
Constrained Proximal Policy Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The problem of constrained reinforcement learning (CRL) holds significant importance
2 as it provides a framework for addressing critical safety satisfaction concerns
3 in the field of reinforcement learning (RL). However, with the introduction of
4 constraint satisfaction, the current CRL methods necessitate the utilization of
5 second-order optimization or primal-dual frameworks with additional Lagrangian
6 multipliers, resulting in increased complexity and inefficiency during implementation.
7 To address these issues, we propose a novel first-order feasible method named
8 Constrained Proximal Policy Optimization (CPPO). By treating the CRL problem
9 as a probabilistic inference problem, our approach integrates the Expectation-
10 Maximization framework to solve it through two steps: 1) calculating the optimal
11 policy distribution within the feasible region (E-step), and 2) conducting a first-
12 order update to adjust the current policy towards the optimal policy obtained in the
13 E-step (M-step). We establish the relationship between the probability ratios and
14 KL divergence to convert the E-step into a convex optimization problem. Further-
15 more, we develop an iterative heuristic algorithm from a geometric perspective to
16 solve this problem. Additionally, we introduce a conservative update mechanism to
17 overcome the constraint violation issue that occurs in the existing feasible region
18 method. Empirical evaluations conducted in complex and uncertain environments
19 validate the effectiveness of our proposed method, as it performs at least as well as
20 other baselines.

21 1 Introduction

22 In recent years, reinforcement learning (RL) has achieved huge success in various aspects (Le
23 et al., 2022; Li et al., 2022; Silver et al., 2018), especially in the field of games. However, due
24 to the increased safety requirements in practice, researchers are starting to consider the constraint
25 satisfaction in RL. Compared with unconstrained RL, constrained RL (CRL) incorporates certain
26 constraints during the process of maximizing cumulated rewards, which provides a framework to
27 model several important topics in RL, such as safe RL (Paternain et al., 2022), highlighting the
28 importance of this problem in industrial applications.

29 The current methods for solving the CRL problem can be mainly classified into two categories:
30 primal-dual method (Paternain et al., 2022; Stooke et al., 2020; Zhang et al., 2020; Altman, 1999) and
31 feasible region method (Achiam et al., 2017; Yang et al., 2020). The primal-dual method introduces
32 the Lagrangian multiplier to convert the constrained optimization problem into an unconstrained dual
33 problem by penalizing the infeasible behaviours, promising the CRL problem to be resolved in a
34 first-order manner. Despite the primal-dual framework providing a way to solve CRL in first-order
35 manner, the update of the dual variable, i.e., the Lagrangian multiplier, tends to be slow and unstable,
36 affecting the overall convergent speed of the algorithms. In contrast, the feasible region method
37 provides a faster learning method by introducing the concept of the feasible region into the trust
38 region method. With either searching in the feasible region (Achiam et al., 2017) or projecting into

39 the feasible region (Yang et al., 2020), the feasible region method can guarantee the generated policies
 40 stay in the feasible region. However, the introduction of the feasible region in the proposed method
 41 relies on computationally expensive second-order optimization using the inverse Fisher information
 42 matrix. This approach can lead to inaccurate estimations of the feasible region and potential constraint
 43 violations, as reported in previous studies (Ray et al., 2019).

44 To address the existing issues mentioned above, this paper proposed the Constrained Proximal
 45 Policy Optimization (CPPO) algorithm to solve the CRL problem in a first-order, easy-to-implement
 46 way. CPPO employs a two-step Expectation-Maximization approach to solve the problem by firstly
 47 calculating the optimal policy (E-step) and then conducting a first-order update to reduce the distance
 48 between the current policy and the optimal policy (M-step), eliminating the usage of the Lagrangian
 49 multiplier and the second-order optimization. The main contributions of this work are summarized as
 50 follows:

- 51 • To our best knowledge, the proposed method is the first **first-order feasible region method**
 52 without using dual variables or second-order optimization, which significantly reduces the
 53 difficulties in tuning hyperparameters and the computing complexity.
- 54 • An **Expectation-Maximization (EM)** framework based on **advantage value** and **probabil-**
 55 **ity ratio** is proposed for solving the CRL problem efficiently. By converting the CRL
 56 problem into a probabilistic inference problem, the CRL problem can be solved in first order
 57 manner without dual variables.
- 58 • To solve the convex optimization problem in E-step, we established the relationship between
 59 the probability ratios and KL divergence, and developed an **iterative heuristic algorithm**
 60 from a geometric perspective.
- 61 • A **recovery update** is developed when the current policy encounters constraint violation. In-
 62 spired by Bang-bang control, this update strategy can improve the performance of constraint
 63 satisfaction and reduce the switch frequency between normal update and recovery update.
- 64 • The proposed method is evaluated in several benchmark environments. The results manifest
 65 its comparable performance over other baselines in complex environments.

66 This paper is organized as follows. Section 2 introduces the concept of constrained Markov decision
 67 process and present an overview of related works in the field. The Expectation-Maximization
 68 framework and the technical details about the proposed constrained proximal policy optimization
 69 method are proposed in Section 3. Section 4 verifies the effectiveness of the proposed method through
 70 several testing scenarios and an ablation study is conducted to show the effectiveness of the proposed
 71 recovery update. Section 5 states the limitations and the broader impact of the proposed method.
 72 Finally, a conclusion is drawn in Section 6.

