
Safe Reinforcement Learning with Contrastive Risk Prediction

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 As safety violations can lead to severe consequences in real-world applications, the
2 increasing deployment of Reinforcement Learning (RL) in safety-critical domains
3 such as robotics has propelled the study of safe exploration for reinforcement
4 learning (safe RL). In this work, we propose a risk preventive training method for
5 safe RL, which learns a binary classifier based on contrastive sampling to predict
6 the probability of a state-action pair leading to unsafe states. Based on the predicted
7 risk probabilities, risk preventive trajectory exploration and optimality criterion
8 modification can be simultaneously conducted to induce safe RL policies. We
9 conduct experiments in robotic simulation environments. The results show the
10 proposed approach outperforms existing model-free safe RL approaches, and yields
11 comparable performance with the state-of-the-art model-based method.

12 1 Introduction

13 Reinforcement Learning (RL) offers a great set of technical tools for many real-world decision
14 making systems, such as robotics, that require an agent to automatically learn behavior policies
15 through interactions with the environments [1]. Conversely, the applications of RL in real-world
16 domains also pose important new challenges for RL research. In particular, many real-world robotic
17 environments and tasks, such as human-related robotic environments [2], helicopter manipulation
18 [3, 4], autonomous vehicle [5], and aerial delivery [6], have very low tolerance for violations of safety
19 constraints, as such violation can cause severe consequences. This raises a substantial demand for
20 safe reinforcement learning techniques.

21 Safe reinforcement learning investigates RL methodologies with critical safety considerations, and
22 has received increased attention from the RL research community. In safe RL, in addition to the
23 reward function [7], an RL agent often deploys a cost function to maximize the discounted cumulative
24 reward while satisfying the cost constraint [8–10]. A comprehensive survey of safe RL categorizes the
25 safe RL techniques into two classes: modification of the optimality criterion and modification of the
26 exploration process [11]. For modification of the optimality criterion, previous works mostly focus
27 on the modification of the reward. Many works [12–17] pursue such modifications by shaping the
28 reward function with penalizations induced from different forms of cost constraints. For modification
29 of the exploration process, safe RL approaches focus on training RL agents on modified trajectory
30 data. For example, some works deploy backup policies to recover from safety violations to safer
31 trajectory data that satisfy the safety constraint [18–20].

32 In this paper, we propose a novel risk preventive training (RPT) method to tackle the safe RL
33 problem. The key idea is to learn a contrastively estimated classification model to predict the risk—
34 the probability of a state-action pair leading to unsafe states, which can then be deployed to modify
35 both the exploration process and the optimality criterion. In terms of exploration process modification,
36 we collect trajectory data in a risk preventive manner based on the predicted probability of risk. A

37 trajectory is terminated if the next state falls into an unsafe region that has above-threshold risk values.
38 Regarding optimality criterion modification, we reshape the reward function by penalizing it with the
39 predicted risk for each state-action pair. Benefiting from the generalizability of risk prediction, the
40 proposed approach can avoid safety constraint violations much early in the training phase and induce
41 safe RL policies, while previous works focus on backup policy and violate more safety constraints
42 by interacting with the environment in the unsafe regions. Moreover, we further deploy a simple
43 unsafe-state augmentation strategy for the proposed method to increase the sample efficiency of the
44 encountered unsafe states and reduce the safety violations of the RL agent in the experiments. We
45 conduct experiments using four robotic simulation environments on MuJoCo [21]. Our model-free
46 approach produces comparable performance with a state-of-the-art model-based safe RL method
47 SMBPO [16] and greatly outperforms other model-free safe RL methods. The main contributions of
48 the proposed work can be summarized as follows:

- 49 • This is the first work that introduces a contrastive sampling based classifier to perform risk
50 prediction and conduct safe RL exploration.
- 51 • With its proficient risk prediction capabilities, the proposed approach possesses the essential
52 capacity to simultaneously modify the exploration process through risk preventive trajectory
53 collection and adjust the optimality criterion through reward reshaping.
- 54 • As a model-free safe RL method, the proposed approach achieves comparable performance
55 to the state-of-the-art model-based safe RL method and outperforms the model-free methods
56 in multiple benchmark robotic simulation environments.

57 2 Related Works

58 Many methods have been developed for safe RL. Garcia and Fernández [11] provided a survey
59 categorizing safe RL methods into categories of modifying the optimality criterion and modifying the
60 exploration process.

