
Safe Exploration in Reinforcement Learning: A Generalized Formulation and Algorithms

Akifumi Wachi
LINE Corporation
akifumi.wachi@linecorp.com

Wataru Hashimoto
Osaka University
hashimoto@is.eei.eng.osaka-u.ac.jp

Xun Shen
Osaka University
shenxun@eei.eng.osaka-u.ac.jp

Kazumune Hashimoto
Osaka University
hashimoto@eei.eng.osaka-u.ac.jp

Abstract

Safe exploration is essential for the practical use of reinforcement learning (RL) in many real-world scenarios. In this paper, we present a generalized safe exploration (GSE) problem as a unified formulation of common safe exploration problems. We then propose a solution of the GSE problem in the form of a meta-algorithm for safe exploration, MASE, which combines an unconstrained RL algorithm with an uncertainty quantifier to guarantee safety in the current episode while properly penalizing unsafe explorations *before actual safety violation* to discourage them in future episodes. The advantage of MASE is that we can optimize a policy while guaranteeing with a high probability that no safety constraint will be violated under proper assumptions. Specifically, we present two variants of MASE with different constructions of the uncertainty quantifier: one based on generalized linear models with theoretical guarantees of safety and near-optimality, and another that combines a Gaussian process to ensure safety with a deep RL algorithm to maximize the reward. Finally, we demonstrate that our proposed algorithm achieves better performance than state-of-the-art algorithms on grid-world and Safety Gym benchmarks without violating any safety constraints, even during training.

1 Introduction

Safe reinforcement learning (RL) is a promising paradigm that enables policy optimizations for safety-critical decision-making problems (e.g., autonomous driving, healthcare, and robotics), where it is necessary to incorporate safety requirements to prevent RL policies from posing risks to humans or objects [14]. As a result, safe exploration has received significant attention in recent years as a crucial issue for ensuring the safety of RL during both the learning and execution phases [6].

Safe exploration in RL has typically been addressed by formulating a constrained RL problem in which the policy optimization is subject to safety constraints [9, 18]. While there have been many attempts under different types of constraint representations (e.g., expected cumulative cost [3], CVaR [29]), satisfying constraints almost surely or with high probability received less attention to date. Imagine safety-critical applications such as planetary exploration where even a single constraint violation may result in catastrophic failure. NASA’s engineers hope Mars rovers to ensure safety at least with high probability [8]; thus, constraint satisfaction “on average” does not fit their purpose.

While several algorithms have addressed this problem with this stricter notion of safety, there are several formulations in terms of how the constraints are represented, including cumulative [32], state [39], and instantaneous constraints [41], which respectively correspond to Problems 1, 2, and 3

as we will discuss shortly in Section 2. Unfortunately, there has been limited discussion on the relationships between these approaches, making it challenging for researchers to acquire a systematic understanding of the field as a whole. If a generalized problem were to be formulated, then the research community could pool their efforts to develop suitable algorithms.

A closer examination of existing algorithms that span the entire theory-to-practice spectrum reveals several areas for improvement. Practical algorithms using deep RL (e.g., [32],[39],[40]) may provide satisfactory performance after convergence, but do not usually guarantee safety during training. In contrast, theoretical studies (e.g., [4], [41]) that guarantee safety with high probability during training often have limitations, such as relying on strong assumptions (e.g., known state transition) or experiencing decreased performance in complex environments. In summary, many algorithms have been proposed in various safe RL formulations, but the creation of a safe exploration algorithm that is both practically useful and supported by theoretical foundations remains an open problem.

Contributions. We first present a generalized safe exploration (GSE) problem and prove its generality compared with existing safe exploration problems. By taking advantage of the tractable form of the safety constraint in the GSE problem, we establish a meta-algorithm for safe exploration, MASE. This algorithm employs an uncertainty quantifier for a high-probability guarantee that the safety constraints are not violated and penalizes the agent *before* safety violation, under the assumption that the agent has access to an “emergency stop” authority. Our MASE is both practically useful and theoretically well-founded, which allows us to optimize a policy via an arbitrary RL algorithm under the high-probability safety guarantee, even during training. We then provide two specific variants of MASE with different uncertainty quantifiers. One is based on generalized linear models (GLMs), for which we theoretically provide high-probability guarantees of safety and near-optimality. The other is more practical, combining a Gaussian process (GP, [27]) to ensure safety with a deep RL algorithm to maximize the reward. Finally, we show that MASE performs better than state-of-the-art algorithms on the grid-world and Safety Gym [28] without violating any safety constraints, even during training.

2 Preliminaries

Definitions. We consider an episodic safe RL problem in a constrained Markov decision process (CMDP, [3]), $M = \langle \mathcal{S}; \mathcal{A}; H; P; r; g; s_1 \rangle$, where \mathcal{S} is a state space, \mathcal{A} is an action space, $H \in \mathbb{Z}_{>0}$ is a (fixed) length of each episode, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0; 1]$ is a state transition probability, $r : \mathcal{S} \times \mathcal{A} \rightarrow [0; 1]$ is a reward function, $g : \mathcal{S} \times \mathcal{A} \rightarrow [0; 1]$ is a safety (cost) function, and $s_1 \in \mathcal{S}$ is an initial state. At each discrete time step, with a given (fully-observable) state s , the agent selects an action a with respect to its policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$, receiving the new state s' , reward r , and safety cost g . Though we assume a deterministic policy, our core ideas can be extended to stochastic policy settings. Given a policy π , the value and action-value functions in a state s at time h are respectively defined as

$$V_{r,h}(s) := \mathbb{E} \sum_{h^0=h}^H r(s_{h^0}; a_{h^0}) \mid s_h = s$$

and $Q_{r,h}(s; a) := \mathbb{E} \sum_{h^0=h}^H r(s_{h^0}; a_{h^0}) \mid s_h = s; a_h = a$, where the expectation \mathbb{E} is taken over the random state-action sequence $f(s_{h^0}; a_{h^0})_{h^0=h}^H$ induced by the policy π . Additionally, $\gamma \in (0; 1]$ is a discount factor for the reward function. In the remainder of this paper, we define $V_{\max} := \frac{1}{1-\gamma}$ and let $T_h : (\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}) \rightarrow (\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R})$ denote the Bellman update operator $T_h(Q)(s; a) := \mathbb{E} [r(s_h; a_h) + \gamma V_Q(s_{h+1}) \mid s_h = s; a_h = a]$, where $V_Q(s) := \max_{a \in \mathcal{A}} Q(s; a)$.

Three common safe RL problems. We tackle safe RL problems where constraints must be satisfied almost surely, even during training. While such problems have garnered attention in the research community, there are several types of formulations, and their relations are yet to be fully investigated.

One of the most popular formulations for safe RL problems involves maximizing $V_r := V_{r,1}(s_1)$ under the constraint that the cumulative cost is less than a threshold, which is described as follows:

Problem 1 (Almost surely safe RL with cumulative constraint [32]).

$$\max V_r \quad \text{subject to} \quad \Pr \sum_{h=1}^H g(s_h; a_h) \leq \gamma \quad P; \quad = 1;$$

where $\gamma \in \mathbb{R}_{>0}$ is a constant representing a threshold, and $g \in (0; 1]$ is a discount factor for g .

Observe that, the expectation is taken regarding the safety constraint in Problem 1. This problem was studied in [2], which is stricter than the conventional one where the expectation is taken with respect to the cumulative safety cost function ($E[\sum_{h=1}^H g(s_h; a_h)]$).

Another popular formulation involves leveraging the state constraints so that safety corresponds to avoiding visits to a set of unsafe states. This type of formulation has been widely adopted by previous studies on safe-critical robotics tasks [38–40, 45], which is written as follows:

Problem 2 (Safe RL with state constraints)

$$\max V_r \quad \text{subject to} \quad E \left[\sum_{h=1}^H I(s_h \in S_{\text{unsafe}}) \right] \leq \beta;$$

where $\beta \in [0, 1]$ is a threshold, $I(\cdot)$ is the indicator function, and $S_{\text{unsafe}} \subseteq S$ is a set of unsafe states.

Finally, some existing studies formulate safe RL problems via an instantaneous constraint, attempting to ensure safety even during the learning stage while aiming for extremely safety-critical applications such as planetary exploration [42] or healthcare [46]. Such studies typically require the agent to satisfy the following instantaneous safety constraint at every time step.

Problem 3 (Safe RL with instantaneous constraints)

$$\max V_r \quad \text{subject to} \quad \Pr \left[g(s_h; a_h) \geq \beta \right] = 0; \quad \forall h \in [1; H];$$

where $\beta \in [0, 1]$ is a time-invariant safety threshold.

3 Problem Formulation

This paper also requires an agent to optimize a policy under a safety constraint, as in the three common safe RL problems. We seek to find the optimal policy π^* of the following problem, which will hereinafter be referred to as the “generalized” safe exploration (GSE) problem:

Problem 4 (GSE problem) Let $\beta_h \in [0, 1]$ denote a time-varying threshold.

$$\max V_r \quad \text{subject to} \quad \Pr \left[g(s_h; a_h) \geq \beta_h \right] = 0; \quad \forall h \in [1; H];$$

This constraint is instantaneous, which requires the agent to learn a policy without a single constraint violation not only after convergence but also during training. We assume that the threshold is myopically known; that is, β_h is known at time h , but unknown before that. Crucially, at every time step h , since s_h is a fully observable state and the agent’s policy is deterministic, we will use a simplified inequality represented as $g(s_h; a_h) < \beta_h$ in the rest of this paper. This constraint is akin to that in Problem 3, with the difference that the safety threshold is time-varying.

