

# Safe Reinforcement Learning with Dead-Ends Avoidance and Recovery

Xiao Zhang, Hai Zhang, Hongtu Zhou, Chang Huang, Di Zhang, Chen Ye\*, Junqiao Zhao\*, *Member, IEEE*,

**Abstract**—Safety is one of the main challenges in applying reinforcement learning to realistic environmental tasks. To ensure safety during and after training process, existing methods tend to adopt overly conservative policy to avoid unsafe situations. However, overly conservative policy severely hinders the exploration, and makes the algorithms substantially less rewarding. In this paper, we propose a method to construct a boundary that discriminates safe and unsafe states. The boundary we construct is equivalent to distinguishing dead-end states, indicating the maximum extent to which safe exploration is guaranteed, and thus has minimum limitation on exploration. Similar to Recovery Reinforcement Learning, we utilize a decoupled RL framework to learn two policies, (1) a task policy that only considers improving the task performance, and (2) a recovery policy that maximizes safety. The recovery policy and a corresponding safety critic are pretrained on an offline dataset, in which the safety critic evaluates upper bound of safety in each state as awareness of environmental safety for the agent. During online training, a behavior correction mechanism is adopted, ensuring the agent to interact with the environment using safe actions only. Finally, experiments of continuous control tasks demonstrate that our approach has better task performance with less safety violations than state-of-the-art algorithms.

## I. INTRODUCTION

Reinforcement learning (RL) has made impressive achievements in long-term control tasks, including Atria games [1], car driving [2] and robot controlling [3]. While RL performs well in games and simulation environment, safety becomes one of the greatest challenge when applying RL to real-world task. In real environments such as autonomous driving tasks, unsafe actions can lead to damage to the agent itself and to the environment, resulting in significant maintenance costs and even human casualties.

To learn a safe policy that satisfies state-wise safety constraint in “safe critical” task, the agent needs to evaluate the safety of each state and avoid entering unsafe states [4]. In some Safe RL algorithms, the agent’s awareness of safety of a state is achieved by the safety critic that evaluates the safety of the task policy in states. The safety critic and a safety threshold together construct a boundary that divides the state space into safe and unsafe subspaces, as shown in Figure 1. The division of the state space depends on the policy, and sub-optimal policies will lead to more states being considered as unsafe, thus limiting agent exploration.

In [5], a recovery policy and a behavior correction mechanism are introduced. The task policy and the recovery policy

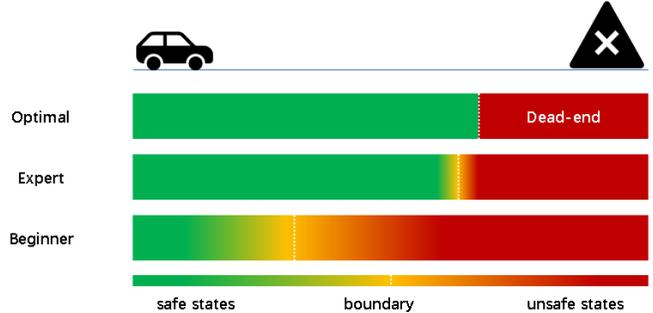


Fig. 1. The agent needs to control the brake of the car to avoid a collision. A safety critic evaluates the safety of the policy in all states and combines it with a threshold value to obtain the boundary that divides the safe and unsafe states. At a given initial speed, no policy can avoid a collision when the car is close enough to an obstacle, and these states are called dead-ends. The boundary that will and only identifies all dead-ends as unsafe states is called the optimal boundary. Since the optimal safe policy is safe in all states except dead-ends, by evaluating the safety of the optimal safe policy, it is possible to distinguish whether a state is a dead-end or not. In contrast, suboptimal safe policies can lead the safety critic to conservatively consider more states as unsafe.

are trained simultaneously to improve task performance and to satisfy safety constraints respectively. However, [5] inappropriately consider the safety of states as the safety of the task policy, which leads to a conservative agent since the agent is prone to misjudge under-explored states as unsafe. This greatly limits exploration, which in turn leads to inadequate collection of trajectories to correct the safety critic.

An example is given to illustrate the problem. A beginner learns to drive. Due to his poor driving skills, he ends up driving off the road several times on narrow roads, which led him to believe that driving on narrow roads is a very dangerous behavior and prevents him from continuing to explore in this environment. To solve this problem, an expert who can drive the car back on the road before the imminent danger occurs should be introduced. Thus, the beginner only needs to focus on the reward cumulating since the expert can take over before the imminent safety violation. In this way, the poor skill of the beginner will no longer prevent risky exploration.

Inspired by this, we propose our safe RL framework *Safe Reinforcement Learning with Dead Ends Avoidance and Recovery* (DEA-RRL), following the Recovery RL framework [5] and decoupled RL framework [6]. A recovery policy and a safety critic are trained with the goal of maximizing safety. When the task policy generates an action, the safety critic judges the safety of the action and decides whether to correct

\*This work is supported by the National Key Research and Development Program of China (No. 2021YFB2501104, No. 2020YFA0711402)

All the authors are with the Department of Computer Science and Technology, Tongji University, China and the MOE Key Lab of Embedded System and Service Computing, Tongji University, China, e-mail: (zhaojunqiao@tongji.edu.cn)

the action before interacting with the environment. Compared with RRL [5], the safety critic in our approach measures the safety of the recovery policy rather than the task policy, which reduces the conservatism and relaxes the constraint on exploration without increasing the risk of violations.

The advantages of our approach are as follows:

- Our method allows task policy to fully explore the environment and significantly improve task performance in complex environments;
- We show theoretically that our approach can learn optimal boundary, which is equivalent to the discovery of dead-ends.
- Our method completely decouples the task policy and safe policy, so that the recovery policy can be plug-and-play within other task policies without fine-tuning.

## II. RELATED WORK

### A. Safe RL

Safe RL addresses two main safety-related problems: asymptotic safety and in-training safety of policies. Asymptotic safety means the safety of policies after convergence, which is commonly achieved by reward penalties and cost constraints. Lagrange Relaxation [7] is the most widely used method due to its simplicity, and other methods such as Trust Regions [8], Lyapunov-based [9], Guide Policy [10] are proposed for their stricter safety guarantees, fewer violations during the training and having faster convergence.

Although these approaches do find safe policies after training, they learn safety by trial-and-error as the same as traditional RL algorithms, which means violations are inevitable before convergence. To ensure in-training safety, state-wise safety constraints and prior knowledge of the environment are utilized [4]. [11] [12] [13] assume a white-box or a black-box environment dynamics model is available, limiting policy optimizing in a safe policy set which is constructed based on the known model. However, these assumptions are hard to meet in real-world problems.