61 **Modification of the optimality criterion.** Since optimizing the conventional reward signal does
62 not ensure the avoidance of safety violations, leading to the exploration of modifying the optimality
63 objective based on risk notions [22, 23], probabilities of visiting risky states [24], etc. Achiam
64 et al. [20] proposed Constrained Policy Optimization (CPO) to update safe policies by optimizing
65 the primal-dual problem in trust regions. Recently, reward shaping techniques [25, 26] have been
66 integrated into safe RL. Tessler et al. [14] introduced Reward Constrained Policy Optimization
67 (RCPO) by penalizing the normal training policy. Thomas et al. [16] reshaped reward functions
68 using a model-based predictor, treating unsafe states as absorbing states to train the RL agent with
69 penalized rewards. Xu et al. [27] developed Constrained Penalized Q-learning (CPQ) using a cost
70 critic to learn constraint values during exploration and penalizing the Bellman operator in policy
71 training to stop the updates for potentially unsafe states.

72 **Modification of the exploration process.** Previous works have optimized safe RL policies by
73 adjusting exploration processes during interaction with the environment. For instance, [28, 3,
74 29] guided exploration based on prior environmental knowledge. Similarly, [30, 31] constrained
75 exploration learning using demonstration data. More recent approaches like [18, 19] focused on
76 utilizing backup policies from safe regions to prevent safety violations. If the agent undertakes a
77 potentially risky action, the task policy is replaced with a guaranteed safe backup policy. Yu et al.
78 [32] defined safe regions as feasible sets and used reachability analysis to expand these sets beyond
79 traditional energy-based methods. Jayant and Bhatnagar [33] introduced a model-based deep RL
80 agent that efficiently learns an ensemble of transition dynamics in an online environment and restricts
81 exploration with a performance ratio.

82 Safe RL is crucial in environments like human-related robotic settings where safety violations can
83 lead to catastrophic failures [2]. Robotic simulation environments such as MuJoCo, developed by
84 Todorov et al. [21], facilitate research in RL applications for robotics. Thomas et al. [16] extended
85 the MuJoCo environment to define safety violations in robotic simulations, making it an ideal test
86 bed for safe RL methods.

87 **3 Preliminary**

88 Reinforcement learning (RL) has been broadly used to train robotic agents by maximizing the
 89 discounted cumulative rewards. The representation of a reinforcement learning problem can be
 90 formulated as a Markov Decision Process (MDP) $M = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$ [7], where \mathcal{S} is the state
 91 space for all observations, \mathcal{A} is the action space for available actions, $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ is the transition
 92 dynamics, $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow [r_{min}, r_{max}]$ is the reward function, and $\gamma \in (0, 1)$ is the discount factor. An
 93 agent can start from a random initial state s_0 to take actions and interact with the MDP environment
 94 by receiving rewards for each action and moving to new states. Such interactions can produce a
 95 transition (s_t, a_t, r_t, s_{t+1}) at each time-step t with $s_{t+1} = \mathcal{T}(s_t, a_t)$ and $r_t = r(s_t, a_t)$, while a
 96 sequence of transitions comprise a trajectory $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{|\tau|+1})$, where $|\tau| + 1$
 97 denotes the length of trajectory τ —i.e., the number of transitions. The goal of RL is to learn an
 98 optimal policy $\pi^* : \mathcal{S} \rightarrow \mathcal{A}$ that can maximize the expected discounted cumulative reward (return):
 99 $\pi^* = \arg \max_{\pi} J_r(\pi) = \mathbb{E}_{\tau \sim \mathcal{D}_{\pi}} [\sum_{t=0}^{|\tau|} \gamma^t r_t]$

100 **3.1 Safe Exploration for Reinforcement Learning**

101 Safe exploration for Reinforcement Learning (safe RL) studies RL with critical safety considerations.
 102 For a safe RL environment, in addition to the reward function, a cost function can also exist to reflect
 103 the risky status of each exploration step. The process of safe RL can be formulated as a Constrained
 104 Markov Decision Process (CMDP) [34], $\hat{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma, c, d)$, which introduces an extra cost
 105 function c and a cost threshold d into MDP. An exploration trajectory under CMDP can be written
 106 as $\tau = (s_0, a_0, r_0, c_0, s_1, \dots, s_{|\tau|+1})$, where the transition at time-step t is $(s_t, a_t, s_{t+1}, r_t, c_t)$, with
 107 a cost value c_t induced from the cost function $c_t = c(s_t, a_t)$. CMDP monitors the safe exploration
 108 process by requiring the cumulative cost $J_c(\pi)$ does not exceed the cost threshold d , where $J_c(\pi)$
 109 can be defined as the expected total cost of the exploration, $J_c(\pi) = \mathbb{E}_{\tau \sim \mathcal{D}_{\pi}} [\sum_{t=0}^{|\tau|} c_t]$ [12]. Safe
 110 RL hence aims to learn an optimal policy π^* that can maximize the expected discounted cumulative
 111 reward subjecting to a cost constraint, as follows:

$$\begin{aligned} \pi^* = \arg \max_{\pi} J_r(\pi) &= \mathbb{E}_{\tau \sim \mathcal{D}_{\pi}} \left[\sum_{t=0}^{|\tau|} \gamma^t r_t \right] \\ \text{s.t. } J_c(\pi) &= \mathbb{E}_{\tau \sim \mathcal{D}_{\pi}} \left[\sum_{t=0}^{|\tau|} c_t \right] \leq d. \end{aligned} \quad (1)$$

112 **4 Method**

113 Robot operations typically have low tolerance for risky/unsafe states and actions, since a robot could
 114 be severely damaged in real-world environments when the safety constraint being violated. Similar to
 115 the work in [9], in this work we adopt a strict setting for the safety constraint such that any “unsafe”
 116 state can cause violation of the safety constraint and the RL agent will terminate an exploration
 117 trajectory when encountering an “unsafe” state. We have the following definition:

118 **Definition 1.** *For a state s and an action a , the value of the cost function $c(s, a)$ can either be 0*
 119 *or 1. When $c(s, a) = 0$, the induced state $\mathcal{T}(s, a)$ is defined as a safe state; when $c(s, a) = 1$, the*
 120 *induced state $\mathcal{T}(s, a)$ is defined as an unsafe state, which triggers the violation of safety constraints*
 121 *and hence causes the termination of the trajectory.*

122 Based on this definition, the cost threshold d in Eq. (1) should be set strictly to 0. The agent is
 123 expected to learn a safe policy π that can operate with successful trajectories containing only safe
 124 states. Towards this goal, we propose a novel risk prediction method for safe RL. The proposed
 125 method deploys a contrastive classifier to predict the probability of a state-action pair leading to
 126 unsafe states, which can be trained during the exploration process of RL and generalized to previously
 127 unseen states.

128 With risk prediction probabilities, a more informative cumulative cost $J_c(\pi)$ can be formed to prevent
 129 unsafe trajectories and reshape the reward in each transition of a trajectory to induce safe RL policies.
 130 Previous safe RL methods in the literature can typically be categorized into two classes: modification
 131 of the optimality criterion and modification of the exploration process [11]. With safety constraints
 132 and risk predictions, the proposed approach (to be elaborated below) has the capacity and is expected
 133 to incorporate the strengths of both categories of safe RL techniques.

134 **4.1 Risk Prediction with Contrastive Classification**

135 Although an RL agent would inevitably encounter unsafe states during the initial stage of the
 136 exploration process in an unknown environment, we aim to quickly learn from the unsafe experience
 137 through statistical learning and generalize the recognition of unsafe trajectories to prevent risk for
 138 future exploration. Specifically, we aim to compute the probability of a state-action pair leading to
 139 unsafe states, i.e., $p(y = 1|s_t, a_t)$, where $y \in \{0, 1\}$ denotes a random variable that indicates whether
 140 (s_t, a_t) leads to an unsafe state $s_u \in S_U$. The set of unsafe states, S_U , can be either pre-given
 141 or collected during initial exploration. However, directly training a binary classifier to make such
 142 predictions is impractical as it is difficult to judge whether a state-action pair is *safe*—i.e., never
 143 leading to unsafe states.

144 For this purpose, we propose to train a contrastive classifier $F_\theta(s_t, a_t)$ with model parameter θ to
 145 discriminate a positive state-action pair (s_t, a_t) in a trajectory that leads to unsafe states (unsafe
 146 trajectory) against random state-action pairs from the overall distribution of any trajectory. Such a
 147 contrastive form of learning can conveniently avoid the impractical identification problem of absolute
 148 negative (i.e., safe) state-action pairs.

149 Specifically, inspired by the noise contrastive estimation based classifier design in the literature
 150 [35, 36], we propose to learn $F_\theta(s_t, a_t)$ as a binary classifier via weighted contrastive sampling by
 151 sampling unsafe state-action pairs as positive samples and sampling general state-action pairs as
 152 contrastive negative samples. Let $p(s_t, a_t|y = 1)$ denote the presence probability of a state-action pair
 153 (s_t, a_t) in a trajectory that leads to unsafe states, and $p(y = 1)$ denote the distribution probability of
 154 unsafe trajectory in the environment. The contrastive classifier $F_\theta(s_t, a_t)$ is then defined as follows:

$$F_\theta(s_t, a_t) = \frac{p(s_t, a_t|y = 1)p(y = 1)}{p(s_t, a_t|y = 1)p(y = 1) + p(s_t, a_t)}, \quad (2)$$

155 where $p(y = 1)$ is used as weight for the positive samples which are only from the unsafe trajectories,
 156 and weight 1 is given to the *contrastively-negative* samples which are from the overall distribution.
 157 This binary classifier identifies the state-action pairs in unsafe trajectories contrastively from general
 158 pairs in the overall distribution.