Importance of the GSE problem. Though our problem may not seem relevant to Problems 1 and 2, we will shortly present and prove a theorem on the relationship between GSE problem and the three common safe RL problems.

Theorem 3.1. Problems 1, 2, and 3 can be transformed into GSE problem (i.e., Problem 4).

See Appendix B for the proof. In other words, the feasible policy space of GSE problem can be identical to those in the other three problems by properly defining the safety cost function and threshold β_h . Crucially, Problem 1 is a special case of GSE problem with $\beta_h = \beta$ for all h , where $\beta_{h+1} = \beta + \gamma (\beta - g(s_h; a_h))$ with $\beta_0 = \beta$. It is particularly beneficial to convert Problems 1 and 2, which have additive constraint structures, to GSE problem, which has an instantaneous constraint. The accurate estimation of the cumulative safety value in Problems 1 and 2 is difficult because they depend on the trajectories induced by a policy. Dealing with the instantaneous constraint in the GSE problem is easier, both theoretically and empirically. Also, especially when the environment is time-varying (e.g., there are moving obstacles), GSE problem is more useful than Problem 3.

Typical CMDP formulations with expected cumulative (safety) cost are out of the scope of the GSE problem. In such problems, the safety notion is milder; hence, although many advanced deep RL algorithms have been actively proposed that perform well in complex environments after convergence,

their performance in terms of safety during learning is usually low, as reported by Stooke et al. [30] or Wachi et al. [43]. Risk-constrained MDPs are also important safe RL problems that are covered by the GSE problem; they have been widely studied by representing risk as a constraint on some conditional value-at-risk [11] or using chance constraints [24, 26].

Dif culties and Assumptions. Theorem 3.1 insists that the GSE problem covers a wide range of safe RL formulations and is worth solving, but the problem is intractable without assumptions. We now discuss the dif culties in solving the GSE problem, and then list the assumptions in this paper.

The biggest dif culty with the GSE problem lies in the fact that there may be no viable safe action given the current state s_h , safety cost g , and threshold b_h . When $b_h = 0$ and $g(s_h; a) = 0$ for all $a \in A$, the agent has no viable action for ensuring safety. The agent needs to guarantee safety, even during training, where little environmental information is available; hence, it is significant for the agent to avoid such situations where there is no action that guarantees safety. Another dif culty is related to the regularity of the safety cost function and the strictness of the safety constraint. In this paper, the safety cost function is unknown a priori; thus, when the safety cost does not exhibit any regularity, the agent can neither infer the safety of decisions nor guarantee safety almost surely.

To address the first dif culty mentioned above, we use Assumptions 3.2 and 3.3.

Assumption 3.2 (Safety margin) There exists $\epsilon > 0$ such that $\Pr[g(s_h; a_h) \leq b_h - \epsilon] \geq 1 - \delta$ for all $s_h \in S$.

Assumption 3.3 (Emergency stop action) Let a be an emergency stop action such that $g(s; a) = 1$ for all $s \in S$. The agent is allowed to execute the emergency stop action and reset the environment if and only if the probability of guaranteed safety is not sufficiently high.

Assumption 3.2 is mild; this is similar to the Slater condition, which is widely adopted in the CMDP literature [13, 25]. We consider Assumption 3.3 is also natural for safety-critical applications because it is usually better to guarantee safety, even with human interventions, if the agent requires help in emergency cases. In some applications (e.g., the agent is in a hazardous environment), however, emergency stop actions should often be avoided because of the expensive cost of human intervention. In such cases, the agent needs to learn a reset policy allowing them to return to the initial state as in Eysenbach et al. [15], rather than asking for human help, which we will leave to future work.

As for the second dif culty, we assume that the safety cost function belongs to a class where uncertainty can be estimated and guarantee the satisfaction of the safety constraint with high probability. We present an assumption regarding an important notion called an uncertainty quantifier:

Assumption 3.4 (Uncertainty quantifier) Let $\hat{g} : S \times A \rightarrow \mathbb{R}$ denote the estimated mean function of safety. There exists an uncertainty quantifier $\eta : S \times A \rightarrow \mathbb{R}$ such that $g(s; a) \leq \hat{g}(s; a) + \eta(s; a)$ for all $(s; a) \in S \times A$, with a probability of at least $1 - \delta$.

4 Method

We propose MASE for the GSE problem, which combines an unconstrained RL algorithm with additional mechanisms for addressing the safety constraints. The pseudo-code is provided in Algorithm 1, and a conceptual illustration can be seen in Figure 1.

The most notable feature of MASE is that safety is guaranteed via the uncertainty quantifier and the emergency stop action (lines 9). The uncertainty quantifier is particularly useful because we can guarantee that the confidence bound contains the true safety cost function, that is, $g(s; a) \leq \hat{g}(s; a) + \eta(s; a)$ for all $s \in S$ and $a \in A$. This means that, if the agent chooses actions such that $\hat{g}(s_h; a_h) + \eta(s_h; a_h) \leq b_h$, then $g(s_h; a_h) \leq b_h$ holds with a probability of at least $1 - \delta$. Regarding the first dif culty mentioned in Section 3, it is crucial that there is at least one safe action. Thus, at every time step, the agent computes a set of actions that are considered to satisfy the safety constraints with a probability of at least $1 - \delta$ given the state s_h and threshold b_h . This is represented as

$$A_h^+ := \{a \in A \mid \hat{g}(s_h; a) + \eta(s_h; a) \leq b_h\}$$

Whenever the agent identifies that at least one action will guarantee safety, the agent is required to choose an action with A_h^+ (line 3). The emergency stop action is executed if and only if there is no

¹The solution in the GSE problem is guaranteed to be a conservative approximation of that in safe RL problems with chance constraints. For more details, see Appendix C.

Algorithm 1 Meta-Algorithm for Safe Exploration (MASE)

```

1: for episode  $t = 1; 2; \dots; T$  do
2:   for time  $h = 1; 2; \dots; H$  do
3:     Take "safe" action  $a_h = \arg \max_{a \in A_h^+} Q(s_h; a)$ . Execute only safe actions
4:     Receive reward  $r(s_h; a_h)$ , safety cost  $b(s_h; a_h)$ , and next state  $s_{h+1}$ 
5:     Update safety threshold  $\epsilon_{h+1}$ 
6:     if  $A_{h+1}^+ = \emptyset$ ; then
7:       Compute  $c(s_h; a_h) = \frac{c}{\min_{a \in A_{h+1}^+} Q(s_{h+1}; a)}$ . Penalty for the emergency stop action
8:       Append  $(s_h; a_h; b(s_h; a_h); s_{h+1})$  to  $D$ 
9:       break (i.e., take action  $a$ ). Execute the emergency stop action
10:    else
11:      Append  $(s_h; a_h; r(s_h; a_h); s_{h+1})$  to  $D$ 
12:    Optimize a policy  $\pi$  based on  $D$  via an (unconstrained) RL algorithm
13:    Update the uncertainty quantifier and rewrite  $D$ 

```

viable action satisfying the safety constraint ($A_{h+1}^+ \neq \emptyset$); that is, the agent is allowed to execute a and start a new episode from an initial safe state (s_{h+1}). The safety cost is upper-bounded by 1 because $b \in [0; 1]$ by definition. Note that MASE proactively avoids unsafe actions by selecting the emergency stop action to take beforehand; this is in contrast to Sun et al. [37], whose method terminates the episode immediately after the agent has already violated a safety constraint.

When safety is guaranteed in the manner described above, the question remains as to how to obtain a policy that maximizes the expected cumulative reward. As such, we first convert the original CMDP M to the following unconstrained MDP

$$\tilde{M} := (h; S; f; A; g; H; P; b; s_1; i)$$

The changes from M lie in the action space and the reward function, as well as in the absence of the safety cost function. First, the action space is augmented so that the agent can execute the emergency stop action a . The second modification concerns the reward function. When executing the emergency stop action a , the agent is penalized as its sacrifice so that the same situation will not occur in future episodes; hence, we modify the reward function as follows:

$$b(s_h; a_h) = \begin{cases} c = \min_{a \in A_{h+1}^+} Q(s_{h+1}; a) & \text{if } A_{h+1}^+ = \emptyset; \\ r(s_h; a_h) & \text{otherwise} \end{cases} \quad (1)$$

where $c \in \mathbb{R}_{>0}$ is a positive scalar representing a penalty for performing the emergency stop. This penalty is assigned to the state-action pair (s_h, a_h) that placed the agent into the undesirable situation at time step $h + 1$, represented as $A_{h+1}^+ = \emptyset$; (see Figure 1).

To show that MASE is a reasonable safe RL algorithm, we express the following intuitions. Consider the ideal situation in which the safety cost function is accurately estimated for any state-action pairs; that is, $b(s; a) = 0$ for all $(s; a)$. In this case, all emergency stop actions are properly executed, and the safety constraint will be violated at the next time step if the agent executes other actions. It is reasonable for the agent to receive a penalty $b(s; a) = 1$ because this state-action pair surely causes a safety violation without the emergency stop action. Unfortunately, however, the safety cost is uncertain and the agent conservatively executes other actions although there are still actions satisfying the safety constraint, especially in the early phase of training; hence, we increase or reduce the penalty according to the magnitude of uncertainty in (1) to avoid an excessively large penalty.