RRL [5] trains a recovery policy using offline data to guarantee the safety of the online training process. [14] combines RRL and Meta-RL [15], improving the generalizability of RRL by enabling pre-trained safety critic and recovery policy to be quickly adapted to different tasks during the fine-tuning phase. Although these methods can achieve in-training safety, they are prone to obtain over-conservative policies due to the reason explained in Section I.

### B. Decoupling Performance and Safety

Decoupling performance and safety offers new ideas for solving the exploration-exploitation dilemma [6], [16], because separating processes of maximizing reward and guaranteeing safety prevents policies from under-performing due to over-conservatism. [17] and [5] introduce an implicit and an explicit safe policy respectively. The task policy no longer needs to consider safety explicitly when being updated. However, the identification of unsafe actions in these methods strongly correlated with the task policy, thus the agent is still prevented from obtaining optimal policies

by exploration-exploitation dilemma. Although [18] achieves full decoupling of exploration and exploitation, it focuses only on the rewards and ignores safety in training.

### C. Dead-Ends Discovery and Avoidance

The concept of dead-ends discovery (DeD) was introduced in [19] and later applied to the medical field [20]. [21] combines risk sensitivity with DeD model, which allows dead-ends to be identified earlier. [22] also proposes “irrecoverable state” with a similar meaning to dead-ends state, preventing dangerous situations from occurring through reward shaping and model-based rollout [23].

We argue that distinguishing dead-ends states from normal states is crucial for improving the performance of safe RL, as it identifies the broadest range of policies that can be explored safely. The safety critic in RRL [5] could not distinguish dead-ends states, as explained later in Section IV-A.

## III. PRELIMINARY

### A. Constraint Markov Decision Processes

We consider safe RL under Constraint Markov Decision Process (CMDPs)  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma, \rho, \mathcal{C}, \gamma_{safe}, \epsilon_{safe})$  [24], where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  the action space,  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  the reward function,  $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  the transition function,  $\gamma \in [0, 1)$  the discount factor for reward,  $\mathcal{C} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  the cost function,  $\gamma_{safe} \in [0, 1)$  the discount factor for cost and  $\epsilon_{safe} \in \mathcal{R}$  the safe threshold. Let  $\Pi$  be the set of Markovian stationary policies. Given policy  $\pi \in \Pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$  maps states to action distributions and  $\pi(a|s)$  denotes the probability of choosing action  $a$  in state  $s$ . The task performance of  $\pi$  is defined as the discounted cumulative reward  $\mathcal{J}(\pi)$ :

$$\mathcal{J}(\pi) = \mathbb{E}_{\tau \sim \pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{R}(s_t, a_t, s_{t+1}) \right] \quad (1)$$

Similar to [25], state-cost function  $V_c$  and action-cost function  $Q_c$  are introduced to indicate the expected cumulative cost of  $\pi$  in  $s$ :

$$V_c^\pi(s) = \mathbb{E}_{\tau \sim \pi, P} \left[ \sum_{t=0}^{\infty} \gamma_{safe}^t \mathcal{C}(s_t, a_t, s_{t+1}) | s_0 = s \right] \quad (2)$$

$$Q_c^\pi(s, a) = \mathbb{E}_{s' \sim P(\cdot | s, a)} V_c^\pi(s') \quad (3)$$

Under the state-wise safety constraint formulation [4], we define the set of safe policies

$$\Pi_c = \{ \pi \in \Pi | \forall s \in \mathcal{S}, V_c^\pi(s) < \epsilon_{safe} \} \quad (4)$$

The objective of CMDPs is to find a policy that maximizes Equation (1) in the set of safe policies  $\Pi_c$ :

$$\pi_{task}^* = \max_{\pi} \mathcal{J}(\pi), \text{ s.t. } \pi \in \Pi_c \quad (5)$$

## B. Safe Markov Decision Processes

In safety-critical tasks, any unsafe action is fatal. Therefore, we define SMDPs, a special case of CMDPs to describe this kind of tasks.

*Definition 1:* Safe Markov Decision Processes (SMDPs)  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma, \rho, \mathcal{C}, \gamma_{safe}, \epsilon_{safe})$ , a special case of CMDPs where episodes terminate after any danger occurred. In SMDPs, the cost function is a binary indicator of the safety of the state:

$$\mathcal{C}(s_t, a_t, s_{t+1}) = \begin{cases} 1, & \text{safety violation} \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

Observe that if  $\gamma_{safe} = 1$ ,  $Q_c^\pi$  indicates the probability of ending up with a failure state in the future with  $\pi$ .  $Q_{\phi,c}^\pi$ , parameterized by  $\phi$ , can be optimized by minimizing the MSE loss:

$$\mathcal{L}_{Q_c}(\phi; \pi) = \frac{1}{2} (Q_{\phi,c}^\pi(s_t, a_t) - (c_t + (1 - c_t) \gamma_{safe} \mathbb{E}_{a_{t+1} \sim \pi(\cdot|s_{t+1})} Q_{\phi,c}^\pi(s_{t+1}, a_{t+1})))^2 \quad (7)$$

where  $(s_t, a_t, s_{t+1}, c_t)$  is the transition sampled from offline data or replay buffer and  $c_t$  is short for  $\mathcal{C}(s_t, a_t, s_{t+1})$ . Equation (7) can be considered to be the policy evaluation of  $\pi$  in terms of safety.

*Definition 2:* The state space is divided into three sub-spaces:

- $\mathcal{S}_{fail}$ : Failure state. Indicates the end of an episode due to a danger.
- $\mathcal{S}_{dead}$ : Dead-ends state. Any dead-ends state will transform into a failure state, regardless of the policy  $\pi$  the agent takes ( $\pi \in \Pi$ ).
- $\mathcal{S}_{safe}$ : Safe state.  $\mathcal{S}_{safe} = \mathcal{S} \setminus (\mathcal{S}_{fail} \cup \mathcal{S}_{dead})$ .

Therefore,

$$\forall s \in \mathcal{S}_{safe}, \exists (a, s') \in (\mathcal{A}, \mathcal{S}_{safe}), P(s, a, s') > 0 \quad (8)$$

The cost function in SMDPs can also be expressed as

$$\mathcal{C}(s_t, a_t, s_{t+1}) = \mathbb{I}(s_{t+1} \in \mathcal{S}_{fail}) \quad (9)$$

Similar to [22], we assume that a safety violation must come fairly soon after entering any dead-ends states region:

*Assumption 1:* There exists a horizon  $H \in \mathbb{N}$ , that any trajectory starting from  $s_0 \in \mathcal{S}_{dead}$  will end up in  $H$  steps.