159 From the definition of $F_\theta(s_t, a_t)$ in Eq.(2), one can derive the probability of interest, $p(y = 1|s_t, a_t)$,
 160 using the Bayes’ theorem, as follows:

$$p(y = 1|s_t, a_t) = \frac{p(s_t, a_t|y = 1)p(y = 1)}{p(s_t, a_t)} = \frac{F_\theta(s_t, a_t)}{1 - F_\theta(s_t, a_t)}, \quad (3)$$

161 where the derivation from the fraction in the top row to the term expressed in F_θ in the second row can
 162 be done easily by dividing both the numerator and denominator of the top row fraction with the same
 163 term $[p(s_t, a_t|y = 1)p(y = 1) + p(s_t, a_t)]$. As the normal output range, $[0, 1]$, of the probabilistic
 164 classifier $F_\theta(s_t, a_t)$ could lead to unbounded values $p(y = 1|s_t, a_t) \in [0, \infty]$ through Eq.(3), we
 165 propose to first rescale the output of classifier $F_\theta(s_t, a_t)$ to the range of $[0, 0.5]$ when calculating
 166 $p(y = 1|s_t, a_t)$ via Eq.(3).

167 Based on the contrastive sampling principle of F_θ , we optimize the contrastive classifier’s parameter
 168 θ using maximum likelihood estimation (MLE) with the following log-likelihood objective function:

$$L(\theta) = \mathbb{E}_{p(s_t, a_t|y=1)p(y=1)} [\log F_\theta(s_t, a_t)] + \mathbb{E}_{p(s_t, a_t)} [\log(1 - F_\theta(s_t, a_t))]. \quad (4)$$

169 By setting the derivative of $L(\theta)$ w.r.t. F_θ to zero, it is easy to verify that the definition of F_θ in Eq.(2)
 170 can achieve the maximum of this MLE objective w.r.t. F_θ .

171 **4.2 Risk Preventive Trajectory**

172 Based on Definition 1, a trajectory terminates when the RL agent encounters an unsafe state and
 173 triggers safety constraint violation. It is however desirable to minimize the number of such safety
 174 violations during the policy training process and learn a good policy in safe regions. The risk
 175 prediction classifier we proposed above provides a convenient tool for this purpose by predicting
 176 the probability of a state-action pair leading to unsafe states, $p(y = 1|s_t, a_t)$. Based on this risk
 177 prediction, we have the following definition for unsafe regions:

178 **Definition 2.** A state-action pair (s_t, a_t) falls into an **unsafe region** if the probability of (s_t, a_t)
 179 leading to unsafe states is greater than a threshold η : $p(y = 1|s_t, a_t) > \eta$, where $\eta \in (0, 1)$.

180 With this definition, an RL agent can pursue risk preventive trajectories to avoid safety violations by
 181 staying away from unsafe regions. Specifically, we can terminate a trajectory before violating the
 182 safety constraint by judging the potential risk—*i.e.*, the probability of $p(y = 1|s_t, a_t)$.

183 Without a doubt, the threshold η is a key for determining the length $T = |\hat{\tau}|$ of an early stopped risk
 184 preventive trajectory $\hat{\tau}$. We make the following assumption for deriving a lemma:

185 **Assumption 1.** For a trajectory $\tau = \{s_0, a_0, r_0, c_0, s_1, \dots, s_H\}$ that leads to an unsafe state
 186 $s_H \in S_U$, the risk prediction probability $p(y = 1|s_t, a_t)$ increases linearly along transition steps
 187 within a base neighborhood of the unsafe region that can be defined through $p(y = 1|s_t, a_t) \geq \eta_b$
 188 with a threshold $\eta_b \in (0, \eta)$.

189 **Lemma 1.** Assume Assumption 1 holds. Let H denote the length of an unsafe trajectory $\tau =$
 190 $\{s_0, a_0, r_0, c_0, s_1, \dots, s_H\}$ that terminates at an unsafe state $s_H \in S_U$. The numbers of transition
 191 steps, T and T_b , along this trajectory to the unsafe region determined by η in Definition 2 and its
 192 neighborhood determined by η_b , respectively, satisfy $T \approx \lfloor \frac{\eta - \eta_b}{1 - \eta_b} H + \frac{1 - \eta}{1 - \eta_b} T_b \rfloor$.