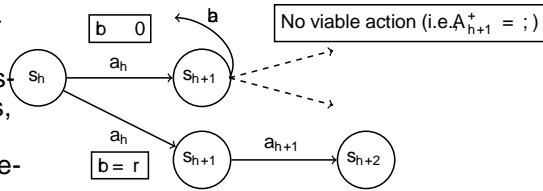


Figure 1: Conceptual illustration of MASE. At every time step, the agent chooses action a within A_h^+ . If there is no safe action at state s_{h+1} satisfying the safety constraint, the emergency stop action is executed and safety violation without the emergency stop action. The agent receives a large penalty $b(s_h; a_h)$.

Unfortunately, however, the safety cost is uncertain and the agent conservatively executes other actions although there are still actions satisfying the safety constraint, especially in the early phase of training; hence, we increase or reduce the penalty according to the magnitude of uncertainty in (1) to avoid an excessively large penalty.

We must carefully consider the fact that the quality of information regarding the modified reward function b is uneven in the replay buffer D . Specifically, in the early phase of training, the

uncertainty quantifier is loose. Hence, the emergency stop action is likely to be executed even if viable actions remain; that is, the agent will receive unnecessary penalties. In contrast, the emergency stop actions in later phases are executed with confidence, as reflected by the tight uncertainty quantifier. Thus, as in line 5, we rewrite the replay buffer while updating Q depending on the model in terms of the safety cost function (as for specific methods to update Q see Sections 5 and 6).

Connections to shielding methods. The notion of the emergency stop action is akin to shielding [20] which has been actively studied in various problem settings including partially-observable environments [0] or multi-agent settings [23]. Thus, MASE can be regarded as a variant of shielding methods (especially, preemptive shielding [21]) that is specialized for the GSE problem. On the other hand MASE does not only block unsafe actions but also provides proper penalties for executing the emergency stop actions based on the uncertainty quantifier, which leads to rigorous theoretical guarantees presented shortly. Such theoretical advantages can be enjoyed in many safe RL problems because of the wide applicability of the GSE problem backed by Theorem 3.1.

Advantages of MASE. Though certain existing algorithms for Problem 3 (i.e., the closest problem to the GSE problem) theoretically guarantee safety during learning, several strong assumptions are needed, such as a known and deterministic state transition and regular safety function ϕ and a known feature mapping function that is linear with respect to transition kernels, reward, and safety as in [5]. Such algorithms have little affinity with deep RL; thus, their actual performance in complex environments tends to be poor. In contrast, MASE is compatible with any advanced RL algorithms, which can also handle various constraint formulations while maintaining the safety guarantee.

Validity of MASE. We conclude this section by presenting the following two theorems to show that our MASE produces reasonable operations in solving GSE problem.

Theorem 4.1. Under Assumption 3.4, MASE guarantees safety with a probability of at least $1 - \epsilon$.

Theorem 4.2. Assume that the safety cost function is estimated for any state-action pairs with an accuracy of better than $\frac{\epsilon}{2}$; that is, $|\hat{c}(s; a) - c(s; a)| < \frac{\epsilon}{2}$ for all $(s; a)$. Set $\epsilon > 0$ to be a sufficiently large scalar such that $\epsilon > \frac{V_{\max}}{2\gamma}$. Then, the optimal policy π^* is identical to that in M .

See Appendix D for the proofs. Unfortunately, obtaining an uncertainty quantifier that works in general cases is highly challenging. To develop a feasible model for the uncertainty quantification, we assume that the safety cost can be modeled via a GLM in Section 5 and via a GP in Section 6.

5 A Provable Algorithm under Generalized Linear CMDP Assumptions

In this section, we focus on CMDPs with generalized linear structures and analyze the theoretical properties of MASE. Specifically, we provide a provable algorithm to use a class of GLMs denoted as \mathcal{F} for modeling $Q_{r,h}^* := Q_{r,h}^*$ and, then provide theoretical results on safety and optimality.

5.1 Generalized Linear CMDPs

We extend the assumption in Wang et al. [44] from unconstrained MDPs to CMDPs settings. Our assumption is based on GLMs as with [44] that makes a strictly weaker assumption than their Linear MDP assumption [19, 46]. As preliminaries, we first list the necessary definitions and assumptions.

Definition 5.1 (GLMs). Let $d \in \mathbb{Z}_{>0}$ be a feature dimension and $\mathcal{B}^d := \{x \in \mathbb{R}^d : \|x\|_2 \leq 1\}$ be the ℓ_2 ball in \mathbb{R}^d . For a known feature mapping function: $S \times A \rightarrow \mathcal{B}^d$ and a known link function $f : [-1, 1] \rightarrow [-1, 1]$, the class of generalized linear model is denoted as $\mathcal{F} := \{f(s; a; \theta) : \theta \in \mathcal{B}^d\}$ where $s; a := (s; a)$.

Assumption 5.2 (Regular link function) The link function $f(\cdot)$ is twice differentiable and is either monotonically increasing or decreasing. Furthermore, there exist absolute constants $\underline{c} < 1$ and $M < 1$ such that $\underline{c} < |f'(x)| < 1$ and $|f''(x)| < M$ for all $|x| \leq 1$.

This assumption on the regular link function is standard in previous studies (e.g., [21]). Linear and logistic models are the special cases of the GLM where the link functions are defined as $f(x) = x$ and $f(x) = \frac{1}{1 + e^{-x}}$. In both cases, the link functions satisfy Assumption 5.2.

We finally make the assumption of generalized linear CMDPs (GL-CMDPs), which extends the notion of the optimistic closure for unconstrained MDP settings in Wang et al. [44].

Assumption 5.3 (GL-CMDP). For any $1 \leq h < H$ and $u \in \mathcal{F}_{\text{up}}$, we have $T_h(u) \in \mathcal{F}$ and $g \in \mathcal{F}$.

Recall that T_h is the Bellman update operator. In Assumption 5.3, with a positive semi-definite matrix $A \in \mathbb{R}^{d \times d}$, $A \succeq 0$ and a fixed positive constant $\kappa_{\max} \in \mathbb{R}_{>0}$, we define

$$F_{\text{up}} := f(s; a) + \min_{i \in \mathcal{I}} V_{\max}(h; s; a; i) + \kappa_{s;a} k_A g : \mathcal{B}^d; 0 \quad \kappa_{\max}; \kappa_A \kappa_{\text{op}}^{-1} g;$$

where $\kappa_A := \frac{1}{\lambda_{\min}(A)}$ is the matrix Mahalanobis seminorm, and κ_{op} is the matrix operator norm. For simplicity, we suppose the same link functions for the Q-function and the safety cost function, but it is acceptable to use different link functions. Note that Assumption 5.3 is a more general assumption than Amani et al [5] that assumes linear transition kernel, reward, and safety cost functions or Wachi et al. [43] that assumes a known transition and GLMs in terms of reward and safety cost functions.

5.2 GLM-MASE Algorithm

We introduce an algorithm GLM-MASE under Assumptions 5.2 and 5.3. Hereinafter, we explicitly denote the episode for each variable. For example, $s_h^{(t)}$ or $a_h^{(t)}$ denote a state or action at the time step t of episode t . We also let $e_h^{(t)} := (s_h^{(t)}; a_h^{(t)})$ for more concise notations.

Uncertainty quantifiers. To actualize MASE in the generalized linear CMDP settings, we first need to consider how to obtain the uncertainty quantifier in terms of the safety cost function. Since we assume $g \in \mathcal{F}$, we can define the uncertainty quantifier based on the existing studies on GLMs, especially in the field of multi-armed bandits [22, 16]. Based on Assumptions 5.2 and 5.3, we now provide a lemma regarding the uncertainty quantifier on safety.

Lemma 5.4. Suppose Assumptions 5.2 and 5.3 hold. Set $\frac{1}{T_H}$. With a universal constant $C \in \mathbb{R}_{>0}$, let $C_g := C^{-1} \frac{1 + M + \kappa_{\max} + d^2 \ln \frac{1 + \kappa_{\max}}{1 - \kappa_{\max}}}{1}$. Define

$$Q_{g,h,t}(s; a) := C_g \kappa_{s;a} k_{h,t}^{-1} \quad \text{with} \quad h_{t,t} := \sum_{\tau=1}^t e_{h,\tau}^{(\cdot)} e_{h,\tau}^{(\cdot)} + I;$$

where I is the identity matrix. Let $b_{h,t}^g \in \mathbb{R}^d$ be the ridge estimate, which is computed by $b_{h,t}^g := \arg \min_{k \in \mathbb{R}^d} \sum_{\tau=1}^t g(s_{h,\tau}^{(\cdot)}; a_{h,\tau}^{(\cdot)}) - f(h e_{h,\tau}^{(\cdot)}; i)^2$. Then, the following inequality holds

$$|g(s_h^{(t)}; a_h^{(t)}) - f(h e_{h,t}^{(t)}; b_{h,t}^g)| \leq Q_{g,h,t}(s_h^{(t)}; a_h^{(t)})$$

for all $(s; a) \in \mathcal{S} \times \mathcal{A}$, with a probability at least $1 - \frac{\delta}{2}$.

