $i = 2$ then $r = -1$, if $i = 3$ then $r = 0$ and if $i = 4$ then -1 . The optimal return of 2.8 is achieved by staying on states with index 1.

An observation is stochastically generated from a state $s_{i,h}$ each time the agent visits the state. This is done by first concatenating two 1-sparse vectors encoding i and h , and adding independent Gaussian noise to each dimension sampled from a 0 mean and 0.1 standard deviation. Let the resultant vector be z^o and it has $H + 5$ dimensions (4 dimensions to encode i and $H + 1$ dimensions to encode h). We concatenate z^o with a 0 vector of size $2^{\log_2 dH+5e} - (H + 5)$. Let Z be the final vector and its dimensionality is given by $k = 2^{\log_2 dH+5e}$. The observation is generated by multiplying Z with a Hadamard matrix of size $k \times k$ in order to mix its dimensions. We use Sylvester's construction to create a Hadamard matrix H_k of size $k \times k$. This is an inductive approach that works for $k = 2^l$ for some $l \in \mathbb{N}$. It defines $H_1 = [1]$ and $H_{2n} = \begin{bmatrix} H_n & H_n \\ H_n & -H_n \end{bmatrix}$ for every $n \in \mathbb{N}$. We borrowed the observation process from Misra et al. (2020).

In addition to the observation, the agent receives $(H + 1)B$ -dimensional safety features $f(s; a)g_{a2A}$ for every action upon visiting a state. We assume that the safety function $f^?$ is given by $f^?(s; a) = w^T(s; a) + b$ for some unknown $(w; b)$. In practice, these safety features can include readings from various safety devices attached to an agent (such as a LIDAR sensor, or temperature readings) which can be different from observation (e.g., camera image) which is used for dynamics. The agent receives a constant safe and unsafe feature for all states of type 3 and 4. States of type 2, have variance in safety features along all dimensions. Visiting these states can help us learn safety function faster, but this comes at the cost of a negative reward associated with visiting these states. Finally, a state $s_{1,h}$ only have variance in safety features along the $f(hB + 1; \dots; hB + 1)g$ dimensions. Visiting $s_{1,h}$ gives higher reward but reveals the safety features more slowly. Therefore, an optimal agent must balance between optimizing reward and minimizing safety oracle calls by exploring safety features quickly.

SABRE Implementation. We use Proximal Policy Optimization (Schulman et al., 2017) (PPO) as our blackbox RL algorithm. PPO is a popular empirical RL method and while not a PAC RL algorithm, is quite effective in practice for problems that do not require strategic exploration. We define the policy class and the value network in PPO using a two-layer feed forward neural network with Leaky ReLU non-linearity.

We implement queries to region of disagreement by reducing it to linear programming. Formally, given a safety labeled dataset $D = \{(s_i; a_i); y_i\}_{i=1}^n$, we solve for whether $(s; a)$ is in $\text{RD}^a(D)$, by first returning False if $a = a_{\text{safe}}$, and otherwise, solving the following two linear programs to compute the value of c_{max} and c_{min} :

$$c_{\text{max}} = \max_{w,b} w^T(s; a) + b; \quad \text{such that;} \\ \delta_i \in [d]; \quad y_i(w^T(s_i; a_i) + b_i) \leq 0; \\ kwk_1 \leq 1; \quad |b| \leq 1;$$

and

$$\begin{aligned}
 c_{\min} &= \min_{w,b} w^T (s; a) + b; \quad \text{such that;} \\
 &\quad \delta_i \geq [d]; \quad y_i(w^T (s_i; a_i) + b_i) \geq 0; \\
 &\quad kwk_1 \leq 1; \quad |b| \leq 1;
 \end{aligned}$$

The given feature $(s; a)$ is in $\text{RD}(D)$ if and only if either $c_{\max} \leq 0$ and $c_{\min} < 0$, or $c_{\max} > 0$ and $c_{\min} \geq 0$.