73 2 Preliminary and Related Work

74 2.1 Constrained Markov Decision Process

75 Constrained Markov Decision Process (CMDP) is a mathematical framework for modelling decision-
 76 making problems subjected to a set of cost constraints. A CMDP can be defined by a tuple
 77 $(\mathcal{S}, \mathcal{A}, P, r, \gamma, \mu, C)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow (0, 1)$
 78 is the transition kernel, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, $\gamma \in (0, 1)$ is the discount factor,
 79 $\mu : \mathcal{S} \rightarrow (0, 1)$ is the initial state distribution, and $C := \{c_i \in C \mid c_i : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}, i = 1, 2, \dots, m\}$
 80 is the set of m cost functions. For simplicity, we only consider a CRL problem with one constraint in
 81 the following paper and use c to represent the cost function. Note that, although we restrict our discus-
 82 sion to the case with only one constraint, the method proposed in this paper can be naturally extended
 83 to the multiple constraint case. However, the result may not as elegant as the one constraint case.
 84 Compared with the common Markov Decision Process (MDP), CMDP introduces a constraint on the
 85 cumulated cost to restrict the agent’s policies. Considering a policy $\pi(s | a) : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1)$, the goal
 86 of MDP is to find the π that maximizes the expected discounted returns $J_r(\pi) = \mathbb{E}_\tau [\sum_{t=0}^{\infty} \gamma^t r(s_t)]$,
 87 where τ is the trajectories generated based on π . Based on these settings, CMDP applied a threshold
 88 d on the expected discounted cost returns $J_c(\pi) = \mathbb{E}_\tau [\sum_{t=0}^{\infty} \gamma^t c(s_t)]$. Thus, the CMDP problem
 89 can be formed as finding policy π^* that $\pi^* = \operatorname{argmax}_\pi J_r(\pi)$ s.t. $J_c(\pi^*) \leq d$. The advan-
 90 tage function A and the cost advantage function A_c is defined as $A(s_t, a_t) = Q(s_t, a_t) - V(s_t)$

91 and $A_c(s_t, a_t) = Q_c(s_t, a_t) - V_c(s_t)$ where $Q(s_t, a_t) = \mathbb{E}_\tau [\sum_{t=0}^{\infty} \gamma^t r \mid s_0 = s_t, a_0 = a_t]$ and
 92 $V(s_t) = \mathbb{E}_\tau [\sum_{t=0}^{\infty} \gamma^t r \mid s_0 = s_t]$ are the corresponding Q-value and V-value for reward function,
 93 and $Q_c(s_t, a_t) = \mathbb{E}_\tau [\sum_{t=0}^{\infty} \gamma^t c \mid s_0 = s_t, a_0 = a_t]$ and $V_c(s_t) = \mathbb{E}_\tau [\sum_{t=0}^{\infty} \gamma^t c \mid s_0 = s_t]$ are the
 94 corresponding Q-value and V-value for cost function. Note that both A and A_c in the batch are
 95 centered to moves theirs mean to 0, respectively.

96 2.2 Related Work

97 2.2.1 Proximal Policy Optimization (PPO)

98 Proximal policy optimization (PPO) (Schulman et al., 2017) is a renowned on-policy RL algorithm for
 99 its stable performance and easy implementation. Based on the first-order optimization methodology,
 100 PPO addresses the challenge of the unconstrained RL problem through the surrogate objective
 101 function that proposed in Trust Region Policy Optimization (TRPO) (Schulman et al., 2015a). With
 102 the clipping and early stop trick, PPO can keep the new policy to stay within the trust region. Thanks
 103 to its stability and superior performance, the PPO algorithm has been employed in various subfields
 104 of RL like multi-agent RL (Yu et al., 2021), Meta-RL (Yu et al., 2020). However, due to the extra
 105 constraint requirements, the direct application of PPO in CRL problems is not feasible. The extra
 106 constraint requirements cause PPO not only restricted by the trust region but also the constraint
 107 feasible region, which significantly increases the challenge in conducting first-order optimization.
 108 Despite the difficulties in the direct application of PPO in CRL, researchers are still searching for a
 109 PPO-like method to solve CRL problems with stable and superior performance.

110 2.2.2 Constrained Reinforcement Learning

111 The current methods for solving the CRL problem can be mainly divided into two categories: primal-
 112 dual method (Paternain et al., 2022; Stooke et al., 2020; Zhang et al., 2020) and feasible region
 113 method (Achiam et al., 2017; Yang et al., 2020). The primal-dual method converts the original
 114 problem into a convex dual problem by introducing the Lagrangian multiplier. By updating the policy
 115 parameters and Lagrangian multiplier iteratively, the policies obtained by the primal-dual method
 116 will gradually converge towards a feasible solution. However, the usage of the Lagrange multiplier
 117 introduces extra hyperparameters into the algorithm and slows down the convergence speed of the
 118 algorithm due to the characteristic of the integral controller. Stooke et al. (2020) tries to solve this
 119 issue by introducing PID control into the update of the Lagrangian multiplier, but this modification
 120 will introduce more hyperparameters and cause the algorithm to be complex. Different from the
 121 primal-dual method, the feasible region method estimates the feasible region within the trust region
 122 using linear approximation and subsequently determines the new policy based on the estimated
 123 feasible region. A representative method is constrained policy optimization (CPO). By converting
 124 the CRL to a quadratically constrained linear program, CPO (Achiam et al., 2017) can solve the
 125 problem efficiently. However, the uncertainties inside the environment may cause an inaccurate cost
 126 assessment, which will affect the estimation of the feasible region and cause the learned policy to fail
 127 to meet the constraint requirements, as shown in Ray et al. (2019). Another issue of CPO is that it
 128 uses the Fisher information matrix to estimate the KL divergence in quadratic approximation, which
 129 is complex in computing and inflexible in network structure.