We additionally introduce a sampling policy for training the safety critic  $\bar{\pi}$  and modify the definition of the set of safe policies as follows:

$$\Pi_{\bar{\pi}} = \{\pi \in \Pi | \forall (s, a) \in (\mathcal{S}_{safe}, \mathcal{A}) \text{ and } \pi(a|s) > 0, Q_c^\pi(s, \pi(s)) < \epsilon_{safe}\} \quad (10)$$

Similar to CMDPs, the objective of SMDPs is to find a policy with Equation (1) in the set of safe policies  $\Pi_{\bar{\pi}}$ :

$$\pi_{task}^* = \max_{\pi} \mathcal{J}(\pi), \text{ s.t. } \pi \in \Pi_{\bar{\pi}} \quad (11)$$

## IV. METHODS

### A. Recovery RL

In most approaches of Safe RL without decoupling,  $\bar{\pi}$  is the same as  $\pi_{task}$ , including RRL [5]. In pretrain phase, RRL trains safety critic  $Q_{\phi,c}^{\pi_{task}}$  and recovery policy  $\pi_{\theta,rec}$  (parameterized by  $\theta$ ) by minimizing  $\mathcal{L}_{Q_c}(\phi; \pi_{task})$  and maximizing  $\mathcal{J}_{\pi_{rec}}(\theta; \pi_{task})$  in Equation (7) and Equation (12).

$$\mathcal{J}_{\pi_{rec}}(\theta; \bar{\pi}) = -\mathbb{E}_{s \sim \mathcal{D}} [Q_c^{\bar{\pi}}(s, \pi_{\theta,rec}(\cdot|s))] \quad (12)$$

In fine-tune phase, unsafe actions will be corrected by  $\pi_{rec}$ :

$$a_t = \begin{cases} a^{\pi_{task}}, & Q_{\phi,c}^{\bar{\pi}}(s_t, a^{\pi_{task}}) < \epsilon_{safe} \\ a^{\pi_{rec}}, & \text{otherwise} \end{cases} \quad (13)$$

Similar to Equation (11) The objective of RRL can be expressed as:

$$\pi_{task}^* = \max_{\pi} \mathcal{J}(\pi), \text{ s.t. } \pi \in \Pi_c^\pi \quad (14)$$

Because  $\pi_{task}$  has not been trained during pretrain phase that it is unable to avoid unsafe situations, resulting in an overestimation of  $Q_{\phi,c}^{\pi_{task}}$  which represents a conservative estimate of safety.

### B. Dead-ends Discovery and Avoidance

Different from RRL, our proposed DEA-RRL ensure safety by distinguishing between dead-ends states and safe states and only prevent the agent from entering dead-ends. In the pretrain phase, DEA-RRL trains safety critic  $Q_{\phi,c}^{\pi_{rec}}$  and recovery policy  $\pi_{\theta,rec}$  by minimizing  $\mathcal{L}_{Q_c}(\phi; \pi_{rec})$  and  $\mathcal{J}_{\pi_{rec}}(\theta; \pi_{rec})$  as in Equation (7) and Equation (12), which is completely decoupled from  $\pi_{task}$ . This is equivalent to solving the optimal Bellman equation [25] for safety. Thus, the optimal recovery policy and the corresponding policy evaluation function can be obtained.

$$\begin{aligned} \pi_{rec}^*(\theta) &= \min_{\theta} \mathbb{E}_{\tau \sim P, \pi_{rec}(\theta)} [Q_c^{\pi_{rec}}(s, \pi_{rec}(s; \theta))] \\ &= \min_{\theta} \mathbb{E}_{\tau \sim P, \pi_{rec}(\theta)} \left[ \sum_{t=0}^{\infty} \gamma_{safe}^t c_t \right] \end{aligned} \quad (15)$$

By using the same behavior correction mechanism as Equation (13) where  $\bar{\pi}$  is  $\pi_{rec}^*$ , the object of DEA-RRL can be expressed as:

$$\pi_{task}^* = \max_{\pi} \mathcal{J}(\pi), \text{ s.t. } \pi \in \Pi_c^{\pi_{rec}^*} \quad (16)$$

### C. Theoretical Proof

We will illustrate the advantages of DEA-RRL over RRL theoretically.

*Theorem 1:*  $\Pi_c^{\pi_{task}}$  and  $\Pi_c^{\pi_{rec}^*}$  are the accessible space for  $\pi$  to explore safely in RRL and DEA-RRL respectively, we have  $\Pi_c^{\pi_{task}} \subseteq \Pi_c^{\pi_{rec}^*}$ .

*Proof:* Since  $Q_c^*$  (the shorthand of  $Q_c^{\pi_{rec}^*}$ ) is the optimal cost value function, we have  $Q_c^*(s, a) \leq Q_c^{\pi_{task}}(s, a)$  for every  $(s, a)$ . As in Equation (10),  $\forall (\pi, s, a) \in (\Pi_c^{\pi_{task}}, \mathcal{S}_{safe}, \mathcal{A})$  and  $\pi(a|s) > 0$ ,

$$Q_c^*(s, \pi(s)) \leq Q_c^{\pi_{task}}(s, \pi(s)) < \epsilon_{safe}$$

therefore  $\pi \in \Pi_c^{\pi_{rec}}$ , which means  $\Pi_c^{\pi_{task}} \subseteq \Pi_c^{\pi_{rec}}$ .

We will show that with appropriate  $\epsilon_{safe}$ ,  $Q_c^*$  can be used to identify dead-ends states.