193 This lemma demonstrates the influence of the risk control threshold η on the length of collected
 194 trajectories. Given η_b (and hence T_b), a larger η value will allow more effective explorations with
 195 longer trajectories to facilitate policy learning, but also tighten the unsafe region and increase the
 196 possibility of violating safety constraints.

197 4.3 Risk Preventive Reward Shaping

198 With Definition 1, the safe RL formulation in Eq. (1) can hardly induce a safe policy since there are
 199 no intermediate costs before encountering an unsafe state. With the risk prediction classifier proposed
 200 above, we can rectify this drawback by defining the cumulative cost function $J_c(\pi)$ using the risk
 201 prediction probabilities, $p(y = 1|s_t, a_t)$, over all encountered state-action pairs. Specifically, we
 202 adopt a reward-like discounted cumulative cost as follows:

$$J_c(\pi) = \mathbb{E}_{\tau \sim \mathcal{D}_\pi} \left[\sum_{t=0}^{|\tau|} \gamma^t p(y = 1|s_t, a_t) \right], \quad (5)$$

203 which uses the predicted risk as the estimated cost. Moreover, instead of solving safe RL as a
 204 constrained discounted cumulative reward maximization problem, we propose to use Lagrangian
 205 relaxation [37] to convert the constrained maximization problem, CMDP, in Eq. (1) into an un-
 206 constrained optimization problem, which is equivalent to shaping the reward function \mathcal{R} with risk
 207 penalties:

$$\min_{\lambda \geq 0} \max_{\pi} [J_r(\pi) - \lambda(J_c(\pi) - d)] \quad (6)$$

$$\iff \min_{\lambda \geq 0} \max_{\pi} [J_r(\pi) - \lambda J_c(\pi)] \quad (7)$$

$$\iff \min_{\lambda \geq 0} \max_{\pi} \mathbb{E}_{\tau \sim \mathcal{D}_\pi} \left[\sum_{t=0}^{|\tau|} \gamma^t (r_t - \lambda p(y = 1|s_t, a_t)) \right] \quad (8)$$

208 where $r_t - \lambda p(y = 1|s_t, a_t)$ can be treated as the risk penalty reshaped reward. The Lagrangian dual
 209 variable λ controls the degree of reward shaping with the predicted risk value.

210 4.4 Risk Preventive Training Algorithm

211 The overall risk preventive RL training procedure for the proposed safe RL method is presented in
 212 Algorithm 1, which trains a contrastive classifier F_θ (line 21) for risk prediction, and performs safe
 213 reinforcement learning by simultaneously enforcing risk preventive trajectory exploration (line 15-17)
 214 and risk preventive reward shaping (line 13).

215 4.5 Data Augmentation for Contrastive Learning

216 As the goal of safe RL is to minimize the encountering of unsafe states, it is desirable to produce
 217 an effective risk predictor with very limited risky state-action pairs. To this end, we propose to
 218 extend RPT by designing a simple *data augmentation* procedure, producing a data augmented
 219 method, RPT+DA, for comparison. The proposed data augmentation solely *enhances the training*

Algorithm 1 Risk Preventive Training

Input: Initial policy π_ϕ , classifier F_θ , trajectory set $D = \emptyset$, set of unsafe state-action pairs S_U , threshold η, η_b ; penalty factor λ , set of unsafe trajectory length $\mathcal{H} = \emptyset$
Output: Trained policy π_ϕ

- 1: **for** $k = 1, 2, \dots, K$ **do**
- 2: $T_b = 0$
- 3: **for** $t = 0, 1, \dots, T_{\max}$ **do**
- 4: Sample transition $(s_t, a_t, r_t, c_t, s_{t+1})$ from the environment with policy π_ϕ .
- 5: **if** $c_t > 0$ **then**
- 6: Add the risky state-action (s_t, a_t) into S_U ; add length t to \mathcal{H} .
- 7: Increase λ if necessary
- 8: Stop trajectory and break.
- 9: **end if**
- 10: Sample next action a_{t+1} as $a_{t+1} = \pi_\phi(\cdot|s_{t+1})$.
- 11: Compute p_t and p_{t+1} via Eq. (3)
- 12: **If** $p_t \geq \eta_b$, then set $T_b = t$
- 13: Penalize reward r_t with p_t : $\hat{r}_t = r_t - \lambda p_t$
- 14: Add transition to the trajectory set, such that: $D = D \cup (s_t, a_t, \hat{r}_t, s_{t+1})$
- 15: **if** $p_{t+1} > \eta$ **then**
- 16: Stop trajectory and break.
- 17: **end if**
- 18: **end for**
- 19: Sample risky state-action pairs from S_U
- 20: Sample transitions from D : $(s_t, a_t, \hat{r}_t, s_{t+1}) \sim D$
- 21: Update classifier F_θ by maximizing $L(\theta)$ in Eq (4)
- 22: Update policy π_ϕ with shaped reward $J_{\hat{r}}(\pi)$ in Eq (8)
- 23: **end for**