130 To address the second-order issue in CRL, several researchers (Zhang et al., 2020; Liu et al., 2022)
 131 proposed the EM-based algorithm in a first-order manner. FOCOPS (Zhang et al., 2020) obtain the
 132 optimal policy from advantage value, akin to the maximum entropy RL, and perform a first-order
 133 update to reduce the KL divergence between the current policy and the optimal policy. Despite its
 134 significant improvement in performance compared to CPO, FOCOPS still necessitates the use of a
 135 primal-dual method to attain a feasible optimal policy, which introduces a lot of hyperparameters for
 136 tuning, resulting in a more complex tuning process. CVPO (Liu et al., 2022) extends the maximum
 137 a posteriori policy optimization (MPO) (Abdolmaleki et al., 2018) method to the CRL problem,
 138 allowing for the efficient calculation of the optimal policy from Q value in an off-policy manner.
 139 However, this algorithm still requires the primal-dual framework in optimal policy calculation and
 140 necessitates additional samplings during the training, increasing the complexity of implementation.
 141 Thus, the development of a simple-to-implement, first-order algorithm with superior performance,
 142 remains a foremost goal for researchers in the CRL subfield.

143 3 Constrained Proximal Policy Optimization (CPPO)

144 As mentioned in Section 2, existing CRL methods often require second-order optimization for
 145 feasible region estimation or the use of dual variables for cost satisfaction. These approaches can be
 146 computationally expensive or result in slow convergence. To address these challenges, we proposed a
 147 two-step approach in an EM fashion named Constrained Proximal Policy Optimization (CPPO), the
 148 details will be shown in this section.

149 3.1 Modelling CRL as Inference

150 Instead of directly pursuing an optimal policy to maximize rewards, our approach involves concep-
 151 tualizing the problem of Constrained Reinforcement Learning (CRL) as a probabilistic inference
 152 problem. This is achieved by assessing the reward performance and constraint satisfaction of state-
 153 action pairs and subsequently increasing the likelihood of those pairs that demonstrate superior
 154 reward performance while adhering to the constraint requirement. Suppose the event of state-action
 155 pairs under policy π_θ can maximize reward is represented by optimality variable O , we assume
 156 the likelihood of state-action pairs being optimal is proportional to the exponential of its advantage
 157 value: $p(O = 1|(s, a)) \propto \exp(A(s, a)/\alpha)$ where α is a temperature parameter. Denote $q(a | s)$
 158 is the feasible posterior distribution estimated from the sampled trajectories under current policy
 159 π , $p_\pi(a | s)$ is the probability distribution under policy π , and θ is the policy parameters. We can
 160 have following evidence lower bound(ELBO) $\mathcal{J}(q, \theta)$ using surrogate function(see Appendix B for
 161 detailed proof)

$$\log p_{\pi_\theta}(O = 1) \geq \mathbb{E}_{s \sim d^\pi, a \sim \pi} \left[\frac{q(a|s)}{p_\pi(a|s)} A(s, a) \right] - \alpha D_{\text{KL}}(q \parallel \pi_\theta) + \log p(\theta) = \mathcal{J}(q, \theta), \quad (1)$$

162 where d^π is the state distribution under current policy π , $p(\theta)$ is a prior distribution of policy
 163 parameters. Considering $q(a | s)$ is a feasible policy distribution, we also have following constraint
 164 (Achiam et al., 2017)

$$J_c(\pi) + \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^\pi, a \sim \pi} \left[\frac{q(a|s)}{p_\pi(a|s)} A_c(s, a) \right] \leq d, \quad (2)$$

165 where d is the cost constraint. By performing iterative optimization of the feasible posterior distri-
 166 bution q (E-step) and the policy parameter θ (M-step), the lower bound $\mathcal{J}(q, \theta)$ can be increased,
 167 resulting in an enhancement in the likelihood of state-action pairs that have the potential to maximize
 168 rewards.

169 3.2 E-Step

170 3.2.1 Surrogate Constrained Policy Optimization

171 As mentioned in the previous section, we will firstly optimize the feasible posterior distribution q to
 172 maximize ELBO in E-step. The feasible posterior distribution q plays a crucial role in determining
 173 the upper bound of the ELBO since the KL divergence is non-negative. Consequently, q needs to be
 174 theoretically optimal to maximize the ELBO. By converting the soft KL constraint in Equation (1)
 175 into a hard constraint and combining the cost constraint in Equation (2), the optimization problem of
 176 q can be expressed as follows:

$$\begin{aligned} & \underset{q}{\text{maximize}} \quad \mathbb{E}_{s \sim d^\pi, a \sim \pi} \left[\frac{q(a|s)}{p_\pi(a|s)} A(s, a) \right] \\ & \text{s.t.} \quad J_c(\pi) + \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^\pi, a \sim \pi} \left[\frac{q(a|s)}{p_\pi(a|s)} A_c(s, a) \right] \leq d, \quad D_{\text{KL}}(q \parallel \pi) \leq \delta, \end{aligned} \quad (3)$$

177 where δ is the reverse KL divergence constraint that determine the trust region. During the E-step, it is
 178 important to note that the optimization is independent of θ , meaning that the policy π_θ remains fixed
 179 to the current sampled policy π . Even we know the closed-form expression of p_{π_θ} , it is impractical to
 180 solve the closed-form expression of q from Equation (3), as we still needs the closed-form expression
 181 of d^π for calculating. Therefore, we opt to represent the solution of q in a non-parametric manner
 182 by calculating the probability ratio $v = \frac{q(a|s)}{p_\pi(a|s)}$ for the sampled state-action pairs, allowing us to