*Lemma 1:* Suppose that Assumption 1 holds and uncertainties are ignored in the environment i.e.  $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \{0, 1\}$ ,

- $\forall (s, \pi) \in (\mathcal{S}_{fail}, \Pi)$ ,  $V_c^\pi(s) = V_c^*(s) = 1$ ,
- $\forall (s, \pi) \in (\mathcal{S}_{dead}, \Pi)$ ,  $V_c^\pi(s) \geq V_c^*(s) \geq \gamma_{safe}^{H-1}$
- $\forall (s, \pi) \in (\mathcal{S}_{safe}, \Pi)$ ,  $V_c^\pi(s) \geq V_c^*(s) = 0$

*Proof:* Since  $V_c^*$  is the optimal value function, we have  $V_c^\pi(s) \geq V_c^*(s)$  for all  $s \in \mathcal{S}$ . By the definition of cost function Equation (9) and state-cost function Equation (2),  $\forall s \in \mathcal{S}_{fail}$ ,  $V_c^\pi(s) = 1$ . Assumption 1 shows that episode would terminate after at most  $H$  steps since the agent reach dead-ends state, we can derive from Equation (2) that

$$\begin{aligned} V_c^\pi(s) &= \mathbb{E}_{\tau \sim \pi, P} \left[ \sum_{t=0}^{\infty} \gamma_{safe}^t \mathcal{C}(s_t, a_t, s_{t+1}) \mid s_0 = s \in \mathcal{S}_{dead} \right] \\ &\geq \sum_{t=0}^{H-2} \gamma_{safe}^t * 0 + \gamma_{safe}^{H-1} * 1 = \gamma_{safe}^{H-1}. \end{aligned} \quad (17)$$

By Equation (8) in Definition 2, for every  $s \in \mathcal{S}_{safe}$ , there always exists at least one action  $a$  such that  $s' \sim P(\cdot | s, a)$ ,  $s' \in \mathcal{S}_{safe}$ . Start at  $s \in \mathcal{S}_{safe}$ , agent would never reach dead-ends state by choosing the safe action, which means that  $V_c^*(s) = 0$ .

*Theorem 2:* Suppose that Assumption 1 holds and uncertainties are ignored in the environment, and let

$$\epsilon_{safe} = \gamma_{safe}^H \quad (18)$$

Then, with behavior correction mechanism showed by Equation (13), the agent will be prevented from reaching dead-ends states.

*Proof:* According to Equation (13), the actions allowed to be performed satisfy:

$$Q_c^\pi(s, a) = \mathbb{E}_{s' \sim P(\cdot | s, a)} \gamma_{safe} V_c^\pi(s') < \gamma_{safe}^H \quad (19)$$

therefore

$$V_c^\pi(s') < \gamma_{safe}^{H-1}, \forall s \in \mathcal{S}_{safe} \quad (20)$$

which ensuring  $s' \in \mathcal{S}_{safe}$ . Also because of the initial state  $s_0 \in \mathcal{S}_{safe}$ , it is ensured that the agent is always explored in safe states.

Theorem 1 shows that DEA-RRL provides larger accessible space for exploration than RRL, reducing conservatism without loss of safety. Lemma 1 and theorem 2 indicate that starting from safe states, both RRL and DEA-RRL are able to prevent agent entering dead-ends states with behavior correction mechanism and appropriate  $\epsilon_{safe}$ . Notice that in RRL, it is not ensured that  $V_c^{\pi_{task}}(s) < \gamma_{safe}^{H-1}$ ,  $\forall s \in \mathcal{S}_{safe}$ , thus RRL will recognize some safe states as dead-ends states due to its over-conservatism on under-explored states.

#### D. Offline Pretrain

The safety critic in RRL is entirely determined by  $\pi_{task}$ , thus  $\pi_{rec}$  is equivalent to a single-step greedy policy that

chooses the action with the smallest  $Q_c^{\pi_{task}}$  in a state. This is referred to as single-step dynamic programming (SSDP) [26]. By contrast, safety critic in DEA-RRL is determined by  $\pi_{rec}$ , giving  $\pi_{rec}$  the ability to combine multiple sub-optimal trajectories into one optimal trajectory, hence it is referred to as multi-steps dynamic programming (MSDP) [26].

The bias in MSDP due to querying the  $Q_c^\pi$  value of out-of-distribution (OOD) actions accumulates over the DP process, making MSDP more sensitive to OOD actions compared with SSDP. In the training of safety critic, we prevent safety critic from overestimating the safety of an action by avoiding querying the  $Q_c$  value of OOD actions. While in the training of recovery policy, we can not completely ignore OOD actions especially in states where the offline data do not contain safe actions, for it is better to try an unknown and possibly safe action than to choose a known but certainly dangerous action.

Inspired by implicit Q-learning (IQL) algorithm [26], we use Expectile Regression to train  $Q_c^{\pi_{rec}}$ , avoiding the impact of OOD actions. Different from IQL, we use Advantage Policy Gradient instead of Advantage Weighted Regression used in [26] to train  $\pi_{rec}$ , allowing an OOD action to be attempted in a state where all known actions are unsafe.

The state-cost function  $V_c^\pi$  and action-cost function  $Q_c^\pi$  is updated by minimizing following loss functions:

$$\mathcal{L}_{V_c^\pi}(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^\tau(Q_{\phi,c}^\pi(s, a) - V_{\psi,c}^\pi(s))] \quad (21)$$

$$\begin{aligned} \mathcal{L}_{Q_c^\pi}(\phi) &= \mathbb{E}_{(s,a,s') \sim \mathcal{D}} [(c(s, a, s') + \\ &(1 - c(s, a, s')) \gamma_{safe} V_{\psi,c}^\pi(s') - Q_{\phi,c}^\pi(s, a))^2] \end{aligned} \quad (22)$$

where  $L_2^\tau$  is the expectile regression loss and  $\mathcal{D}$  the offline dataset.

As illustrated in Algorithm 1, we provide the agent with  $\mathcal{D}$  that contain several transitions from both safe trajectories and unsafe trajectories. Safety critic and recovery are trained by minimizing corresponding objective functions without task policy anymore.

---

#### Algorithm 1 DEA-RRL Pretrain Offline

---

- 1: Input: offline dataset  $\mathcal{D}$ , pretraining steps  $N$
  - 2: **for**  $steps, \leftarrow 1, N$  **do**
  - 3:   Sample a mini-batch  $(s_t, a_t, s_{t+1}, c_t)$  from  $\mathcal{D}$
  - 4:   Update  $\phi$  by minimizing  $\mathcal{L}_{Q_c^\pi}(\phi)$  (Equation (22))
  - 5:   Update  $\psi$  by minimizing  $\mathcal{L}_{V_c^\pi}(\psi)$  (Equation (21))
  - 6:   Update  $\theta$  by maximizing  $\mathcal{J}_{\pi_{rec}}(\theta)$  (Equation (12))
  - 7: **end for**
- 

#### E. Online Training

Any of the RL algorithms can be used to train  $\pi_{task}$ . In our work, we utilize Soft Actor Critic algorithm (SAC) [27]. The process of online fine-tuning is illustrated in Algorithm 2, and we give some remarks on online training.

- We relabel all actions with the action proposed by  $\pi_{task}$  as the same as [5], which is important to achieve decoupled safe reinforcement learning, as it prompts

the agent to view behavior correction as part of the environment.