220 *of contrastive classifier for risk prediction*, with no additional interaction with the environment or
221 trajectory generation. Specifically, we perform data augmentation only for the data sampled from the
222 set of risky states S_U . For each sampled risky state-action pair (s_t, a_t) , we propose to produce an
223 augmented state \hat{s}_t by adding a random Gaussian noise sampled from the standard normal distribution
224 $\mathcal{N}(0, 1)$ to each entry of the observed data s_t . We can repeat this process to generate multiple (e.g.,
225 n) augmented states for each s_t . In our experiments, we used $n = 3$. Together with a_t , each \hat{s}_t can
226 be used to form an additional risky state-action pair (\hat{s}_t, a_t) for training the contrastive classifier. The
227 hypothesis is that without any prior information about the environment, the training of the proposed
228 contrastive classifier highly depends on the data collected during the agent’s interactions with the
229 environment, especially on the limited number of observed unsafe states. By using the proposed data
230 augmentation technique above, we expect to improve the unsafe states’ sample efficiency and the
231 generalizability of the approach on discriminating unsafe states and hence reduce the possible safety
232 violations during the exploration process.

233 5 Experiment

234 5.1 Experimental Settings

235 **Experimental Environments** Following the experimental setting in [16], we adopted four robotics
236 simulation environments, *Ant*, *Cheetah*, *Hopper*, and *Humanoid*, based on the MuJoCo simulator
237 [21]. For *Ant* and *Hopper*, a robot violates the safety constraint when it falls over. For *Cheetah*, a
238 robot violates the safety constraint when its head flips on the ground, which is modified from the
239 HalfCheetah environment with extra safety constraint [16]. For *Humanoid*, the human-like robot
240 violates the safety constraint when the head of the robot falls to the ground. The RL agent cumulates
241 returns by operating in the environment. As shown in Algorithm 1, the RL trajectory terminates when
242 either the RL agent encounters safety violation, the maximum length is reached, or the preventive
243 trajectory break takes place.

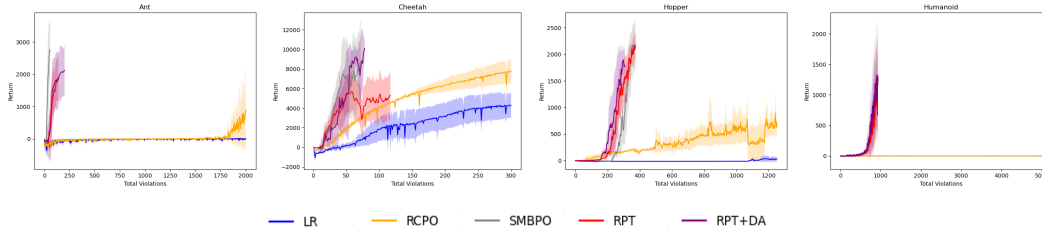


Figure 1: For each method, each plot presents the undiscounted return vs. the total number of violations. The curve shows the mean of the return over five runs, while the shadow shows the standard deviation.

244 **Comparison Methods** We compare the proposed Risk Preventive Training (RPT) approach with
 245 three state-of-the-art safe RL methods: SMBPO [16], RCPO [14], and LR [12].

246 5.2 Experimental Results

247 We compared all the five methods (LR, RCPO, SMBPO, RPT, and RPT+DA) by running each method
 248 five times with random seeds in each of the four MuJoCo environments. The performance of each
 249 method is evaluated by presenting the corresponding return vs. the total number of violations obtained
 250 in the training process. The results for all the methods are presented on the left side of Figure 1,
 251 one plot for each robotic simulation environment. The curve for each method shows the learning
 252 ability of the RL agent with limited safety violations. From the plots, we can see RPT, RPT+DA and
 253 SMBPO achieve large returns with a small number of violations on all the four robotic tasks, and
 254 largely outperform the other two methods, RCPO and LR, which have much smaller returns even
 255 with large numbers of safety violations. The proposed model-free RPT produces slightly inferior
 256 performance than the model-based SMBPO on *Ant* and *Cheetah*, where RPT requires more examples
 257 of unsafe states to yield good performance at the initial training stage. Nevertheless, RPT outperforms
 258 SMBPO on both *Hopper* and *Humanoid* with smaller number of safety violations. As a model-free
 259 safe RL method, RPT produces an overall comparable performance with the model-based method
 260 SMBPO. With data augmentation, RPT+DA further improves the performance of RPT on all the four
 261 environments, which demonstrates the efficacy of our simple unsafe-state augmentation strategy.