183 avoid explicitly parameterizing q and instead leverage the probability ratio to guide the optimization
 184 process. After relaxing the reverse KL divergence constraint with the estimated reverse KL divergence
 185 calculated through importance sampling, we can obtain

$$\begin{aligned} & \underset{v}{\text{maximize}} && \mathbb{E}_{s \sim d^\pi, a \sim \pi} [vA(s, a)] \\ & \text{s.t.} && \mathbb{E}_{s \sim d^\pi, a \sim \pi} [vA_c(s, a)] \leq d' \\ & && \mathbb{E}_{\substack{s \sim d^\pi \\ a \sim \pi}} [v \log v] \leq \delta. \end{aligned} \quad (4)$$

186 where d' the scaled cost margin $d' = (1 - \gamma)(d - J_c(\pi))$. Although Equation (4) is convex
 187 optimization problem that can be directly solved through existing convex optimization algorithm, the
 188 existence of non-polynomial KL constraint tends to cause the optimization to be computationally
 189 expensive. To overcome this issue, the following proposition is proposed to relax Equation (4) into
 190 an linear optimization problem with quadratic constraint.

191 **Proposition 3.1.** *Denote v as the probability ratios $\frac{q(a|s)}{p_\pi(a|s)}$ calculated from sampled trajectories. If
 192 there are a sufficient number of sampled v , we have $\mathbb{E}[v] = 1$ and $\mathbb{E}[v \log v] \leq \text{Var}(v - 1)$.*

193 With Proposition 3.1, the relationship between reverse KL divergence and l^2 -norm of vector $v - 1$
 194 is constructed. Also, consider that the expectation of v equals 1, the optimization variable can be
 195 changed from v to $v - 1$. Let \bar{v} denote the vector consists of $v - 1$ and replace the reverse KL
 196 divergence constraint with the l^2 -norm constraint, Equation (4) can be rewritten in the form of vector
 197 multiplication

$$\begin{aligned} & \underset{\bar{v}}{\text{maximize}} && \bar{v} \cdot \mathbf{A} \\ & \text{s.t.} && \bar{v} \cdot \mathbf{A}_c \leq Nd', \|\bar{v}\|_2 \leq 2N\delta' \\ & && \mathbb{E}(\bar{v}) = 0, \bar{v} > -1 \text{ element-wise,} \end{aligned} \quad (5)$$

198 where \mathbf{A} and \mathbf{A}_c are the advantage value vectors for reward and cost (for all sampled state-action pairs
 199 in one rollout) respectively, N is the number of state-action pair samples, δ' is l^2 -norm constraint, and
 200 the element-wise lower bound of \bar{v} is -1 , as $v > 0$. Thus, the optimal feasible posterior distribution
 201 q expressed through \bar{v} can be obtained by solving the aforementioned optimization problem.

202 *Remark 3.2.* By replacing the non-polynomial KL constraint with an l^2 -norm constraint, the original
 203 optimization problem in Equation (4) can be reformulated as a geometric problem. This reformulation
 204 enables the use of the proposed heuristic method to efficiently solve the problem **without the need**
 205 **for dual variables**.

206 *Remark 3.3.* Our proposed method builds upon the idea presented in CVPO (Liu et al., 2022) of
 207 treating the CRL problem as a probabilistic inference problem. However, our approach improves
 208 upon their idea in two significant ways. Firstly, the probabilistic inference problem in our method is
 209 constructed based on **advantage value**, which is more effective in reducing the bias in estimating the
 210 cost return, compared to the Q-value used in CVPO. Secondly, while CVPO tries to directly calculate
 211 the value of $q(a|s)$, our method employs the **probability ratio** v to represent q . By replacing $q(a|s)$
 212 with v , our method only needs to find a vector of v whose elements are positive and $\mathbb{E}[v] = 1$, thereby
 213 negating the need to sample multiple actions in one state to calculate the extra normalizer that ensures
 214 q is a valid distribution. This results in a significant reduction in computational complexity.

215 3.2.2 Recovery update

216 Although the optimal solution q in Section 3.2.1 is applicable when the current policy is out of
 217 the feasible region, the inconsistent between optimal q and π_θ and the inaccurate cost evaluations
 218 tends to result in the generation of infeasible policies, as demonstrated in Ray et al. (2019) where
 219 CPO fail to satisfy constraint. To overcome this issue, a recovery update strategy is proposed for
 220 pushing the agent back to the feasible region. This strategy aims to minimize costs while preserving
 221 or minimizing any reduction in overall reward return. In the event that it is not possible to recover
 222 from the infeasible region without compromising the reward return, the strategy aims to identify an
 223 optimal policy within the feasible region that minimizes the adverse impact on the reward return. The
 224 optimization problem in recovery update can be expressed as

$$\begin{aligned} & \text{if } \bar{v} \cdot \mathbf{A} \geq 0 \text{ not exists when } \bar{v} \cdot \mathbf{A}_c \leq Nd': && \underset{\bar{v}}{\text{maximize}} \bar{v} \cdot \mathbf{A} \\ & \text{else:} && \underset{\bar{v}}{\text{minimize}} \bar{v} \cdot \mathbf{A}_c \\ & \text{s.t.} && \|\bar{v}\|_2 \leq 2N\delta', \mathbb{E}(\bar{v}) = 0, \bar{v} > -1 \text{ element-wise.} \end{aligned} \quad (6)$$

225 Figure 1 illustrates the recovery update strategy from the perspective of geometry. The blue, red,
 226 and yellow arrows represent the direction of minimizing cost, maximizing reward and the recovery
 227 update, respectively. The reward preservation region is defined by the zero reward boundary, which is
 228 depicted as the dashed line perpendicular to the red arrow. As a result, the semi-circle encompassing
 229 the red arrow indicates a positive increment in reward. Case 1 and Case 3 illustrate the case when the
 230 reward preservation region has an intersection with the feasible region. In these cases, we choose
 231 the direction of minimizing cost within the reward preservation region, e.g., the recovery update
 232 direction is coincident with the dashed line in Case 1, and the recovery update direction is coincident
 233 with the blue arrow in Case 3. Case 2 shows the case when there is no intersection between the
 234 reward preservation region and the feasible region. In this case, the direction with the least damage
 235 to reward is chosen. If we use an angle α to represent the direction of update, then we can have
 236 $\alpha = \text{Clip}(\alpha, \max(\theta_f, \theta_A + \pi/2), \pi)$, where θ_A represents the direction of \mathbf{A} , θ_f is the minimum
 237 angle that can point toward the feasible region.