- Since  $\pi_{rec}$  ignores task performance, unsafe actions result in low reward by behavior correction, which enables  $\pi_{task}$  to learn to avoid unsafe actions with reward feedback only.
- $\pi_{rec}$  is a Gaussian policy, we use the action mean rather than sampling over the distribution during behavior correction for safety.
- We choose not to fine-tune  $Q_c^{\pi_{rec}}$  and  $\pi_{rec}$  in online training, because a fixed behavior correction strategy provides the agent a stable MDP dynamic model, thus improving the stability of training.
- Finally, we use a small  $\epsilon_{safe}$  to improve the safety of the algorithm due to aleatoric uncertainty and epistemic uncertainty,

---

#### Algorithm 2 DEA-RRL Training Online

---

```

1: Input: safety critic  $Q_c^{\pi_{rec}}$ , recovery policy  $\pi_{rec}$ , task
   horizon  $H$ , training steps  $N$ 
2: Initialize replay buffer  $\mathcal{D}_{task} \leftarrow \emptyset$ 
3:  $s_0 \leftarrow env.reset()$ 
4: for  $steps, \leftarrow 1, N$  do
5:   for  $t \in \{1, \dots, H\}$  do
6:     if  $c_{t-1} = 1$  or  $t = H$  then
7:        $s_t \leftarrow env.reset()$ 
8:     end if
9:     Sample  $a^{\pi_{task}}, a^{\pi_{rec}}$  from  $\pi_{task}, \pi_{rec}$ 
10:    if  $Q_c^{\pi_{rec}}(s_t, a^{\pi_{task}}) < \epsilon_{safe}$  then
11:       $a_t = a^{\pi_{task}}$ 
12:    else
13:       $a_t = a^{\pi_{rec}}$  (behavior correction)
14:    end if
15:    Execute  $a_t$ , Observe  $s_{t+1}, r_t, c_t$ 
16:     $\mathcal{D}_{task} \leftarrow \mathcal{D}_{task} \cup (s_t, a^{\pi_{task}}, s_{t+1}, r_t)$ 
17:    Train  $\pi_{task}$  on  $\mathcal{D}_{task}$  by maximizing  $\mathcal{J}(\pi)$ 
18:  end for
19: end for

```

---

## V. EXPERIMENTS

In the experiments, we investigate whether our approach can:

- exceed state-of-the-art algorithms in terms of task performance, post-training safety and in-training safety;
- improve the safety of the training process with minimal impact on task performance;
- obtain a task-independent behavior correction strategy that can ensure safety of other trained policies in testing.

### A. Domains

Experiments were conducted under the standard safety reinforcement learning test environment Safety Gym [28]. As shown in Figure 2, We selected the two simple environments (StaticEnv, DynamicEnv) set up in [29] and three more complex environments (PointGoal1, CarGoal1, DoggoGoal1) pre-defined in Safety Gym. Unlike the general Safety Gym

setup, we require that any violation of safety constraints will result in immediate termination of the episode, which corresponds to real world tasks. This setup significantly increases the difficulty of the task and poses an extreme challenge to the agent’s capacity to balance exploration and exploitation.

### B. Offline Data Collection

Unlike general offline RL where the training results are influenced by the performance of the behavior policy that used for data sampling, Recovery RL framework requires offline data to contain a wide coverage of unsafe trajectories, allowing  $Q_c$  to learn to identify unsafe actions efficiently. The offline dataset we used contains 2M transitions including: (1) 1M transitions sampled from the replay buffer of SAC training, and (2) 1M transitions obtained by interacting with the environment using random actions. The latter is included because in our experiments we find that the failed transitions in the SAC replay buffer tend to share similar characteristics. The random sampling data can diversify the failed transitions in the offline data. Inspired by [30], we filter out the parts of the data that are less relevant to safety by keeping only 100 transitions before violations.

### C. Evaluation Metric

The average cumulative reward (ACR) over episodes is used as a measure of the task performance, and a higher return indicates a better task performance. The default reward functions provided by Safety Gym are used in our experiments, which define algorithm-independent tasks. We use the average rate of constraint violation (AVR) in testing as a measure of asymptotic safety and the number of constraint violations (TV) in training as a measure that indicating in-training safety. To highlight the differences between our method and RRL, we use the ratio of steps that behavior correction mechanism is used (ARR) as a measure of the extent to which the behavior correction mechanism intervene in training and testing phase.

### D. Comparisons with Baselines

- Unconstrained Baseline [27]: We use SAC as a baseline for unconstrained methods, optimizing task performance and ignoring safety constraints.
- Worst-Case Soft Actor Critic (WCSAC) [29]: A combined Lagrangian relaxation and risk-sensitive approach that maximizing

$$\mathcal{J}(\pi) - \lambda(\mathbb{E}_{\tau \sim P, \pi}[CVaR_\alpha(Q_c^{\pi_{task}}(s, a))] - \epsilon_{safe})$$

where  $CVaR_\alpha \doteq \mathbb{E}_{p^\pi}[Q_c^\pi | Q_c^\pi \geq F_C^{-1}(1 - \alpha)]$  and  $F_C$  is the CDF of  $p^\pi(Q_c^\pi | s, a)$ , updating policy parameters and  $\lambda$  via dual gradient descent.

- RRL [5]: A *semi-decoupled* approach that preventing agent using unsafe action for which  $Q_c^{\pi_{task}}(s, a) \geq \epsilon_{safe}$  by behavior correction mechanism.
- Implicit Q-Learning (IQL) [26]: An offline RL algorithm that optimizing critic and actor by expectile regression and advantage weighted regression.

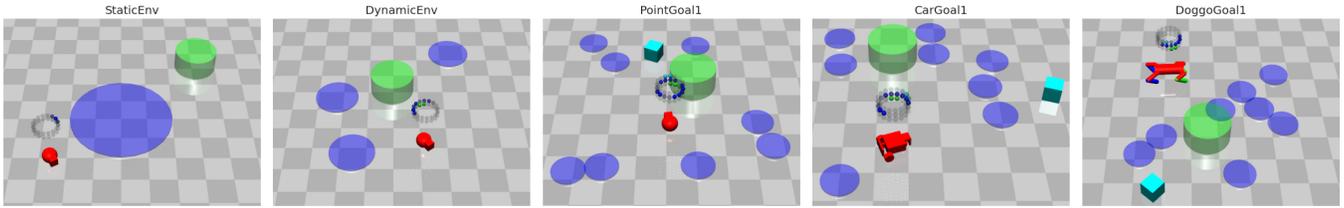


Fig. 2. The experiments in this paper were conducted in the safety gym standard test environment. The figure shows five task environments from left to right, including StaticEnv, DynamicEnv, PointGoal1, CarGoal1 and DoggoGoal1. These environments are based on Safety Gym, where the task is to navigate a Point robot from its initial position to the target area (green area) and avoid entering unsafe areas (blue areas) during the process. PointGoal1, CarGoal1 and DoggoGoal1 can be seen as upgraded versions of DynamicEnv, which include more unsafe areas, larger action space and more complex robots.