For the purpose of qualifying uncertainty in GLMs, the weighted norm of $(i.e., \kappa_{s;a} k_{h,t}^{-1})$ plays an important role. Because we assume that the Q-function and safety cost function share the same feature, we have a similar lemma on the uncertainty quantifier regarding the Q-function as follows:

Lemma 5.5. Suppose Assumptions 5.2 and 5.3 hold. Let $b_{h,t}^Q \in \mathbb{R}^d$ denote the ridge estimate; that is, $b_{h,t}^Q := \arg \min_{k \in \mathbb{R}^d} \sum_{\tau=1}^t y_{h,\tau}^{(\cdot)} - f(h e_{h,\tau}^{(\cdot)}; i)^2$, where $y_{h,\tau}^{(\cdot)} := r(s_{h,\tau}^{(\cdot)}; a_{h,\tau}^{(\cdot)}) + \max_{a \in \mathcal{A}} Q_{r,h+1}^{(\cdot)}(s_{h+1}^{(\cdot)}; a^0)$ for all τ with

$$Q_{r,h}^{(t)}(s; a) := \min_{i \in \mathcal{I}} V_{\max}(h; s; a; i) + C_{Q_{\Rightarrow}}(s; a)$$

that is initialized with $b_{r,h}^{(0)} = 0$ for all $h \in H$ and $Q_{r,h+1}^{(t)} = 0$ for all $1 \leq t \leq T$. Then, with a universal constant $C_{Q_{\Rightarrow}} \in \mathbb{R}_{>0}$, the following inequalities holds

$$|Q_{r,h}^{(t)}(s_h^{(t)}; a_h^{(t)}) - f(h e_{h,t}^{(t)}; b_{h,t}^Q)| \leq C_{Q_{\Rightarrow}}(s_h^{(t)}; a_h^{(t)})$$

for all $(s; a) \in \mathcal{S} \times \mathcal{A}$, with a probability at least $1 - \frac{\delta}{2}$.

Note that $(s_h^{(t)}; a_h^{(t)})$ is the β -uncertainty quantifier with respect to the safety cost function. One of the biggest advantages of the generalized linear CMDPs is that the magnitude of uncertainty for the Q-function is proportional to that for the safety cost function. Hence, by exploring the Q-function

based on the optimism in the face of the uncertainty principle [16], the safety cost function is also explored simultaneously, which contributes to the efficient exploration of state-action spaces.

Integration into MASE. The GLM-MASE is an algorithm to integrate the uncertainty quantifiers inferred by the GLM into the MASE sequence. Detailed pseudo code is presented in Appendix E.

To deal with the safety constraint, GLM-MASE leverages the upper bound inferred by the GLM; that is, for all h and t , the agent takes only actions that satisfy

$$f_h(s_h^{(t)}; a_h^{(t)}) + b_{h,t} \leq b_h.$$

By Lemma 5.4, such state-action pairs satisfy the safety constraint $f_h(s_h^{(t)}; a_h^{(t)}) \leq b_h$, for all h and t , with a probability at least $1 - \delta$. If there is no action satisfying the safety constraint (i.e., $A_h^+ = \emptyset$), the emergency stop action is taken, and then the agent receives a penalty defined in (1).

As for policy optimization, we follow the optimism in the face of the uncertainty principle. Specifically, the policy π is optimized so that the upper-confidence bound of the Q-function characterized by b is maximized; that is, for any state s , the policy is computed as follows:

$$\pi_h^{(t)}(s) = \arg \max_{a \in A} Q_{bh}^{(t)}(s; a).$$

Intuitively, this equation enables us to 1) solve the exploration and exploitation dilemma by incorporating the optimistic estimates of the Q-function and 2) make the agent avoid generating trajectories to violate the safety constraint via the modified reward function.

Theoretical results. We now provide two theorems regarding safety and near-optimality. For both theorems, see Appendix E for the proofs.

Theorem 5.6. Suppose the assumptions in Lemma 5.4 hold. Then the MASE satisfies $g(s_h^{(t)}; a_h^{(t)}) \leq b_h$ for all $t \in [1; T]$ and $h \in [1; H]$, with a probability at least $1 - \delta$.

Theorem 5.7. Suppose the assumptions in Lemmas 5.4 and 5.5 hold. Let C_1 and C_2 be positive, universal constants. Also, with a sufficiently large t^* , let t^* denote the smallest integer satisfying $\lambda_{\min}(\Sigma) \geq C_1 t^* H d + C_2 t^* H \ln \frac{1}{\delta} + 2C_g t^*$, where $\lambda_{\min}(\cdot)$ is the minimum eigenvalue of the second moment matrix. Then, the policy $\pi^{(t)}$ obtained by GLM-MASE at episode t satisfies

$$\sum_{t=t^*}^T \sum_{h=1}^H \sum_{i=1}^I \mathbb{E} \left[\frac{1}{d^3(T-t)} \right]$$

with probability at least $1 - \delta$.

Theorem 5.6 shows that the GLM-MASE guarantees safety with high probability for every time step and episode, which is a variant of Theorem 4.1 under the generalized linear CMDP assumption and corresponding uncertainty quantifier. Theorem 5.7 demonstrates the agent's ability to act near-optimally after a sufficiently large number of episodes. The proof is based on the following idea. After t^* episodes, the safety cost function and the Q-function are estimated with an accuracy better than $\frac{\epsilon}{2}$. Then, based on Theorem 4.2, the optimal policy π^* is identical to that in M; thus, the agent achieves a near-optimal policy by leveraging the (well-estimated) optimistic Q-function.

6 A Practical Algorithm

Though we established an algorithm backed by theory under the generalized linear CMDP assumption in Section 5, it is often challenging to obtain proper feature mapping functions in complicated environments. Thus, in this section, we propose a more practical algorithm combining a GP-based estimator to guarantee safety with unconstrained deep RL algorithms to maximize the reward.

Guaranteeing safety via GPs. As shown in the previous sections, the uncertainty quantifier plays a critical role in MASE. To qualify the uncertainty in terms of the safety cost function, we consider modeling it as a GP $\mathbb{P}(z) = \mathcal{GP}(\mu(z); k(z; z^0))$, where $z := [s; a]$, $\mu(z)$ is a mean function, and $k(z; z^0)$ is a covariance function. The posterior distribution $\mathbb{P}(y|z)$ is computed based on $2 \times Z_{>0}$ observations at state-action pairs $(z_1; z_2; \dots; z_n)$ with safety measurements $y_n := f(y_1; y_2; \dots; y_n)$, where $y_n := g(z_n) + N_n$ and $N_n \sim \mathcal{N}(0; \sigma^2)$ is zero-mean Gaussian noise

with a standard deviation of $2R_0$. We consider episodic RL problems, and so $t \leq h$ for episode t and time step t , although the equality does not hold because of the episode cutoffs. Using the past measurements, the posterior mean, variance, and covariance are computed analytically as $\mu_n(z) = k_n^>(z)(K_n + I)^{-1}y_n$, $\Sigma_n(z; z) = k_n(z; z)$, and $k_n(z; z^0) = k(z; z^0) k_n^>(z)(K_n + I)^{-1}k_n(z^0)$, where $k_n(z) = [k(z_1; z); \dots; k(z_n; z)]^>$ and K_n is the positive definite kernel matrix. We now present a theorem on the safety guarantee.

Theorem 6.1. Assume $k_k^2 \leq B$ and $N_n \geq 1$ for all $n \geq 1$. Set $\beta_n^{1=2} := B + 4! \frac{P}{n + 1 + \ln(1 =)}$ and construct the uncertainty quantifier by

$$(\sigma; a) := \beta_n \Sigma_n(s; a); \quad \delta(s; a) \leq S \cdot A; \quad (2)$$

where β_n is the information capacity associated with kernel k . Then, MASE based on (2) satisfies the safety constraint $\mathbb{P}(s_h^{(t)}; a_h^{(t)}) \leq b_h$ for all t and h with a probability of at least $1 - \epsilon$.

See Appendix F for the proofs. Theorem 6.1 guarantees that the safety constraint is satisfied by combining the GP-based uncertainty quantifier in (2) and the emergency stop action.

Maximizing reward via deep RL. The remaining task is to optimize the policy via the modified reward function in (1), whereby the agent is penalized for emergency stop actions. This problem is decoupled from the safety constraint and can be solved as the following unconstrained RL problem:

$$\pi := \arg \max V_b; \quad (3)$$

There are many excellent algorithms for solving (3) such as trust region policy optimization (TRPO, [31]) and twin delayed deep deterministic policy gradient (TD3, [37]). One of the key benefits of our MASE is such compatibility with a broad range of unconstrained (deep) RL algorithms.

7 Experiments

We conduct two experiments. The first is on Safety Gym [26] where an agent must maximize the expected cumulative reward under a safety constraint with additive structures as in Problems 1 and 2. The safety cost function is binary (i.e., 1 for an unsafe state-action pair and 0 otherwise), and the safety threshold is set to $\tau = 20$. The reason for choosing Safety Gym is that this benchmark is complex and elaborate, and has been used to evaluate a variety of excellent algorithms. The second is a grid world where a safety constraint is instantaneous as in Problem 3. Due to the page limit, we present the settings and results of the grid-world experiment in Appendix H.

To solve the Safety Gym tasks, we implement the practical algorithm presented in Section 6 as follows. First, we convert the problem into a CSE problem by defining $\tilde{g}_h := g^{-1} \left(1 - \prod_{h=0}^{h-1} g(s_{h^0}; a_{h^0}) \right)$ and enforcing the safety constraint represented as $\tilde{g}_h(s; a) \leq b_h$ for every time step h in each episode. Second, to infer the safety cost, we use deep GP to conduct training and inference, especially when dealing with high-dimensional input spaces. Finally, as for the policy optimization, we leverage the TRPO algorithm. We delegate other details to Appendix G.