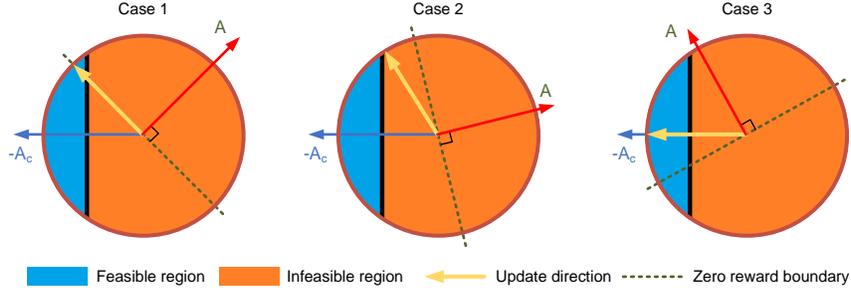


Figure 1: The illustration of recovery update.

238 To further improve the constraint satisfaction performance, a switching mechanism inspired by bang-
 239 bang control (Lasalle, 1960) is introduced. As shown in Figure 2, the agent will initially conduct
 240 normal update in Section 3.2.1; when the agent violates the cost constraint, it will switch to recovery
 241 update to reduce the cost until the cost is lower than the lower switch cost. By incorporating this
 242 switching mechanism, a margin is created between the lower switch cost and the cost constraint.
 243 This margin allows for a period of normal updates before the recovery update strategy is invoked.
 244 As a result, this mechanism prevents frequent switching between the two strategies, leading to
 245 improved performance in both reward collection and cost satisfaction. This switching mechanism
 246 effectively balances the exploration of reward-maximizing actions with the need to maintain constraint
 247 satisfaction.

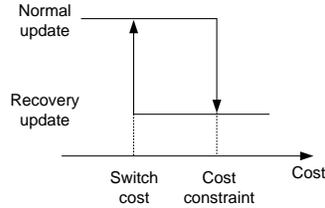


Figure 2: The switch mechanism inspired by bang-bang control. Once the current policy violates the cost constraint, the agent will switch to recovery update until it reaches the switch cost.

248 3.3 Heuristic algorithm from geometric interpretation

249 Section 3.2 and Section 3.4 provide a framework for solving CRL problem in theory. However,
 250 solving Equation (5) and Equation (6) in Section 3.2 is a tricky task in practice. To reduce the
 251 computation complexity, an iterative heuristic algorithm is proposed to solve this optimization
 252 problem from geometric interpretation. Recall Equation (5), the l_2 -norm can be interpreted as a radius
 253 constraint from the geometric perspective. Additionally, both the objective function and the cost

254 function are linear, indicating that the optimal solution lies on the boundary of the feasible region. By
 255 disregarding the element-wise bounds in Equation (5), we can consider the optimization problem as
 256 finding a optimal angle θ' on the A - A_c plane, in accordance with Theorem 3.4. The optimal solution
 257 can be expressed as $\bar{\mathbf{v}} = 2N\delta'(\cos\theta'\hat{\mathbf{A}}_c + \sin\theta'\hat{\mathbf{A}})$, where $\hat{\mathbf{A}}$ and $\hat{\mathbf{A}}_c$ are the orthogonal unit vectors
 258 of \mathbf{A} and \mathbf{A}_c respectively. Considering Assumption 3.5, we proposed a iterative heuristic algorithm
 259 to solve Equation (5) by firstly calculating the optimal angle θ' regardless the element-wise bound
 260 and obtain a initial solution $\bar{\mathbf{v}}$, then clip $\bar{\mathbf{v}}$ according to the element-wise bound and mask the clipped
 261 value, and iteratively update the rest unmasked elements according to aforementioned steps until all
 262 elements in $\bar{\mathbf{v}}$ are satisfy the element-wise bound. The detailed steps are outlined in Appendix C. For
 263 the recovery update in Section 3.2.2, the same algorithm can be used to find the angle that satisfy
 264 $\bar{\mathbf{v}} \cdot \mathbf{A}_c = Nd'$ or $\bar{\mathbf{v}} \cdot \mathbf{A} = 0$.

265 **Theorem 3.4.** *Given a feasible optimization problem of the form:*

$$\begin{aligned} & \underset{\bar{\mathbf{v}}}{\text{maximize}} \quad \bar{\mathbf{v}} \cdot \mathbf{A} \\ & \text{s.t.} \quad \bar{\mathbf{v}} \cdot \mathbf{A}_c \leq D, \quad \|\bar{\mathbf{v}}\|_2 \leq 2N\delta' \\ & \quad \quad \mathbb{E}(\bar{\mathbf{v}}) = \mathbb{E}(\mathbf{A}) = \mathbb{E}(\mathbf{A}_c) = 0 \end{aligned}$$

266 where $\bar{\mathbf{v}}$, \mathbf{A} , and \mathbf{A}_c are N -dimensional vectors, then the optimal solution $\bar{\mathbf{v}}$ will lie in the A - A_c
 267 plane determined by \mathbf{A}_c and \mathbf{A} .