To better compare the ability of each method to balance safety and task performance, we let each method has similar asymptotic safety achieving by using different  $\epsilon_{safe}$ .

### E. Results

1) *Main Results:* The performance and safety of all methods are showed in Figure 3 and Table I. The results show that our method has high asymptotic safety and in-training safety in all tasks and similar cumulative reward to the unconstrained method in most tasks, suggesting that our method can substantially improve the safety of the algorithm with less impact on performance.

WCSAC obtains a policy with higher safety than SAC after training for all tasks, however, since WCSAC can only learn information related to the safety of the environment during training, a large number of constraint violations are inevitable in the early stages of training.

IQL avoids interaction with the environment by training completely offline so that safety constraints are not violated during training, but a safe policy cannot be obtained after training.

RRL learns information about the safety of the environment in advance through pre-training and therefore maintains a low number of constraint violations throughout the training process. As shown in Figure 3, the ratio of steps RRL using behavior correction during training process is much higher than that of DEA-RRL, indicating a large amount of intervention in the training process of  $\pi_{task}$ , which is an important reason why RRL can only learn sub-optimal policies in complex environments. The training curve of RRL has large variance, which is caused by the fact that  $\pi_{task}$  is training in an unstable environment. In contrast, DEA-RRL does not fine-tune  $\pi_{task}$  and  $Q_c^{\pi_{rec}}$  during online training, providing a stable MDPs for training  $\pi_{task}$ .

We also remark that DEA-RRL has lower return than SAC in DoggoGoal1 because the environment are more difficult to explore with the introduction of the behavior correction mechanism, and it takes longer for agents to learn the optimal policy.

2) *Ablations:* We design ablation experiments to investigate the effect of settings in pre-training and online training on safety and task performance.

We first compare the effects of using different training methods in offline training on the experimental results

which are shown in Figure 4. RRL (MSDP) is a multi-step dynamic programming version of RRL, significantly reducing conservatism of RRL. However, RRL (MSDP) is still highly conservative due to the OOD actions and has high recovery ratio and low cumulative reward. DEARRL (IQL) extract a policy from safety critic using advantage weighted regression. Since DEARRL (IQL) will only select known actions, it cannot guarantee safety in states where the dataset does not contain safe actions, resulting in an unsafe policy.

We then compared the performance and safety of RRL and DEA-RRL under different  $\epsilon_{safe}$ . As Figure 5 shows, as  $\epsilon_{safe}$  decreases, RRL and DEA-RRL have the same performance, i.e. the average reward and the proportion of constraint violations decrease. Notice that with equal  $\epsilon_{safe}$ , DEA-RRL has a similar safety and far superior task performance than RRL. DEA-RRL substantially reduces the influence of the behavior correction mechanism on  $\pi_{task}$ , and eventually enables algorithm to find a better  $\pi_{task}$ .

3) *Decoupled Framework:* In DEA-RRL,  $\pi_{rec}$  and  $Q_c^{\pi_{rec}}$  are completely decoupled from  $\pi_{task}$ , so that the pre-trained  $\pi_{rec}$  and  $Q_c^{\pi_{rec}}$  can be directly combined with other  $\pi_{task}$  trained to improve the safety of these algorithms in testing. We combine the behavior correction strategy obtained by RRL and DEA-RRL pre-training with the other policies obtained through SAC and IQL training and test them on 500 random episodes, the results of which are shown in Table II. The results show that DEA-RRL can be combined directly with any RL algorithm to substantially improve the safety of these algorithms in testing at the expense of slightly lower return. Notice that in some environment  $s$ ,  $\pi_{task}$  and  $\pi_{rec}$  obtained by IQL (with DEA-RRL) even achieve the performance and safety of the online training method when trained completely offline.

Finally, we find that although the average episode length in testing increases significantly with the introduction of the behavior correction mechanism compared to using SAC and IQL only, there is only a slight improvement in the return. The reason is that the introduction of the behavior correction mechanism actually changes the MDPs model so that policies trained using SAC and IQL that maximize the returns of the original MDPs model have lower performance under the new MDPs model.

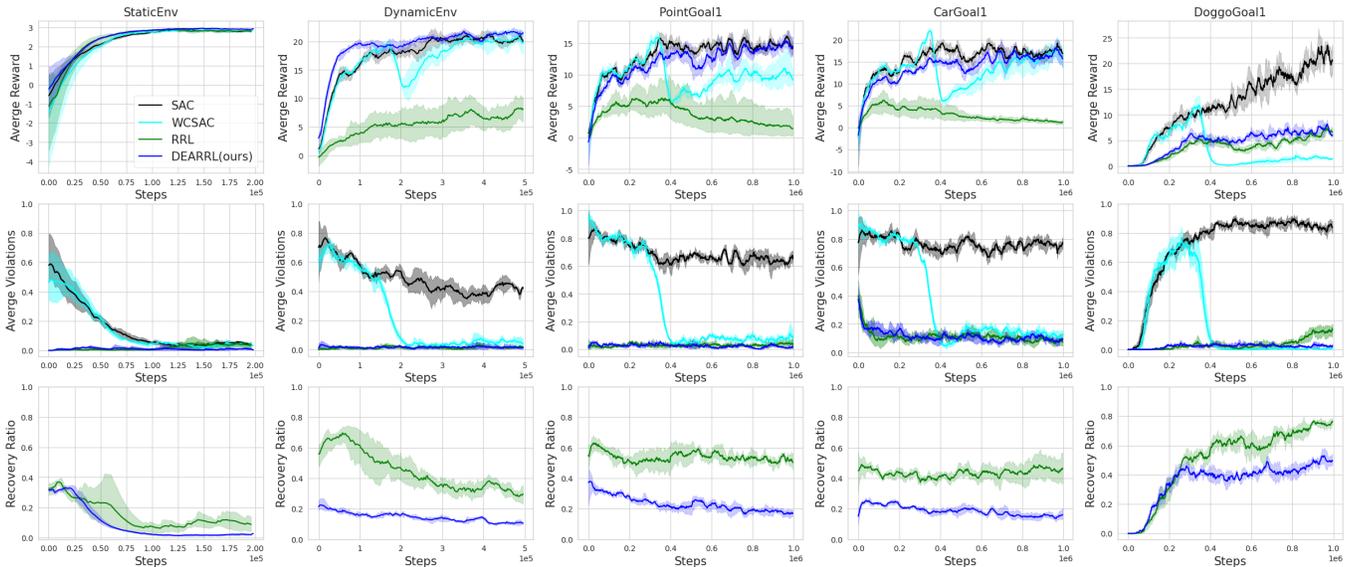


Fig. 3. Training Curve. Means (solid lines) and variances (shaded) of the training curves for the four algorithms under different tasks, with the first row showing the average cumulative reward, the second row showing the average proportion of constraint violations over update steps and the third row showing the average use of behavior correction. We adapted  $\epsilon_{safe}$  for RRL and DEA-RRL to make the two algorithms have similar safety in training. For each method we used random seeds for training.

