

Robot Reinforcement Learning on the Constraint Manifold

Puze Liu¹, Davide Tateo¹, Haitham Bou-Ammar² and Jan Peters¹

¹ Department of Computer Science, Technische Universität Darmstadt, Germany

² Huawei R&D London, United Kingdom

{puze, davide}@robot-learning.de,

haitham.ammam@huawei.com, jan.peters@tu-darmstadt.de

Abstract: Reinforcement learning in robotics is extremely challenging due to many practical issues, including safety, mechanical constraints, and wear and tear. Typically, these issues are not considered in the machine learning literature. One crucial problem in applying reinforcement learning in the real world is Safe Exploration, which requires physical and safety constraints satisfaction throughout the learning process. To explore in such a safety-critical environment, leveraging known information such as robot models and constraints is beneficial to provide more robust safety guarantees. Exploiting this knowledge, we propose a novel method to learn robotics tasks in simulation efficiently while satisfying the constraints during the learning process.

Keywords: Robot Learning, Constrained Reinforcement Learning, Safe Exploration

1 Introduction

Despite the notable success of Deep Reinforcement Learning (RL) in solving complex tasks in the discrete world, video games, as well as continuous control problems in simulation [1, 2, 3, 4], applying RL in the real world remains a challenging task. One important factor that cannot be neglected in real-world applications is the necessity of satisfying constraints. Many practical considerations can be formulated in the form of constraints, such as safety and mechanical viability. For example, in the robot manipulation task, the robot should not take actions that damage the environment and can not take actions that exceed its feasible range. However, typical RL algorithms, which maximize the cumulative reward by continuous trial and error, do not take into account the satisfaction of constraints during the exploration process. Exploring the environment while meeting the constraints is a challenging problem.

Safe exploration is an significant field of RL which requires to comply with the constraints during the whole learning process [5]. There are several safe exploration frameworks in the literature: a possible direction is proposed in [6, 7] that relies on prior knowledge (policies, value functions) to initialize the system in a safe region and gradually increase the area of exploration using new information obtained from the environment. Other approaches rely on the definition of a safe policy [8, 9], which tries to pull the agent back to a safe state. However, these choices require excessive work in defining such policy, and safe policies could conflict against each other when multiple constraints are violated. Finally, other methods incorporate model information of constraints with model-free RL algorithms and do not require

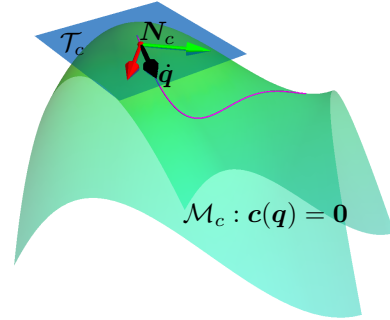


Figure 1: Acting on the Tangent Space of the Constraint Manifold. The constraint set $c(q) = 0$ is a differentiable manifold \mathcal{M}_c embedded in the original state space. We use a set of basis vectors N_c to represent the span of tangent space \mathcal{T}_c . The tangent space velocity/acceleration can be determined by a coordinate based on the bases, and the control action is determined based on the tangent space velocity/acceleration, the resulting trajectory is maintained on the constraint manifold.

the definition of a manual policy [10, 11, 12]. In these approaches, the agent tries to find the feasible action using constrained optimization techniques at each time step.

In this paper, we propose a novel method, Acting on the Tangent Space of the Constraint Manifold (ATACOM), which the agent explores in the tangent space of the constraint manifold, as shown in Figure 1. The proposed method convert the constrained RL problem to a typical unconstrained RL problem. This method allows us to utilize *any* model-free RL algorithms while maintaining the constraints below the tolerance. Furthermore, ATACOM can handle both equality and inequality constraints. For example, in the task of a robot wiping a table, the end-effector should move on the surface of the table (equality constraints) while the joint positions and velocities are within its joint limits (inequality constraints). In addition, for tasks with equality constraints, our method explores the *lower-dimensional* manifold embedded in the original action space. To test our method, we demonstrate three different tasks, *CircleMoving*, *PlanarAirHockey*, and *IiwaAirHockey*[13], with different combinations of equality and inequality constraints. We test five state-of-the-art model-free RL algorithms (PPO, TRPO, DDPG, TD3, SAC) in each environment. The result shows that all algorithms can learn the policy efficiently while maintaining the constraints below the tolerance.