268 **Assumption 3.5.** If the optimization problem in Theorem 3.4 has a optimal solution $\bar{\mathbf{v}}_{\text{opt}} =$
 269 $[\bar{v}_1, \bar{v}_2, \dots]$, and the same problem with element-wise lower bound constraint b has a optimal
 270 solution $\bar{\mathbf{v}}'_{\text{opt}} = [\bar{v}'_1, \bar{v}'_2, \dots]$, then $\bar{v}'_t = b$ where $\bar{v}_t \leq b$.

271 *Remark 3.6.* By utilizing the proposed heuristic algorithm, the optimal solution to Equation (5) can
 272 be obtained in just a few iterations. The time complexity of each iteration is $O(n)$, where n represents
 273 the number of unmasked elements. As a result, the computational complexity is significantly reduced
 274 compared to conventional convex optimization methods.

275 3.4 M-Step

276 After determining the optimal feasible posterior distribution q to maximize the upper bound of ELBO,
 277 an M-step is implemented to maximize ELBO by updating policy parameters θ in a supervised
 278 learning manner. Recall the definition of ELBO in Equation (1) in Section 3.1, by dropping the part
 279 that independent from θ , we will obtain following optimization problem

$$\underset{\theta}{\text{maximize}} \quad -\alpha D_{\text{KL}}(q \parallel \pi_\theta) + \log p(\theta). \quad (7)$$

280 Note that if we assume $p(\theta)$ is a Gaussian distribution, then $\log p(\theta)$ can be converted into $D_{\text{KL}}(\pi \parallel$
 281 $\pi_\theta)$ (see Appendix B for details). Using the same trick in Section 3.2.1 to convert soft KL constraint
 282 to hard KL constraint, the supervised learning problem in M-step can be expressed as

$$\begin{aligned} & \underset{\theta}{\text{minimize}} \quad D_{\text{KL}}(q \parallel \pi_\theta) \\ & \text{s.t.} \quad D_{\text{KL}}(\pi_\theta \parallel \pi) \leq \delta, \end{aligned} \quad (8)$$

283 Note that $D_{\text{KL}}(\pi_\theta \parallel \pi)$ is chosen to lower than δ so that the current policy π can be reached during
 284 the E-step in next update iteration to achieve robust update.

285 For Equation (7), it is a common practice for researchers to directly minimize the KL divergence,
 286 like CVPO (Liu et al., 2022) and MPO (Abdolmaleki et al., 2018). However, recall Equation (6), it
 287 is evident that the value of surrogate reward and cost are deeply connected to the projection of \mathbf{v}
 288 onto the A - A_c plane, while KL divergence can hardly reflect this kind of relationship between \mathbf{v} and
 289 surrogate value. Consequently, we choose to replace the original KL objective function
 290 with the l^2 -norm $\mathbb{E}[\|v - p_{\pi_\theta}/p_\pi\|_2]$, where v is the optimal probability ratio obtained in E-step and
 291 p_{π_θ}/p_π is the probability ratio under policy parameter θ . With this replacement, the optimization
 292 problem can be treated as a fixed-target tracking control problem. This perspective enables us to plan
 293 tracking trajectories that can consistently satisfy the cost constraint, enhancing the ability to maintain
 294 cost satisfaction throughout the learning process. The optimization problem after replacement can be
 295 rewritten as

$$\underset{\theta}{\text{minimize}} \quad \mathbb{E} \left[\left\| v - \frac{p_{\pi_\theta}}{p_\pi} \right\|_2 \right] \quad \text{s.t.} \quad D_{\text{KL}}(\pi_\theta \parallel \pi) \leq \delta, \quad (9)$$

296 To ensure the tracking trajectories can satisfy cost constraint at nearly all locations, we calculated the
 297 several recovery $\bar{\mathbf{v}}'$ under different δ' and guide $\frac{p_{\pi_\theta}}{p_\pi}$ to different $\bar{\mathbf{v}}$ according to the l_2 -norm of $\frac{p_{\pi_\theta}}{p_\pi}$,
 298 so that even $\|\frac{p_{\pi_\theta}}{p_\pi}\|_2$ is much smaller than $2N\delta'$, the new policy can still satisfy the cost constraint.
 299 Moreover, inspired by the proportional navigation (Yanushevsky, 2018), we also modify the recovery
 300 update gradient from $(v - \frac{p_{\pi_\theta}}{p_\pi})\frac{\partial\pi_\theta}{\partial\theta}$ to $((\beta(v - \frac{p_{\pi_\theta}}{p_\pi}) + (1 - \beta)\mathbf{A}'_c)\frac{\partial\pi_\theta}{\partial\theta})$ to reduce the cost during
 301 the tracking, where \mathbf{A}'_c is the projection of $v - \frac{p_{\pi_\theta}}{p_\pi}$ on cost advantage vector \mathbf{A}_c . In according with
 302 Theorem 3.7, the lower-bound clipping mechanism similar with PPO is applied on updating $\frac{p_{\pi_\theta}}{p_\pi}$ in
 303 M-step to satisfy the forward KL constraint (see Appendix C for details).

304 **Theorem 3.7.** For a probability ratio vector $\bar{\mathbf{v}}$, if the variance of $\bar{\mathbf{v}}$ is constant, then the upper bound
 305 of the approximated forward KL divergence $D_{\text{KL}}(\pi_\theta \parallel \pi)$, will decrease as the element-wise lower
 306 bound of $\bar{\mathbf{v}}$ increase.