Baselines and metrics. We use the following four algorithms as baselines. The first is TRPO, which is a safety-agnostic deep RL algorithm that purely optimizes a policy without safety consideration. The second and third are CPO [1] and TRPO-Lagrangian [28], which are well-known algorithms for solving CMDPs. The final algorithm is Sauté RL [2], which is a recent, state-of-the-art algorithm for solving safe RL problems where constraints must be satisfied almost surely. We employ the following three metrics to evaluate MASE and the aforementioned four baselines: 1) the expected cumulative reward, 2) the expected cumulative safety, and 3) the maximum cumulative safety. We execute each algorithm with five random seeds and compute the means and confidence intervals.

Results. The experimental results are summarized in Figure 2. The figures show that TRPO, TRPO-Lagrangian, and CPO successfully learn the policies, but violate the safety constraints during training and even after convergence. Sauté RL is much safer than those three algorithms, but the safety constraint is not satisfied in some episodes, and the performance of the policy significantly deteriorates in terms of the cumulative reward during training. MASE obtains better policies in a smaller number of samples compared with Sauté RL, while also satisfying the safety constraints with respect to both the average and the worst-case. Note that, after convergence, the policy obtained by MASE performs worse than those obtained by the baseline algorithms in terms of reward, as shown in

(a) Average episode return. (b) Average episode safety. (c) Maximum episode safety.

(d) Average episode return. (e) Average episode safety. (f) Maximum episode safety.

Figure 2: Experimental results on Safety Gym (Top: PointGoal1, Bottom: CarGoal1). The proposed MASE satisfies the safety constraint in every episode and achieves better performance in terms of the reward than the state-of-the-art method called Sauté RL. Conventional methods (i.e., TRPO, TRPO-Lagrangian, and CPO) repeatedly violate the safety constraint, especially in the early phase of training. Shaded areas represent confidence intervals across 100 different random seeds.

(a) Average episode return. (b) Average episode safety. (c) Maximum episode safety.

Figure 3: Experimental results on Safety Gym (CarGoal1) with the emergency stop action. The box plots show the converged performance. Though MASE performs worse than other baselines in terms of reward, the acquired policy is still near-optimal. As for safety, while baselines violate the safety constraint in most of the episodes, MASE guarantees the satisfaction of the severe safety constraint.

Figure 3. The emergency stop action is a variant of resetting actions that are common in episodic RL settings, which prevent the agent from exploring the state-action spaces since the uncertainty quantifier is sometimes quite conservative. We consider that this is a reason why the converged reward performance of MASE is worse than other methods. However, because we require the agent to solve difficult problems where safety is guaranteed at every time step and episode, we consider that this result is reasonable, and further performance improvements are left to future work.

8 Conclusion

In this article, we first introduced the MASE problem and proved that it is more general than three common safe RL problems. We then proposed MASE to optimize a policy under safety constraints that allow the agent to execute an emergency stop action at the sacrifice of a penalty based on the -uncertainty quantifier. As a specific instance of MASE, we first presented GLM-MASE to theoretically guarantee the near-optimality and safety of the acquired policy under generalized linear CMDP assumptions. Finally, we provided a practical MASE and empirically evaluated its performance in comparison with several baselines on the Safety Gym and grid-world.

Acknowledgments and Disclosure of Funding

We would like to thank the anonymous reviewers for their helpful comments. This work is partially supported by JST CREST JPMJCR201 and by JSPS KAKENHI Grant 21K14184.

References

- [1] Achiam, J., Held, D., Tamar, A., and Abbeel, P. (2017). Constrained policy optimization. In International Conference on Machine Learning (ICML)
- [2] Alshiekh, M., Bloem, R., Ehlers, R., Könighofer, B., Niekum, S., and Topcu, U. (2018). Safe reinforcement learning via shielding. AAAI Conference on Artificial Intelligence (AAAI)
- [3] Altman, E. (1999). Constrained Markov decision processes, volume 7. CRC Press.
- [4] Amani, S., Alizadeh, M., and Thrampoulidis, C. (2019). Linear stochastic bandits under safety constraints. In Neural Information Processing Systems (NeurIPS)
- [5] Amani, S., Thrampoulidis, C., and Yang, L. (2021). Safe reinforcement learning with linear function approximation. International Conference on Machine Learning (ICML)
- [6] Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., and Mané, D. (2016). Concrete problems in AI safety. arXiv preprint arXiv:1606.06565
- [7] Auer, P. and Ortner, R. (2007). Logarithmic online regret bounds for undiscounted reinforcement learning. In Neural Information Processing Systems (NeurIPS)
- [8] Bajracharya, M., Maimone, M. W., and Helmick, D. (2008). Autonomy for mars rovers: Past, present, and future. Computer 41(12):44–50.
- [9] Brunke, L., Greeff, M., Hall, A. W., Yuan, Z., Zhou, S., Panerati, J., and Schoellig, A. P. (2022). Safe learning in robotics: From learning-based control to safe reinforcement learning. Review of Control, Robotics, and Autonomous Systems 5(4):411–444.
- [10] Carr, S., Jansen, N., Junges, S., and Topcu, U. (2023). Safe reinforcement learning via shielding under partial observability. AAAI Conference on Artificial Intelligence
- [11] Chow, Y., Ghavamzadeh, M., Janson, L., and Pavone, M. (2017). Risk-constrained reinforcement learning with percentile risk criteria. Journal of Machine Learning Research (JMLR) 18(1):6070–6120.
- [12] Chowdhury, S. R. and Gopalan, A. (2017). On kernelized multi-armed bandits. International Conference on Machine Learning (ICML)
- [13] Ding, D., Zhang, K., Basar, T., and Jovanovic, M. (2020). Natural policy gradient primal-dual method for constrained Markov decision processes. Neural Information Processing Systems (NeurIPS)
- [14] Dulac-Arnold, G., Levine, N., Mankowitz, D. J., Li, J., Paduraru, C., Gowal, S., and Hester, T. (2021). Challenges of real-world reinforcement learning: definitions, benchmarks and analysis. Machine Learning
- [15] Eysenbach, B., Gu, S., Ibarz, J., and Levine, S. (2018). Leave no trace: Learning to reset for safe and autonomous reinforcement learning. International Conference on Learning Representations (ICLR).
- [16] Filippi, S., Cappe, O., Garivier, A., and Szepesvári, C. (2010). Parametric bandits: The generalized linear case. Neural Information Processing Systems (NeurIPS)
- [17] Fujimoto, S., van Hoof, H., and Meger, D. (2018). Addressing function approximation error in actor-critic methods. International Conference on Machine Learning (ICML)
- [18] García, J. and Fernández, F. (2015). A comprehensive survey on safe reinforcement learning. Journal of Machine Learning Research (JMLR) 16(1):1437–1480.

- [19] Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. (2020). Provably efficient reinforcement learning with linear function approximation. *Conference on Learning Theory (COLT)*
- [20] Könighofer, B., Bloem, R., Junges, S., Jansen, N., and Serban, A. (2020). Safe reinforcement learning using probabilistic shields. *International Conference on Concurrency Theory: CONCUR*
- [21] Li, L., Chu, W., Langford, J., and Schapire, R. E. (2010). A contextual-bandit approach to personalized news article recommendation. *International Conference on World Wide Web (WWW)*
- [22] Li, L., Lu, Y., and Zhou, D. (2017). Provably optimal algorithms for generalized linear contextual bandits. *International Conference on Machine Learning (ICML)*, pages 2071–2080.
- [23] Melcer, D., Amato, C., and Tripakis, S. (2022). Shield decentralization for safe multi-agent reinforcement learning. *Neural Information Processing Systems (NeurIPS)*
- [24] Ono, M., Pavone, M., Kuwata, Y., and Balam, J. (2015). Chance-constrained dynamic programming with application to risk-aware robotic space exploration. *Autonomous Robots* 39(4):555–571.
- [25] Paternain, S., Chamon, L., Calvo-Fullana, M., and Ribeiro, A. (2019). Constrained reinforcement learning has zero duality gap. *Neural Information Processing Systems (NeurIPS)*
- [26] Pfrommer, S., Gautam, T., Zhou, A., and Sojoudi, S. (2022). Safe reinforcement learning with chance-constrained model predictive control. *Learning for Dynamics and Control Conference (L4DC)*, pages 291–303.
- [27] Rasmussen, C. E. (2003). Gaussian processes in machine learning. *Summer school on machine learning*, pages 63–71. Springer.
- [28] Ray, A., Achiam, J., and Amodei, D. (2019). Benchmarking safe exploration in deep reinforcement learning. *OpenAI*.
- [29] Rockafellar, R. T., Uryasev, S., et al. (2000). Optimization of conditional value-at-risk. *Journal of Risk*, 2:21–42.
- [30] Salimbeni, H. and Deisenroth, M. (2017). Doubly stochastic variational inference for deep Gaussian processes. *Neural Information Processing Systems (NeurIPS)*
- [31] Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. (2015). Trust region policy optimization. *International Conference on Machine Learning (ICML)*
- [32] Sootla, A., Cowen-Rivers, A. I., Jafferjee, T., Wang, Z., Mguni, D. H., Wang, J., and Ammar, H. (2022). Sauté RL: Almost surely safe reinforcement learning using state augmentation. *International Conference on Machine Learning (ICML)*
- [33] Stooke, A., Achiam, J., and Abbeel, P. (2020). Responsive safety in reinforcement learning by pid Lagrangian methods. *International Conference on Machine Learning (ICML)*, pages 9133–9143.
- [34] Strehl, A. L. and Littman, M. L. (2005). A theoretical analysis of model-based interval estimation. *International Conference on Machine Learning (ICML)*
- [35] Strehl, A. L. and Littman, M. L. (2008). An analysis of model-based interval estimation for Markov decision processes. *Journal of Computer and System Sciences* 74(8):1309–1331.
- [36] Sui, Y., Gotovos, A., Burdick, J. W., and Krause, A. (2015). Safe exploration for optimization with Gaussian processes. *International Conference on Machine Learning (ICML)*
- [37] Sun, H., Xu, Z., Fang, M., Peng, Z., Guo, J., Dai, B., and Zhou, B. (2021). Safe exploration by solving early terminated MDP. *ParXiv preprint arXiv:2107.04200*