The advantage of ATACOM can be summarized as follows: (i) can deal both with **equality and inequality constraints**. All of the constraints at each time step are maintained below the tolerance during the whole learning process. (ii) does not require an initial feasible policy, the agent can **learn from scratch**. (iii) requires **no manual safe backup policy** to move the system back into the safe region. (iv) can be applied to any model-free RL algorithm, using both **deterministic and stochastic policies**. (v) can focus the exploration on the **lower-dimensional manifold** instead of exploring in the original action space for equality constrained problem. (vi) have **better learning performance** as the inequality constraints restrict to a smaller feasible state-action space. As a downside, our method requires: (i) differentiable constraint functions. (ii) a sufficient accurate invertible dynamics model of the robot or a well-performed tracking controller.

Related Work. In the last decades, Constrained Markov Decision Processes (CMDP) [14] has attracted a lot of interest from RL researchers, to solve constrained control problems. Under this framework, several different forms of constraints have been studied. One important form of constraint is the expected cost below a threshold. Many works maximize the expected return while maintaining the expected cost below a threshold [11, 15, 16, 17, 18, 19, 20, 21]. Different types of constrained optimization techniques are applied in the policy update process. Achiam et al. proposed a trust-region method Constrained Policy Optimization (CPO) inspired from Trust Region Policy Optimization (TRPO) [19]. Liu et al. proposed the interior point method for policy optimization [16]. Another type of approach is to adapt the Lagrangian relaxation method for the constrained RL setting [14, 17, 15, 18]. Lastly, Chow et al. proposed a method to generate the Lyapunov function that guarantees constraints satisfaction [11, 22]. These approaches focus on the constraint of the cumulative cost and require an initial feasible policy. However, this cumulative cost criterion cannot ensure safety for tasks where avoiding catastrophic failures is crucial, e.g., car crashing.

Other approaches focus on the state dependant constraints, which should be fulfilled at every time step. To meet this requirement, safe exploration methods can be employed. Garcia, et al., proposed a method based on a risk function and a baseline agent, where the control action is sampled based on the evaluation of the risk [6]. The shielding [8] and backup policy [9] frameworks interfere with the control action to pull the system back to the safe states. These approaches require a manual defined safe policy. Berkenkamp, et al. [7], Wachi, et al. [23], Koller, et al. [24], and Hewing, et al. [25] proposed model-based approaches to ensure the safety. These approaches start from an initial feasible policy and progressively increase the safe region based on the learned dynamics model. Recent methods also try to incorporate the model and the constraint information with the model-free RL algorithms. Dalal, et al., added a safe layer which analytically finds the closest action w.r.t the policy derived one [10]. Cheng, et al., proposed a barrier function method to guarantee safety during the exploration [12]. Finally, other approaches has also address the safety issue from different perspectives, such as the policy composition [26, 27] and reachability-based approach [28, 29, 30].

Our approach considers the second group of constraints. However, different from other comparable methods, ATACOM does not require an initial policy, it can learn from scratch. In addition, our method does not require a backup policy either, as the constraint violations are forecasted and corrected at each step. Furthermore, our method is not specifically restricted to any learning algorithm.

2 Learning on the Constraint Manifold

In this section, we discuss ATACOM in detail. We first introduce the mathematical notation used in this paper. Then, to demonstrate the core concept, we start with a simple scenario that the constraint on the subset of the state variable \mathbf{q} and the action can be formulated as a function of the state velocity $\mathbf{a} = \Lambda(\dot{\mathbf{q}})$. Next, considering the continuity of velocity (sampling over the velocity does not ensure the continuity), we convert the original state constraint to a viability constraint that incorporates the velocity of the constraint. The action is chosen as a function of the acceleration $\mathbf{a} = \Lambda(\ddot{\mathbf{q}})$. From a robotics point of view, this \mathbf{a} can be the torque applied to each joint, and Λ is the inverse dynamics model. Then, to cope with the velocity limit, we add the viability condition to the acceleration. Lastly, we discuss some important practical issues of ATACOM, such as the error correction, the tangent space convention that determines the null space bases.