307 Apart from E-step and M-step introduced in Section 3.2 and Section 3.4, our method shares the same
 308 Generalized Advantage Estimator (GAE) technique (Schulman et al., 2015b) with PPO in calculating
 309 the advantage value A and A_c . The main steps of CPPO are summarized in Appendix C.

310 4 Experiment

311 In this section, Safety Gym (Ray et al., 2019) benchmark environments and Circle environment
 312 (Achiam et al., 2017) are used to verify and evaluate the performance of the proposed method. Five
 313 test scenarios, namely CarPush, PointGoal, PointPush, PointCircle, and AntCircle are evaluated.
 314 The detailed information about the test scenarios can be seen in Appendix D. Three algorithms
 315 are chosen as the benchmarks to compare the learning curves and the constraint satisfaction: CPO
 316 (Achiam et al., 2017), PPO-Lagrangian method (simplified as PPO_lag), and TRPO-Lagrangian
 317 method (simplified as TRPO_lag) (Ray et al., 2019). CPO is chosen as the representative of the
 318 feasible region method. PPO_lag and TRPO_lag are treated as the application of the primal-dual
 319 method in first-order optimization and second-order optimization. TRPO and PPO are also used in
 320 this section as unconstrained performance references. For a fair comparison, all of the algorithms
 321 use the same policy network and critic network. The detail of the hyperparameter setting is listed in
 322 Appendix E.

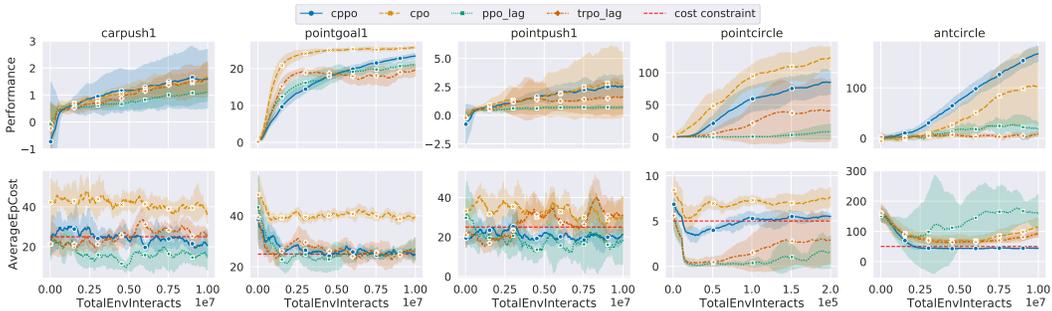


Figure 3: The learning curves for comparison, CPPO is the method proposed in this paper.

323 **Performance and Constraint Satisfaction:** Figure 3 compares the learning curves of the proposed
 324 method and other benchmark algorithms in terms of the episodic return and the episodic cost. The
 325 first row records the undiscounted episodic return for performance comparison, and the second row is
 326 the learning curves of the episodic cost for constraint satisfaction analysis, where the red dashed line
 327 indicates the cost constraint. The learning curves for the Push and Goal environments are averaged
 328 over 6 random seeds, while those for the Circle environments are averaged over 4 random seeds.
 329 The curve itself represents the mean value, and the shadow indicates the standard deviation. In
 330 terms of performance comparison, it was observed that CPO can achieve the highest reward return
 331 in PointGoal and PointCircle. The proposed CPPO method, on the other hand, achieves similar or
 332 even higher reward return in the remaining test scenarios. However, when considering constraint
 333 satisfaction, CPO fails to satisfy the constraint in all four tasks due to approximation errors, as

334 previously reported in Ray et al. (2019). In contrast, CPPO successfully **satisfies the constraint**
 335 in all five environments, showing the effectiveness of the proposed **recovery update**. Referring to
 336 the learning curves in Circle scenarios, it can be seen that the primal-dual based CRL methods, i.e.,
 337 PPO_lag and TRPO_lag, suffer from the slow and unstable update of the dual variable, causing the
 338 conservative performance in PointCircle and slow cost satisfaction in AntCircle. On the other hand,
 339 CPPO can achieve a faster learning speed in Circle environment by **eliminating the need for the**
 340 **dual variable**. Overall, the experimental results demonstrate the effectiveness of CPPO in solving
 341 the CRL problem.

342 **Ablation Study:** An ablation study was conducted to investigate the impact of the recovery update in
 343 CPPO. Figure 4 presents the reward performance and cost satisfaction of CPPO with and without the
 344 recovery update in the PointCircle environment. The results indicate that without the recovery update,
 345 CPPO achieves higher reward performance; however, the cost reaches 15, which significantly violates
 346 the cost constraint. In contrast, when the recovery update is applied, CPPO successfully satisfies
 347 the constraint, thereby demonstrating the importance of the recovery update in ensuring constraint
 satisfaction.

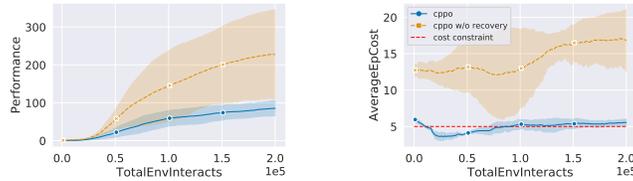


Figure 4: The comparison between CPPO with and without recovery update in PointCircle.