- [38] Thananjeyan, B., Balakrishna, A., Nair, S., Luo, M., Srinivasan, K., Hwang, M., Gonzalez, J. E., Ibarz, J., Finn, C., and Goldberg, K. (2021). Recovery RL: Safe reinforcement learning with learned recovery zone. *IEEE Robotics and Automation Letters* 6(3):4915–4922.
- [39] Thomas, G., Luo, Y., and Ma, T. (2021). Safe reinforcement learning by imagining the near future. In *Neural Information Processing Systems (NeurIPS)*
- [40] Turchetta, M., Kolobov, A., Shah, S., Krause, A., and Agarwal, A. (2020). Safe reinforcement learning via curriculum induction. *Neural Information Processing Systems (NeurIPS)* volume 33, pages 12151–12162.
- [41] Wachi, A. and Sui, Y. (2020). Safe reinforcement learning in constrained Markov decision processes. *International Conference on Machine Learning (ICML)*
- [42] Wachi, A., Sui, Y., Yue, Y., and Ono, M. (2018). Safe exploration and optimization of constrained MDPs using Gaussian processes. *AAAI Conference on Artificial Intelligence (AAAI)*
- [43] Wachi, A., Wei, Y., and Sui, Y. (2021). Safe policy optimization with local generalized linear function approximations. In *Neural Information Processing Systems (NeurIPS)*
- [44] Wang, Y., Wang, R., Du, S. S., and Krishnamurthy, A. (2020). Optimism in reinforcement learning with generalized linear function approximation. *International Conference on Learning Representations (ICLR)*
- [45] Wang, Y., Zhan, S. S., Jiao, R., Wang, Z., Jin, W., Yang, Z., Wang, Z., Huang, C., and Zhu, Q. (2023). Enforcing hard constraints with soft barriers: Safe reinforcement learning in unknown stochastic environments. *International Conference on Machine Learning* pages 36593–36604. PMLR.
- [46] Yang, L. and Wang, M. (2019). Sample-optimal parametric Q-learning using linearly additive features. In *International Conference on Machine Learning (ICML)*

Appendices

A Limitations and Potential Negative Societal Impacts

We will first discuss limitations and potential negative societal impacts regarding our work.

A.1 Limitations

One of the major limitations of this study is that emergency stop actions are allowed for agents. Emergency stop actions should often be avoided because of the expensive cost of human intervention in many applications (e.g., the agent is in a hazardous or remote environment). In future work, we will investigate an algorithm that requires the agent to learn a reset policy allowing them to return to the initial state as in [15], rather than asking for human intervention via emergency stop actions.

Another limitation is how to construct the uncertainty quantifier. In our experiment, because we used a computationally inexpensive deep GP algorithm [30] and the uncertainty quantifier is updated at the end of the episode (see Line 13 in Algorithm 1), the computational time of the GP part was much smaller than the RL part in our experiment settings. However, since GP is generally a computationally expensive algorithm, GP can be a computational bottleneck in some cases.

A.2 Potential Negative Societal Impacts

We believe that safety is an essential requirement for applying RL in many real problems. While we have not found any potential negative societal impact of our proposed method, we must remain aware that RL algorithms are vulnerable to misuse and ours is no exception.

B Proof of Theorem 3.1

We first present lemmas regarding the relationship between the GSE problem and Problems 1, 2, and 3. After that, we present the proof for the Theorem 3.1 in Appendix B.4 by combining them.

B.1 Relationship between the GSE problem and Problem 1

Lemma B.1. Problem 1 can be transformed into the GSE problem.

Proof. We first utilize a safety state augmentation technique presented in Sootla [32] by defining a new variable h such that

$$h := \gamma^h \prod_{h^0=1}^h g(s_{h^0}; a_{h^0}) \quad ; \quad \forall h \in [1; H] \quad (4)$$

This new variable h means the remaining safety budget associated with the discount factor γ which is updated as follows:

$$h_{t+1} = \gamma^{-1} (h_t g(s_h; a_h)) \quad \text{with} \quad h_0 = 1 \quad (5)$$

By (4), the necessary and sufficient condition for satisfying the constraint in Problem 1 is

$$h_t \geq 0; \quad \forall h \in [1; H] \quad (6)$$

By (5), we have

$$h_{t+1} \geq 0; \quad \forall h \in [1; H] \iff h_t g(s_h; a_h) \geq 0; \quad \forall h \in [1; H] \quad (7)$$

In summary, by introducing the new variable h , Problem 1 is rewritten to the following problem:

$$\max V_r \quad \text{subject to} \quad \Pr[g(s_h; a_h) \geq h \mid P;] = 1; \quad \forall h \in [1; H] \quad (8)$$

The aforementioned problem (8) is a special case of the GSE problem with $h_t := h$. Therefore, we obtain the desired lemma. \square

B.2 Relationship between the GSE problem and Problem 2

Lemma B.2. Problem 2 can be transformed into the GSE problem.

Proof. The following chain of inequalities holds:

$$\begin{aligned}
 & \mathbb{E} \left[\sum_{h=1}^H g^h(s_h, a_h) \mid P; \gamma \right] \\
 \leq & \mathbb{E} \left[\sum_{h=0}^{H-1} g^h(s_h, a_h) \mid P; \gamma \right] + \gamma \sum_{h=1}^H \mathbb{E} \left[\sum_{i=1}^h g^i(s_i, a_i) \mid P; \gamma \right] \\
 \leq & \mathbb{E} \left[\sum_{h=0}^{H-1} g^h(s_h, a_h) \mid P; \gamma \right] + \gamma \sum_{h=1}^H \mathbb{E} \left[\sum_{i=1}^h g^i(s_i, a_i) \mid P; \gamma \right] \\
 \leq & \mathbb{E} \left[\sum_{h=0}^{H-1} g^h(s_h, a_h) \mid P; \gamma \right] + \gamma \sum_{h=1}^H \mathbb{E} \left[\sum_{i=1}^h g^i(s_i, a_i) \mid P; \gamma \right]
 \end{aligned}$$

Because we assume Markov property, we simply define the safety cost function A as

$$g(s_h; a) := \mathbb{E}[I(s_h \in S_{\text{unsafe}}) \mid P; \gamma] + \gamma \sum_{i=1}^h A(s_i, a_i)$$

We now set

$$b_h := \mathbb{E} \left[\sum_{i=1}^h g^i(s_i, a_i) \mid P; \gamma \right]$$

then the Problem 2 can be transformed into

$$\Pr[g(s; a) \leq b_h \mid P; \gamma] = 1:$$

Finally, we obtained the desired lemma. □

B.3 Relationship between the GSE problem and Problem 3

Lemma B.3. Problem 3 can be transformed into the GSE problem.

Proof. Set $b_h = \gamma$ for all h . Then, the GSE problem is identical to Problem 3. □

B.4 Summary

Proof. Combining Lemma B.1, B.2, and B.3, we obtain the desired Theorem 3.1. □

C Connections to Safe RL Problems with Chance Constraints

As a strongly related formulation to Problem 2, policy optimization under joint chance-constraints has been studied especially in the field of control theory such as Ono et al. [24] and Pfrommer et al. [26], which is written as

Problem 5 (Safe RL with joint chance constraints) Let $\gamma \in \mathbb{R}_0$ be a constant representing a safety threshold. Also, let $S_{\text{unsafe}} \subseteq S$ denote a set of unsafe states. Find the optimal policy such that

$$\max V_r \quad \text{subject to} \quad \Pr \left[\sum_{h=1}^H I(s_h \in S_{\text{unsafe}}) \leq \gamma \right] = 1:$$

Lemma C.1. Problem 2 is a conservative approximation of Problem 5.

Proof. This lemma mostly follows from Theorem 1 in Ono et al. [24]. Regarding the constraint in Problem 5, we have the following chain of equations:

$$\Pr \left[\bigcap_{h=1}^H s_h \geq S_{\text{unsafe}} \mid P; \right] = \Pr [s_1 \geq S_{\text{unsafe}} \mid P;] \quad (9)$$

$$= E \left[I(s_1 \geq S_{\text{unsafe}}) \mid P; \right] = E \left[I(s_1 \geq S_{\text{unsafe}}) \mid P; \right] \quad (10)$$

In the first step, we used Boole's inequality (i.e. $\Pr[A \cap B] \geq \Pr[A] + \Pr[B] - 1$). The final term in (10) is the LHS of the constraint in Problem 2, which implies that Problem 2 is a conservative approximation of Problem 5. Therefore, the GSE problem is also a conservative approximation of Problem 5. \square

Corollary C.2. Suppose we solve the GSE problem by properly defining the safety function $g_h(\cdot)$ and the threshold b_h . Then, the obtained policy is a conservative solution of Problem 5.