We implement $\text{RD}(D)$ by mapping each policy π to a *known safe policy* π_D and defining $\text{RD}(D) = \{ (s; a) \mid \pi_D(s; a) \geq g \}$. Let $(s; a) \in \text{SurelySafe}(D)$ denote set membership in the set of surely safe state-action pairs. We define π_D given as

$$\pi_D(a \mid s) = \frac{\mathbb{1}_{(s; a) \in \text{SurelySafe}(D)} g(s; a)}{\sum_{a' \in A} \mathbb{1}_{(s; a') \in \text{SurelySafe}(D)} g(s; a')};$$

When the agent visits a state s , we allow the agent to receive the safety features for all actions $f(s; a) \mathbb{1}_{(s; a) \in \text{SurelySafe}(D)}$. We test for set membership $(s; a) \in \text{SurelySafe}(D)$ by first returning True if $a = a_{\text{safe}}$, and otherwise, solving the linear programs described above using the feature $(s; a)$ and returning True if and only if $c_{\max} \leq 0$ and $c_{\min} \geq 0$.

Our implementation of SABRE closely follows Algorithm 1. One notable departure for efficiency is that we update the region of disagreement incrementally, rather than in batch.

Hyperparameters. We list hyperparameters that we use for SABRE for Table 2 and for the PPO baseline in Table 3. Note that we do *not* provide explicit values for the hyperparameters $\epsilon_{\text{explore}}$, R_{explore} , and R in SABRE. Instead, in practice these are implicitly determined by the other hyperparameters relating to the PPO blackbox algorithm. For example, the greater the total number of iterations of PPO is, the stronger guarantees we can make for the algorithm with respect to $\epsilon_{\text{explore}}$ and R_{explore} ; *i.e.* the greater the number of iterations, the smaller the implicit $\epsilon_{\text{explore}}$ and R_{explore} are. The hyperparameters N , B , and m in Table 2 correspond to those in Algorithm 1. The remaining hyperparameters in Table 2 are for the PPO blackbox algorithm, which are the same as the hyperparameters used by the unsafe PPO baseline that are listed in Table 3. Note that we use the same hyperparameter values for the PPO blackbox algorithm both when it is called on line 6 for safety learning, or line 9 for return optimization, in Algorithm 1. The hyperparameters used by PPO are the following:

- *Iterations of PPO* denotes the number of episodes used by each call to the PPO. Note that each time we call PPO in Algorithm 1 we start with a randomly initialized policy.
- *Gradient Clip* Is a maximum absolute value of all gradients for PPO updates; any gradients with absolute values above this threshold are clipped, which helps regularize and avoid extreme updates.
- *Number of PPO Updates* refers to the number of iterations of gradient descent optimization that we perform for each sampled batch. We use Adam optimization for training the policy parameters.
- *Ratio clipping coefficient* is a parameter used by PPO to control probability ratios, and limit how much the policy can change in each update.
- *Entropy coefficient* is the coefficient of entropy regularization for PPO, which encourages greater exploration and helps avoid convergence to a sub-optimal deterministic policy.

G.1 Additional Experiments

We also compare SABRE with a baseline SABRE-RO, which is based on the naive baseline algorithm described in Section 3.1. Like SABRE, SABRE-RO only ever takes actions that are known to be safe. In addition, to make the comparison fair, it labels data offline in batches at the same frequency as SABRE. However, instead explicitly exploring safety to directly learn F^* , it greedily tries to optimize return following the PPO algorithm (under the safety constraint). Specifically, each time it labels data, it sequentially labels $(s; a)$ for all previously observed s and $a \in A$ where $s \in \text{RD}^d(D)$, updating D after each call to the labeling oracle. We compare each method for various values of N and m .

Hyperparameter	Value
N	5
B	1
m	100
Iterations of PPO	1000
Learning Rate	0.001
Batch Size	32
Gradient Clip	20
Number of PPO Updates	10
Ratio Clipping Coefficient	0.1
Entropy Coefficient	0.01

Table 2: Hyperparameters for SABRE

Hyperparameter	Value
Iterations of PPO	6500
Learning Rate	0.001
Batch Size	32
Gradient Clip	20
Number of PPO Updates	10
Ratio Clipping Coefficient	0.1
Entropy Coefficient	0.01

Table 3: Hyperparameters for PPO

Results . We present mean return and cumulative oracle calls in Figure 3. We observed that SABRE-RO optimizes the environment reward which doesn't lead to the negative reward that SABRE accumulates early in the training. Further, if we consider the goal of achieving at least a return of 1.9 which is $1/2$ the optimal return of 2.8, then SABRE achieves this for $N = 2$ and $m = 100$, with least number of oracle calls of around 70. This is possible since SABRE can quickly learn the safety function by reaching the safe but low reward path, and then using it to optimize the environment reward. However, overall SABRE-RO is a strong baseline and almost as competitive as the SABRE. We leave a more detailed study and questions of lower bound for *passive safety* learning for future work.