TABLE I

| Environments | SAC    |       |       | WCSAC         |              |       | RRL   |          |            | DEA-RRL(ours) |              |            | IQL           |       |
|--------------|--------|-------|-------|---------------|--------------|-------|-------|----------|------------|---------------|--------------|------------|---------------|-------|
|              | ACR    | AVR   | TV    | ACR           | AVR          | TV    | ACR   | AVR      | TV         | ACR           | AVR          | TV         | ACR           | AVR   |
| StaticEnv    | 2.754  | 0.043 | 1642  | <b>2.936</b>  | <b>0.008</b> | 1453  | 2.779 | 0.038    | 311        | 2.882         | 0.023        | <b>119</b> | 2.630         | 0.095 |
| DynamicEnv   | 17.434 | 0.519 | 3330  | 19.004        | 0.069        | 1451  | 7.479 | <b>0</b> | 92         | <b>20.784</b> | 0.018        | <b>86</b>  | 18.044        | 0.565 |
| PointGoal1   | 12.609 | 0.753 | 12375 | 7.811         | 0.055        | 5739  | 1.155 | 0.078    | 408        | <b>14.079</b> | <b>0.035</b> | <b>236</b> | 12.740        | 0.76  |
| CarGoal1     | 16.727 | 0.73  | 14850 | <b>18.672</b> | <b>0.035</b> | 6199  | 1.827 | 0.043    | 1006       | 15.430        | 0.091        | <b>786</b> | 11.357        | 0.94  |
| DoggoGoal1   | 16.956 | 0.908 | 31450 | 1.055         | <b>0.013</b> | 10825 | 6.969 | 0.22     | <b>546</b> | 7.002         | 0.048        | 706        | <b>18.889</b> | 0.838 |

The table records the average cumulative reward (ACR) and the average violation rate (AVR) for the policy trained by, and the number of violations during training (TV) for the SAC, WCSAC, RRL, DEA-RRL and IQL. Since IQL is trained completely offline, the number of violations during training is 0. The ACR and AVR are obtained by testing  $\pi_{task}$  on 500 random episodes on average.

TABLE II

| Environments | SAC(with RRL) |       |       | IQL(with RRL) |       |       | SAC(with DEA-RRL) |       |       | IQL(with DEA-RRL) |       |       |
|--------------|---------------|-------|-------|---------------|-------|-------|-------------------|-------|-------|-------------------|-------|-------|
|              | ACR           | AVR   | ARR   | ACR           | AVR   | ARR   | ACR               | AVR   | ARR   | ACR               | AVR   | ARR   |
| StaticEnv    | 2.379         | 0     | 0.417 | 1.308         | 0.005 | 0.426 | 2.673             | 0     | 0.009 | <b>2.947</b>      | 0     | 0.031 |
| DynamicEnv   | 1.928         | 0.004 | 0.671 | 1.488         | 0.008 | 0.764 | <b>18.769</b>     | 0.032 | 0.242 | <b>19.425</b>     | 0.078 | 0.235 |
| PointGoal1   | 5.462         | 0.116 | 0.546 | 3.669         | 0.064 | 0.581 | <b>10.596</b>     | 0     | 0.275 | <b>9.926</b>      | 0.04  | 0.322 |
| CarGoal1     | 0.093         | 0.034 | 0.614 | 0.144         | 0.094 | 0.602 | 8.165             | 0.028 | 0.334 | <b>10.783</b>     | 0.116 | 0.297 |
| DoggoGoal1   | 0.751         | 0     | 0.987 | 0.656         | 0     | 0.941 | <b>13.285</b>     | 0.002 | 0.780 | 9.407             | 0     | 0.846 |

We compare the average cumulative reward (ACR) and average violation rate (AVR) of the algorithms by directly combining a pre-trained behavior correction mechanism with  $\pi_{task}$  trained using SAC and IQL, where the parameters of RRL and DEA-RRL and  $\epsilon_{safe}$  are set as the same as in the previous experiment. All data are derived from the average of 500 random episodes of test results.

## VI. CONCLUSION

In this paper, we propose DEA-RRL, a fully decoupled safe reinforcement learning framework. The decoupling of task performance and safety is achieved by pre-training the recovery policies that maximize safety. The task policies are free to explore without considering safety at all, and safety in exploration is achieved by combining recovery policies with behavior correction. We show that, our approach is able to identify the dead-end states in determined MDPs,

which is the maximum range that a policy can explore safely. Using appropriate safety thresholds to prevent policies from going into dead-ends can ensure safety without compromising reward improvement. We demonstrate in a series of experiments that our approach balances task performance and safety well, suggesting that our approach has potential to be used in realistic environmental tasks.

Future work will consider the following aspects: (1) introducing model-based reinforcement learning to estimate uncertainty, allowing  $\epsilon_{safe}$  to be adaptive, and (2) obtaining

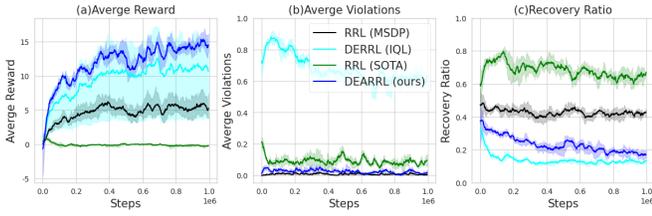


Fig. 4. Ablations of offline pretraining. Average reward, average proportion of constraint violations and average proportion of use of behavioral correction for different methods. For each method we used same  $\epsilon_{safe} = 0.7$  and random seeds for training.