Remark C.3. It is extremely challenging to directly solve the Problem 5 characterized by joint chance-constraints without approximation, as discussed in Ono et al. [24]. Practically, it would be a promising and realistic approach to solve Problem 5 by converting it into a GSE problem.

More detailed explanations for the aforementioned remark are as follows. It is extremely challenging to directly solve the Problem 5 characterized by joint chance-constraints. Most of the previous work does not directly deal with this type of constraint and uses some approximations or assumptions. For example, Pfrommer et al. [26] assume a known linear time-invariant dynamics. Also, Ono et al. [24] approximate the joint chance constraint as in the above procedure and obtain

$$\Pr \left[\bigcap_{h=1}^H s_h \geq S_{\text{unsafe}} \mid P; \right] \leq E \left[\sum_{h=1}^H I(s_h \geq S_{\text{unsafe}}) \mid P; \right]$$

This is a conservative approximation with an additive structure, which is easier to solve than the original joint chance constraint. Ono et al. [24] deals with the above constraints with additive safety structure. By additionally transforming the conservatively-approximated problem into the GSE problem, the problem would become easier to handle because the safety constraint is instantaneous.

D Proof of Theorems 4.1 and 4.2

D.1 Proof of Theorem 4.1

Proof. Recall that, at every time step t , the MASE chooses actions satisfying

$$g(s_h; a_h) + (s_h; a_h) \geq b_h \quad (11)$$

By definition of the ϵ -uncertainty quantifier, the actions that are conservatively chosen based on \hat{g} also satisfy the safety constraint, with a probability of at least $1 - \epsilon$; that is,

$$g(s_h; a_h) \geq b_h \quad (12)$$

In addition, when there is no action satisfying (11), the emergency stop action is executed; that is, safety is guaranteed with a probability of $1 - \epsilon$. Hence, the desired theorem is now obtained. \square

D.2 Proof of Theorem 4.2

Proof. Assumption 3.2 implies that there exists a policy that satisfies a more conservative safety constraint written as

$$g(s_h; a_h) \geq b_h \quad ; \quad \forall s_h \in [1; H] \quad (13)$$

By combining (13) and the assumption $(s; a) \geq \frac{1}{2} ; (s; a) \geq S_{\text{A}}$, we have

$$(s_h; a_h) + (s_h; a_h) \geq g(s_h; a_h) + b_h \quad ; \quad \forall s_h \in [1; H] \quad (14)$$

Algorithm 2 GLM-MASE

```

1: for episode = 1; 2; ...; T do
2:   for time h = 1; ...; H do
3:     Take action  $a_h = \pi(s_h)$  within  $A_h^+$ 
4:     Receive reward  $r(s_h; a_h)$  and next state  $s_{h+1}$ 
5:     Receive safety cost  $c(s_h; a_h)$  and update safety threshold  $d_{h+1}$ 
6:     if  $A_h^+ = \emptyset$ ; then
7:       Compute  $b(s_h; a_h) = \frac{c}{\min_{a \in A_h} (s_{h+1}; a)}$ 
8:       Append  $(s_h; a_h; b(s_h; a_h); s_{h+1})$  to  $D$ 
9:       break (i.e., emergency stop action)
10:    else
11:      Append  $(s_h; a_h; r(s_h; a_h); s_{h+1})$  to  $D$ 
12:    Update  $Q_{h;t}^g$  and  $Q_{h;t}^Q$ 
13:    Compute the optimistic Q-estimate

```

$$Q_{bh}^{(t)}(s; a) = \min_{i \in \{g, Q\}} \{ V_{\max}; f(s; a; Q_{h;t}^i) \} + C_{\alpha} \gamma (s; a) g$$

```

14:   Optimize the policy by

```

$$\pi_h^{(t)}(s) = \arg \max_{a \in A_h} Q_{bh}^{(t)}(s; a)$$

```

15:   Update  $(s; a) := C_g \cdot k_{s;a} \cdot k_{h;t}^{-1}$  and then rewrite  $D$ 

```

which guarantees that there exists a policy that conservatively satisfies the safety constraint via the ϵ -uncertainty quantifier $(\cdot; \cdot)$ at every time step, with a probability of at least $1 - \epsilon$.

When we set ϵ to be a sufficiently large scalar such that $\epsilon > \frac{1}{2} \frac{V_{\max}}{r}$, the penalty b satisfies

$$\begin{aligned}
 b(s_h; a_h) &= \frac{c}{\min_{a \in A_h} (s_{h+1}; a)} \\
 &< \frac{\frac{1}{2} \frac{1}{r} V_{\max}}{\frac{1}{2}} \\
 &< \frac{1}{r} V_{\max}
 \end{aligned}$$

This means that, when the constraint violation happens even a single time, the value by a policy obtained in \mathcal{M} becomes negative because $\max_s V_r(s) = V_{\max}$.

Under Assumption 3.2, after convergence, the optimal policy π^* will not violate the safety constraint, and thus the emergency stop action will not be executed. In this case, the modified (unconstrained) MDP \mathcal{M}^* is identical to the original CMDP \mathcal{M} . Therefore, we now obtain the desired theorem. \square

E Supplementary materials regarding GLM-MASE

E.1 Pseudo-code for GLM-MASE

We first present the pseudo-code for GLM-MASE in Algorithm 2.

E.2 Proofs of Lemmas 5.4 and 5.5

Proof. See Lemma 1 (and Lemma 7) in Wang et al. [44]. \square

E.3 Preliminary Lemmas

Lemma E.1. Suppose the assumptions in Lemma 5.4 and 5.5 hold. Let C_1 and C_2 be positive, universal constants. Also, with a sufficiently large t , let t^* denote the smallest integer satisfying

$$\min(\cdot) \leq \frac{C_1}{t} \frac{1}{\ln t} \leq \frac{C_2}{t} \frac{1}{\ln t} \leq \frac{1}{2} C_g^{-1} \quad (15)$$

where $\min(\cdot)$ is the minimum eigenvalue of the second moment matrix. Then, we have

$$(s_h^{(t)}; a_h^{(t)}) \geq \frac{1}{2} \quad (16)$$

Proof. By Proposition 1 of Li et al. [22],

$$\min(s_h^{(t)}) \geq \min(\cdot) t H \geq C_1 \frac{\rho}{t H d} \geq C_2 \frac{\rho}{t H \ln^{-1}}.$$

By combining the aforementioned inequality with (15), we have

$$\min(s_h^{(t)}) \geq 2 C_g^{-1}.$$

Using the definition of $\max(s_h^{(t)}) = \frac{1}{\min(s_h^{(t)})}$, the following chain of equations hold for all $t \geq [t^?; T]$ and $h \geq [1; H]$:

$$\begin{aligned} (s_h^{(t)}; a_h^{(t)}) &= C_g \cdot k_{s;a} \cdot k_{h;t}^{-1} \\ &= C_g \cdot \max(s_h^{(t)}) \\ &= C_g \cdot \frac{1}{\min(s_h^{(t)})} \\ &\geq \frac{1}{2}. \end{aligned}$$

Therefore, we have the desired lemma. \square

E.4 Proof of Theorem 5.6

Proof. By definition, the GLM-MASE chooses actions satisfying

$$f(h e_h^{(t)}; b_{h;t}^{(t)}) \geq (s_h^{(t)}; a_h^{(t)}) \geq b_h. \quad (17)$$

By Lemma 5.4, the actions that are conservatively chosen based on (17) also satisfy the actual safety constraint, with a probability at least $1 - \delta$; that is,

$$g(s_h^{(t)}; a_h^{(t)}) \geq b_h. \quad (18)$$

In addition, when there is no action satisfying (17), the emergency stop action is executed where no unsafe action will be executed. Hence, the desired theorem is now obtained. \square

E.5 Proof of Theorem 5.7

By Assumption 3.2, the optimal policy π^* satisfies

$$g(s_h; \pi^*(s_h)) \geq b_h \quad ; \quad \forall h \geq [1; H]. \quad (19)$$

Thus, the set of state-action pairs that are potentially visited by π^* are written as

$$\mathcal{F}(s; a) \geq \mathcal{S} \cup \mathcal{A} \mid g(s; a) \geq b_h \geq g;$$

which satisfies the following chain of inequalities:

$$\begin{aligned} \mathcal{F}(s; a) \geq \mathcal{S} \cup \mathcal{A} \mid g(s; a) \geq b_h \geq g &\supseteq \mathcal{F}(s; a) \geq \mathcal{S} \cup \mathcal{A} \mid (s; a) \in \mathcal{S} \cup \mathcal{A} \mid b_h \geq g \\ &\supseteq \mathcal{F}(s; a) \geq \mathcal{S} \cup \mathcal{A} \mid (s; a) \in \mathcal{S} \cup \mathcal{A} \mid b_h \geq g. \end{aligned}$$

The state and action spaces in the last line represent the set of state-action pairs that may be visited by the policy obtained by the MASE algorithm. We used Lemma 5.4 in the first line and $(s; a) \in \mathcal{S} \cup \mathcal{A} \mid b_h \geq g$ in the second line.

By Lemma 8 in Strehl and Littman [34], the total regret can be decomposed into two parts as follows:

$$\sum_{t=t^?}^h V_r^{\pi^*} - V_r^{\pi} = R(T) + \sum_{t=t^?}^h V_{\max}; \quad (20)$$

where the first term is the regret under the condition that the confidence bound based on δ -uncertainty quantifier successfully contains the true safety function. Also, the second term is the regret under the opposite condition, which occurs with a probability δ .