262 6 Conclusion

263 Inspired by the increasing demands for safe exploration of Reinforcement Learning, we proposed a
 264 novel model-free risk preventive training method, RPT, to perform safe RL by learning a contrastive-
 265 sampling based binary classifier to predict the probability of a state-action pair leading to unsafe
 266 states. Based on risk prediction, we produce a systematic scheme to collect risk preventive trajectories
 267 that terminate early without triggering safety constraint violations. Moreover, the predicted risk
 268 probabilities are also used as penalties to perform reward shaping for learning safe RL policies. A
 269 simple data augmentation strategy has also been deployed to improve the efficiency of the observed
 270 unsafe-states for RPT. We compared the proposed approach with a few state-of-the-art safe RL meth-
 271 ods using four robotic simulation environments. The proposed approach demonstrates comparable
 272 performance with the state-of-the-art model-based method and outperforms the model-free safe RL
 273 methods.

274 References

- 275 [1] J. Kober, J. A. Bagnell, and J. Peters, “Reinforcement learning in robotics: A survey,” *The*
 276 *International Journal of Robotics Research*, 2013.
- 277 [2] L. Brunke, M. Greeff, A. W. Hall, Z. Yuan, S. Zhou, J. Panerati, and A. P. Schoellig, “Safe
 278 learning in robotics: From learning-based control to safe reinforcement learning,” *Annual*
 279 *Review of Control, Robotics, and Autonomous Systems*, 2021.

- 280 [3] J. A. Martín H and J. d. Lope, “Learning autonomous helicopter flight with evolutionary
281 reinforcement learning,” in *International Conference on Computer Aided Systems Theory*.
282 Springer, 2009.
- 283 [4] R. Koppejan and S. Whiteson, “Neuroevolutionary reinforcement learning for generalized
284 control of simulated helicopters,” *Evolutionary intelligence*, 2011.
- 285 [5] L. Wen, J. Duan, S. E. Li, S. Xu, and H. Peng, “Safe reinforcement learning for autonomous
286 vehicles through parallel constrained policy optimization,” in *IEEE International Conference
287 on Intelligent Transportation Systems (ITSC)*. IEEE, 2020.
- 288 [6] A. Faust, I. Palunko, P. Cruz, R. Fierro, and L. Tapia, “Automated aerial suspended cargo
289 delivery through reinforcement learning,” *Artificial Intelligence*, 2017.
- 290 [7] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- 291 [8] O. Mihatsch and R. Neuneier, “Risk-sensitive reinforcement learning,” *Machine learning*, 2002.
- 292 [9] A. Hans, D. Schneegaß, A. M. Schäfer, and S. Udluft, “Safe exploration for reinforcement
293 learning,” in *ESANN*, 2008.
- 294 [10] Y. J. Ma, A. Shen, O. Bastani, and J. Dinesh, “Conservative and adaptive penalty for model-based
295 safe reinforcement learning,” in *Proceedings of the AAAI Conference on Artificial Intelligence
296 (AAAI)*, 2022.
- 297 [11] J. Garcia and F. Fernández, “A comprehensive survey on safe reinforcement learning,” *Journal
298 of Machine Learning Research (JMLR)*, pp. 1437–1480, 2015.
- 299 [12] A. Ray, J. Achiam, and D. Amodei, “Benchmarking safe exploration in deep reinforcement
300 learning,” *arXiv preprint arXiv:1910.01708*, 2019.
- 301 [13] L. Shen, L. Yang, S. Chen, B. Yuan, X. Wang, D. Tao *et al.*, “Penalized proximal policy
302 optimization for safe reinforcement learning,” *arXiv preprint arXiv:2205.11814*, 2022.
- 303 [14] C. Tessler, D. J. Mankowitz, and S. Mannor, “Reward constrained policy optimization,” in
304 *International Conference on Learning Representations (ICLR)*, 2019.
- 305 [15] Y. Hu, W. Wang, H. Jia, Y. Wang, Y. Chen, J. Hao, F. Wu, and C. Fan, “Learning to utilize
306 shaping rewards: A new approach of reward shaping,” in *Advances in Neural Information
307 Processing Systems (NeurIPS)*, 2020.
- 308 [16] G. Thomas, Y. Luo, and T. Ma, “Safe reinforcement learning by imagining the near future,” in
309 *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- 310 [17] J. Zhang, B. Cheung, C. Finn, S. Levine, and D. Jayaraman, “Cautious adaptation for reinforce-
311 ment learning in safety-critical settings,” in *International Conference on Machine Learning
312 (ICML)*, 2020.
- 313 [18] B. Thananjeyan, A. Balakrishna, S. Nair, M. Luo, K. Srinivasan, M. Hwang, J. E. Gonzalez,
314 J. Ibarz, C. Finn, and K. Goldberg, “Recovery rl: Safe reinforcement learning with learned
315 recovery zones,” *IEEE Robotics and Automation Letters*, 2021.
- 316 [19] O. Bastani, S. Li, and A. Xu, “Safe reinforcement learning via statistical model predictive
317 shielding,” in *Robotics: Science and Systems*, 2021.
- 318 [20] J. Achiam, D. Held, A. Tamar, and P. Abbeel, “Constrained policy optimization,” in *International
319 conference on machine learning (ICML)*, 2017.
- 320 [21] E. Todorov, T. Erez, and Y. Tassa, “Mujoco: A physics engine for model-based control,” in
321 *IEEE/RSJ international conference on intelligent robots and systems*. IEEE, 2012.
- 322 [22] R. A. Howard and J. E. Matheson, “Risk-sensitive markov decision processes,” *Management
323 science*, 1972.