348

349 5 Limitations and Boarder Impact

350 Although our proposed method has shown its ability in test scenarios, there still exist some limitations.
 351 Firstly, CPPO method is an on-policy constrained RL, which suffers from lower sampling efficiency
 352 compared to other off-policy algorithms, potentially limiting its applicability in real-world scenarios.
 353 Additionally, the convergence of our method is not yet proven. However, we believe that our work
 354 will offer researchers a new EM perspective for using PPO-like algorithms to solve the problem
 355 of constrained RL, thereby leading to the development of more efficient and stable constrained RL
 356 algorithms.

357 6 Conclusion

358 In this paper, we have introduced a novel first-order Constrained Reinforcement Learning (CRL)
 359 method called CPPO. Our approach avoids the use of the primal-dual framework and instead treats the
 360 CRL problem as a probabilistic inference problem. By utilizing the Expectation-Maximization (EM)
 361 framework, we address the CRL problem through two key steps: the E-step, which focuses on deriving
 362 a theoretically optimal policy distribution, and the M-step, which aims to minimize the difference
 363 between the current policy and the optimal policy. Through the non-parametric representation of the
 364 policy using probability ratios, we convert the CRL problem into a convex optimization problem
 365 with a clear geometric interpretation. As a result, we propose an iterative heuristic algorithm that
 366 efficiently solves this optimization problem without relying on the dual variable. Furthermore, we
 367 introduce a recovery update strategy to handle approximation errors in cost evaluation and ensure
 368 constraint satisfaction when the current policy is infeasible. This strategy mitigates the impact of
 369 approximation errors and strengthens the capability of our method to satisfy constraints. Notably, our
 370 proposed method does not require second-order optimization techniques or the use of the primal-dual
 371 framework, which simplifies the optimization process. Empirical experiments have been conducted to
 372 validate the effectiveness of our proposed method. The results demonstrate that our approach achieves
 373 comparable or even superior performance compared to other baseline methods. This showcases
 374 the advantages of our method in terms of simplicity, efficiency, and performance in the field of
 375 Constrained Reinforcement Learning.

376 **References**

- 377 Abdolmaleki, A., Springenberg, J. T., Tassa, Y., Munos, R., Heess, N., and Riedmiller, M. Maximum a posteriori
378 policy optimisation. *arXiv preprint arXiv:1806.06920*, 2018.
- 379 Achiam, J., Held, D., Tamar, A., and Abbeel, P. Constrained policy optimization. In *International conference on*
380 *machine learning*, pp. 22–31. PMLR, 2017.
- 381 Altman, E. *Constrained Markov decision processes*, volume 7. CRC press, 1999.
- 382 Lasalle, J. The ‘bang-bang’ principle. *IFAC Proceedings Volumes*, 1(1):503–507, 1960. ISSN 1474-6670. 1st
383 International IFAC Congress on Automatic and Remote Control, Moscow, USSR, 1960.
- 384 Le, N., Rathour, V. S., Yamazaki, K., Luu, K., and Savvides, M. Deep reinforcement learning in computer vision:
385 a comprehensive survey. *Artificial Intelligence Review*, 55(4):2733–2819, 2022.
- 386 Li, Q., Peng, Z., Feng, L., Zhang, Q., Xue, Z., and Zhou, B. Metadrive: Composing diverse driving scenarios for
387 generalizable reinforcement learning. *IEEE transactions on pattern analysis and machine intelligence*, 2022.
- 388 Liu, Z., Cen, Z., Isenbaev, V., Liu, W., Wu, S., Li, B., and Zhao, D. Constrained variational policy optimization
389 for safe reinforcement learning. In *International Conference on Machine Learning*, pp. 13644–13668. PMLR,
390 2022.
- 391 Paternain, S., Calvo-Fullana, M., Chamon, L. F., and Ribeiro, A. Safe policies for reinforcement learning via
392 primal-dual methods. *IEEE Transactions on Automatic Control*, 2022.
- 393 Ray, A., Achiam, J., and Amodei, D. Benchmarking safe exploration in deep reinforcement learning. *arXiv*
394 *preprint arXiv:1910.01708*, 7:1, 2019.
- 395 Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. Trust region policy optimization. In *International*
396 *conference on machine learning*, pp. 1889–1897. PMLR, 2015a.
- 397 Schulman, J., Moritz, P., Levine, S., Jordan, M., and Abbeel, P. High-dimensional continuous control using
398 generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015b.
- 399 Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms.
400 *arXiv preprint arXiv:1707.06347*, 2017.
- 401 Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran,
402 D., Graepel, T., et al. A general reinforcement learning algorithm that masters chess, shogi, and go through
403 self-play. *Science*, 362(6419):1140–1144, 2018.
- 404 Stooke, A., Achiam, J., and Abbeel, P. Responsive safety in reinforcement learning by pid lagrangian methods.
405 In *International Conference on Machine Learning*, pp. 9133–9143. PMLR, 2020.
- 406 Yang, T.-Y., Rosca, J., Narasimhan, K., and Ramadge, P. J. Projection-based constrained policy optimization.
407 *arXiv preprint arXiv:2010.03152*, 2020.
- 408 Yanushevsky, R. *Modern missile guidance*. CRC Press, 2018.
- 409 Yu, C., Velu, A., Vinitsky, E., Wang, Y., Bayen, A., and Wu, Y. The surprising effectiveness of ppo in cooperative,
410 multi-agent games. *arXiv preprint arXiv:2103.01955*, 2021.
- 411 Yu, T., Quillen, D., He, Z., Julian, R., Hausman, K., Finn, C., and Levine, S. Meta-world: A benchmark and
412 evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*, pp. 1094–1100.
413 PMLR, 2020.
- 414 Zhang, Y., Vuong, Q., and Ross, K. First order constrained optimization in policy space. *Advances in Neural*
415 *Information Processing Systems*, 33:15338–15349, 2020.