The first term in (20) is upper-bounded based on Wang et al. [44] as follows:

$$R(T) \leq O\left(H^{\rho} \frac{1}{(T-t^?) \ln((T-t^?)H)}\right) + H^{-1} M + d^2 \ln \frac{1 + \max_{(T-t^?)H} (T-t^?)d \ln 1 + \frac{(T-t^?)}{d}}{1} \Theta(H^{\rho} d^{\beta} (T-t^?)):$$

As for the second term in (20), set $\epsilon = \frac{1}{TH}$ and then we have

$$\begin{aligned} \sum_{t=t^?}^T V_{\max} &= \sum_{t=t^?}^T V_{\max} \frac{1}{TH} \\ &= \sum_{t=t^?}^T \frac{1}{1} \frac{H}{1} \frac{1}{TH} \\ &= \sum_{t=t^?}^T H \frac{1}{TH} \\ &= O(1): \end{aligned}$$

In summary, the regret can be upper bounded by $\Theta(H^{\rho} d^{\beta} (T-t^?))$. We now obtain the desired theorem.

F Proof of Theorem 6.1

Lemma F.1 (δ -uncertainty quantifier). Assume $kgk_k^2 \leq B$ and $N_n \geq 1$ for all $n \geq 1$. Set

$$(s; a) := \sum_{j=1}^n n(s; a)_j; \quad \delta(s; a) \geq S \cup A \quad (21)$$

with $\frac{1}{n} := B + 4! \frac{\rho}{n+1 + \ln^{-1}}$, where n is the information capacity associated with kernel k . Then, δ is a δ -uncertainty quantifier.

Proof. Recall the assumption that $kgk_k^2 \leq B$ and $N_n \geq 1$, $\delta n \geq 1$. Also, set

$$\frac{1}{n} := B + 4! \frac{\rho}{n+1 + \ln^{-1}}:$$

By Theorem 2 in Chowdhury and Gopalan [12], we have

$$jg(s; a) \leq n(s; a)_j \leq n \sum_{j=1}^n n(s; a)_j; \quad \delta(s; a) \geq S \cup A$$

for all $n \geq 1$, with probability at least $1 - \frac{1}{n}$. Now, define

$$(s; a) := \sum_{j=1}^n n(s; a)_j; \quad \delta(s; a) \geq S \cup A \quad (22)$$

then we interpret that $\delta : S \cup A \rightarrow \mathbb{R}$ is a δ -uncertainty quantifier based on GP. \square

F.1 Proof of Theorem 6.1

Proof. Based on Lemma F.1, when there is at least one safe action (i.e., $A_n^+ \neq \emptyset$), the satisfaction of the safety constraint is guaranteed based on the δ -uncertainty quantifier with a probability at least $1 - \frac{1}{n}$. Also, if there is no safe action (i.e., $A_n^+ = \emptyset$), the emergency stop action is executed. In both cases, MASE guarantees the satisfaction of the safety constraint with probability at least $1 - \frac{1}{n}$. \square

Table 1: Hyper-parameters for Safety Gym experiments.

	NAME	VALUE
COMMON PARAMETERS	NETWORK ARCHITECTURE	[64;64]
	ACTIVATION FUNCTION	tanh
	LEARNING RATE (CRITIC)	$5 \cdot 10^{-3}$
	LEARNING RATE (POLICY)	$3 \cdot 10^{-4}$
	LEARNING RATE (PENALTY)	$3 \cdot 10^{-2}$
	DISCOUNT FACTOR (REWARD)	0.99
	DISCOUNT FACTOR (SAFETY)	0.99
	STEPS PER EPOCH	10;000
	NUMBER OF GRADIENT STEPS	80
	NUMBER OF EPOCHS	500
	TARGET KL	0.01
TRPO & CPO	DAMPING COEFFICIENT	0.1
	BACKTRACK COEFFICIENT	0.8
	BACKTRACK ITERATIONS	10
	LEARNING MARGIN	FALSE
MASE	PENALTY FOR EMERGENCY STOP ACTIONS	1
	DEEP GP NETWORK ARCHITECTURE	[16;16]
	NUMBER OF INDUCING POINTS	128
	KERNEL FUNCTION	RADIAL BASIS FUNCTION

G Details of Safety-Gym Experiment

We present the details regarding our experiments using Safety-Gym. Our experimental setting is based on Sootla et al. [32], which is slightly different from the original Safety Gym in that the obstacles (i.e., unsafe region) are replaced deterministically so that the environment is solvable and there is a viable solution. In this experiment, we used a machine with Intel(R) Xeon(R) Silver 4210 CPU, 128GB RAM, and NVIDIA A100 GPU. For a fair comparison, we basically used the same hyper-parameter as in Sootla et al. [32], which is summarized in Table 1.

In our experiment, when the agent identified that there was no safe action based on the GP-based uncertainty quantifier, we simply terminated the current episode (i.e., resetting) immediately after the emergency stop action and started the new episode. The frequency of the emergency stop actions is shown in Table 2.

Table 2: Frequency of the emergency stop actions.

TASK	TOTAL	LAST 100 EPISODES
POINTGOAL1	154/500	24/100
CARGOAL1	397/500	46/100

The emergency stop action is a variant of so-called resetting actions that are common in episodic RL settings, which prevent the agent from exploring the state-action spaces since the uncertainty quantifier is sometimes quite conservative. We consider that this is the reason why the reward performance of our MASE is worse than other methods (e.g., TRPO-Lagrangian, CPO). However, because we require the agent to solve more difficult problems where safety is guaranteed at every time step and episode, we consider that this result is reasonable to some extent. Though it is better for an algorithm for such a severe safety constraint to have a comparable performance as CPO, we will leave it for future work.

We also conducted an experiment to compare the performance of MASE with the early-terminated MDP (ET-MDP, [37]) algorithm. The ET-MDP is an algorithm to execute emergency stop actions *immediately after* safety constraints are violated. Figure 4 shows the experimental results. The ET-MDP and MASE exhibit similar learning curves on the average episode reward and average episode safety. However, while ET-MDP violated the safety constraint in most episodes (i.e., almost

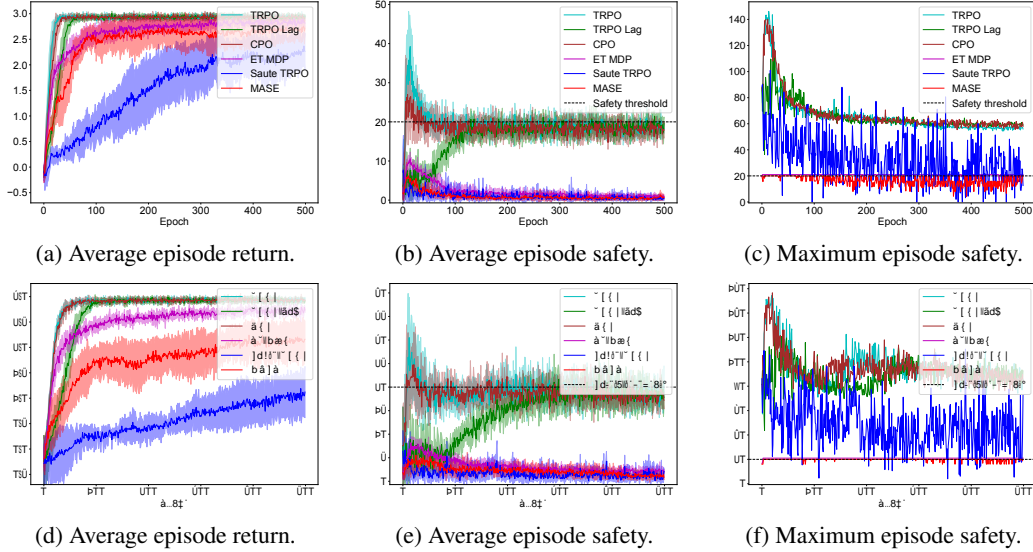


Figure 4: Experimental results on Safety Gym (Top: PointGoal1, Bottom: CarGoal1) with an additional implementation of the Early-terminated MDP (ET-MDP) algorithm (Sun et al. [37]).

all episodes are terminated after an unsafe action is executed), MASE did not violate any safety constraint.

H Grid-world Experiment

We also conduct an experiment using the grid-world problem as in Wachi and Sui [41]. Experimental settings are based on their original implementation (https://github.com/akifumi-wachi-4/safe_near_optimal_mdp). We consider a 20 × 20 square grid in which reward and safety functions are randomly generated. There are two types regarding the safety threshold: one is time-invariant as in [41] and the other is time-variant as in the GSE problem.

We run SNO-MDP [41] and MASE in 100 randomly generated environments, and we compute the reward collected by the algorithms and count the number of episodes in which the safety constraint is violated at least once. The reward is normalized with respect to that by SNO-MDP.

Table 3: Experimental results for grid-world experiments.

	TIME-INVARIANT SAFETY THRESHOLD			TIME-VARIANT SAFETY THRESHOLD		
	REWARD		SAFETY VIOLATION	REWARD		SAFETY VIOLATION
SNO-MDP [41]	1:0	0:0	0	1:0	0:0	87
MASE (OURS)	1:0	0:0	0	2:4	1:0	0

The experimental results are shown in Table 3. When the safety threshold is time-invariant, MASE behaves identically with SNO-MDP; thus, the performance of the MASE is comparable with that of SNO-MDP. When the safety threshold is time-variant, SNO-MDP cannot deal with it by nature; hence, the safety constraint is not satisfied in most of the episodes. In contrast, our MASE satisfies the safety constraint in every episode, which also contributes to the larger reward.