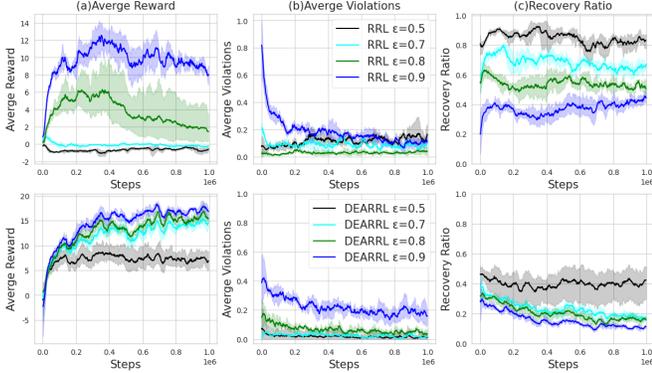


Fig. 5. Ablations of online training. Average reward, average proportion of constraint violations and average proportion of use of behavioral correction for different  $\epsilon_{safe}$  choices for RRL and DEARRL in PointGoal1. For each method we used random seeds for training.

an offline safe reinforcement learning algorithm that involves behavior correction mechanism in offline training process of  $\pi_{task}$  to improve task performance.

## REFERENCES

- [1] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al., “Human-level control through deep reinforcement learning,” *nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [2] A. Kendall, J. Hawke, D. Janz, P. Mazur, D. Reda, J.-M. Allen, V.-D. Lam, A. Bewley, and A. Shah, “Learning to drive in a day,” in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 8248–8254.
- [3] C. Bodnar, A. Li, K. Hausman, P. Pastor, and M. Kalakrishnan, “Quantile qt-opt for risk-aware vision-based robotic grasping,” *arXiv preprint arXiv:1910.02787*, 2019.
- [4] W. Zhao, T. He, R. Chen, T. Wei, and C. Liu, “State-wise safe reinforcement learning: A survey,” *arXiv preprint arXiv:2302.03122*, 2023.
- [5] B. Thananjeyan, A. Balakrishna, S. Nair, M. Luo, K. Srinivasan, M. Hwang, J. E. Gonzalez, J. Ibarz, C. Finn, and K. Goldberg, “Recovery rl: Safe reinforcement learning with learned recovery zones,” *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4915–4922, 2021.
- [6] L. Schäfer, F. Christianos, J. Hanna, and S. V. Albrecht, “Decoupling exploration and exploitation in reinforcement learning,” in *ICML 2021 Workshop on Unsupervised Reinforcement Learning*, 2021.
- [7] S. Ha, P. Xu, Z. Tan, S. Levine, and J. Tan, “Learning to walk in the real world with minimal human effort,” *arXiv preprint arXiv:2002.08550*, 2020.
- [8] J. Achiam, D. Held, A. Tamar, and P. Abbeel, “Constrained policy optimization,” in *International conference on machine learning*. PMLR, 2017, pp. 22–31.

- [9] T.-Y. Yang, J. Rosca, K. Narasimhan, and P. J. Ramadge, “Projection-based constrained policy optimization,” *arXiv preprint arXiv:2010.03152*, 2020.
- [10] D. Kim, Y. Kim, K. Lee, and S. Oh, “Safety guided policy optimization,” in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 2462–2467.
- [11] R. Cheng, G. Orosz, R. M. Murray, and J. W. Burdick, “End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 3387–3395.
- [12] Y. S. Shao, C. Chen, S. Kousik, and R. Vasudevan, “Reachability-based trajectory safeguard (rts): A safe and fast reinforcement learning safety layer for continuous control,” *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 3663–3670, 2021.
- [13] A. Wachi and Y. Sui, “Safe reinforcement learning in constrained markov decision processes,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 9797–9806.
- [14] M. Luo, A. Balakrishna, B. Thananjeyan, S. Nair, J. Ibarz, J. Tan, C. Finn, I. Stoica, and K. Goldberg, “Mesa: Offline meta-rl for safe adaptation and fault tolerance,” *arXiv preprint arXiv:2112.03575*, 2021.
- [15] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *International conference on machine learning*. PMLR, 2017, pp. 1126–1135.
- [16] W. F. Whitney, M. Bloesch, J. T. Springenberg, A. Abdolmaleki, K. Cho, and M. Riedmiller, “Decoupled exploration and exploitation policies for sample-efficient reinforcement learning,” *arXiv preprint arXiv:2101.09458*, 2021.
- [17] K. Srinivasan, B. Eysenbach, S. Ha, J. Tan, and C. Finn, “Learning to be safe: Deep rl with a safety critic,” *arXiv preprint arXiv:2010.14603*, 2020.
- [18] L. Zhang, Z. Yan, L. Shen, S. Li, X. Wang, and D. Tao, “Safety correction from baseline: Towards the risk-aware policy in robotics via dual-agent reinforcement learning,” in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 9027–9033.
- [19] M. Fatemi, S. Sharma, H. Van Seijen, and S. E. Kahou, “Dead-ends and secure exploration in reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 1873–1881.
- [20] M. Fatemi, T. W. Killian, J. Subramanian, and M. Ghassemi, “Medical dead-ends and learning to identify high-risk states and treatments,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 4856–4870, 2021.
- [21] T. W. Killian, S. Parbhoo, and M. Ghassemi, “Risk sensitive dead-end identification in safety-critical offline reinforcement learning,” *arXiv preprint arXiv:2301.05664*, 2023.
- [22] G. Thomas, Y. Luo, and T. Ma, “Safe reinforcement learning by imagining the near future,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 13 859–13 869, 2021.
- [23] M. Janner, J. Fu, M. Zhang, and S. Levine, “When to trust your model: Model-based policy optimization,” *Advances in neural information processing systems*, vol. 32, 2019.
- [24] E. Altman, *Constrained Markov decision processes*. CRC press, 1999, vol. 7.
- [25] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, “Deep reinforcement learning: A brief survey,” *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, 2017.
- [26] I. Kostrikov, A. Nair, and S. Levine, “Offline reinforcement learning with implicit q-learning,” *arXiv preprint arXiv:2110.06169*, 2021.
- [27] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” in *International conference on machine learning*. PMLR, 2018, pp. 1861–1870.
- [28] A. Ray, J. Achiam, and D. Amodei, “Benchmarking safe exploration in deep reinforcement learning,” *arXiv preprint arXiv:1910.01708*, vol. 7, no. 1, p. 2, 2019.
- [29] Q. Yang, T. D. Simão, S. H. Tindemans, and M. T. Spaan, “Wcsac: Worst-case soft actor critic for safety-constrained reinforcement learning,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 12, 2021, pp. 10 639–10 646.
- [30] S. Feng, H. Sun, X. Yan, H. Zhu, Z. Zou, S. Shen, and H. X. Liu, “Dense reinforcement learning for safety validation of autonomous vehicles,” *Nature*, vol. 615, no. 7953, pp. 620–627, 2023.