- 324 [23] M. Sato, H. Kimura, and S. Kobayashi, “Td algorithm for the variance of return and mean-
325 variance reinforcement learning,” *Transactions of the Japanese Society for Artificial Intelligence*,
326 2001.
- 327 [24] P. Geibel and F. Wyszotzki, “Risk-sensitive reinforcement learning applied to control under
328 constraints,” *Journal of Artificial Intelligence Research (JAIR)*, 2005.
- 329 [25] M. Dorigo and M. Colombetti, “Robot shaping: Developing autonomous agents through
330 learning,” *Artificial intelligence*, 1994.
- 331 [26] J. Randoøv and P. Alstrøm, “Learning to drive a bicycle using reinforcement learning and
332 shaping,” in *International Conference on Machine Learning (ICML)*, 1998.
- 333 [27] H. Xu, X. Zhan, and X. Zhu, “Constraints penalized q-learning for safe offline reinforcement
334 learning,” in *AAAI Conference on Artificial Intelligence (AAAI)*, 2022.
- 335 [28] K. Driessens and S. Džeroski, “Integrating guidance into relational reinforcement learning,”
336 *Machine Learning*, 2004.
- 337 [29] Y. Song, Y.-b. Li, C.-h. Li, and G.-f. Zhang, “An efficient initialization approach of q-learning
338 for mobile robots,” *International Journal of Control, Automation and Systems*, 2012.
- 339 [30] P. Abbeel, A. Coates, and A. Y. Ng, “Autonomous helicopter aerobatics through apprenticeship
340 learning,” *The International Journal of Robotics Research*, 2010.
- 341 [31] J. Tang, A. Singh, N. Goehausen, and P. Abbeel, “Parameterized maneuver learning for au-
342 tonomous helicopter flight,” in *IEEE international conference on robotics and automation*
343 *(ICRA)*, 2010.
- 344 [32] D. Yu, H. Ma, S. Li, and J. Chen, “Reachability constrained reinforcement learning,” in
345 *International Conference on Machine Learning (ICML)*. PMLR, 2022.
- 346 [33] A. K. Jayant and S. Bhatnagar, “Model-based safe deep reinforcement learning via a constrained
347 proximal policy optimization algorithm,” *Advances in Neural Information Processing Systems*
348 *(NeurIPS)*, pp. 24 432–24 445, 2022.
- 349 [34] E. Altman, *Constrained Markov decision processes: stochastic modeling*. Routledge, 1999.
- 350 [35] M. Gutmann and A. Hyvärinen, “Noise-contrastive estimation: A new estimation principle for
351 unnormalized statistical models,” in *Proceedings of the international conference on artificial*
352 *intelligence and statistics (AISTATS)*, 2010.
- 353 [36] A. Mnih and Y. W. Teh, “A fast and simple algorithm for training neural probabilistic language
354 models,” in *International conference on machine learning (ICML)*, 2012.
- 355 [37] D. P. Bertsekas, “Nonlinear programming,” *Journal of the Operational Research Society*, 1997.