Definitions We consider the CMDP with continuous state-action space. A CMDP is a tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma, \mathcal{C})$, where \mathcal{S} is a state space, \mathcal{A} is an action space, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is a transition kernel, γ is a discount factor, and $\mathcal{C} : \{c_i : \mathcal{S} \rightarrow \mathbb{R} | i \in 1, \dots, k\}$ is a set of *immediate state-constraint* functions.

Assumption In this paper, we decompose the state variable $\mathbf{s} \in \mathcal{S}$ into the directly controllable state $\mathbf{q} \in \mathcal{Q}$ and uncontrollable state $\mathbf{x} \in \mathcal{X}$, i.e., $\mathbf{s} = [\mathbf{q} \ \mathbf{x}]^\top$. We assume that the constraints $\mathbf{c}(\mathbf{q}) \leq \mathbf{0}$ are known and depend purely on the controllable state. In addition, we assume that the action \mathbf{a} can be determined based on the i -th order time derivative of the controllable state, i.e., $\mathbf{a} = \Lambda(\mathbf{q}^{(i)})$, $i \in \{1, 2, \dots\}$. For example, we can determine the joint torque using an inverse dynamics model or send the desired positions/velocities obtained via integration to a tracking controller (e.g., PID controller). The general form of the constrained reinforcement learning problem can be formulated as

$$\max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{s}_t, \mathbf{a}_t} \left[\sum_{t=0}^T \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right], \quad \text{s.t.} \quad \mathbf{c}(\mathbf{q}_t) \leq \mathbf{0}.$$

2.1 State Constraints

The state constraints are defined as

$$\mathbf{f}(\mathbf{q}) = \mathbf{0}, \quad \mathbf{g}(\mathbf{q}) \leq \mathbf{0}, \quad (1)$$

where $\mathbf{f} : \mathbb{R}^Q \rightarrow \mathbb{R}^F$, $\mathbf{g} : \mathbb{R}^Q \rightarrow \mathbb{R}^G$ are two C^2 mappings for F equality and G inequality constraints, and $F < Q$. We add the slack variables $\boldsymbol{\mu} \in \mathbb{R}^G$ in inequality constraints to convert the original constraints (1) into equality constraints

$$\mathbf{c}(\mathbf{q}, \boldsymbol{\mu}) = [\mathbf{f}(\mathbf{q}) \quad \mathbf{g}(\mathbf{q}) + \frac{1}{2}\boldsymbol{\mu}^2]^\top = \mathbf{0}. \quad (2)$$

The constraint set (2) is a $(F + G)$ dimensional manifold embedded in $(Q + G)$ dimensional space. We calculate the time derivative of (2)

$$\dot{\mathbf{c}}(\mathbf{q}, \boldsymbol{\mu}, \dot{\mathbf{q}}, \dot{\boldsymbol{\mu}}) = \begin{bmatrix} \mathbf{J}_f(\mathbf{q}) & \mathbf{0} \\ \mathbf{J}_g(\mathbf{q}) & \text{diag}(\boldsymbol{\mu}) \end{bmatrix} \begin{bmatrix} \dot{\mathbf{q}} \\ \dot{\boldsymbol{\mu}} \end{bmatrix} = \mathbf{J}_c(\mathbf{q}, \boldsymbol{\mu}) \begin{bmatrix} \dot{\mathbf{q}} \\ \dot{\boldsymbol{\mu}} \end{bmatrix}, \quad (3)$$

with the Jacobians $\mathbf{J}_f \in \mathbb{R}^{F \times Q}$ and $\mathbf{J}_g \in \mathbb{R}^{G \times Q}$ of $\mathbf{f}(\mathbf{q})$ and $\mathbf{g}(\mathbf{q})$, respectively. Both Jacobians are combined into the Jacobian Matrix $\mathbf{J}_c(\mathbf{q}, \boldsymbol{\mu}) \in \mathbb{R}^{(F+G) \times (Q+G)}$ of the complete constraint set.

We can find the null space matrix $\mathbf{N}_c(\mathbf{q}, \boldsymbol{\mu}) = \text{Null}[\mathbf{J}_c(\mathbf{q}, \boldsymbol{\mu})] \in \mathbb{R}^{(Q+G) \times (Q-F)}$ via SVD [31] or QR [32] decomposition, such that $\mathbf{J}_c(\mathbf{q}, \boldsymbol{\mu})\mathbf{N}_c(\mathbf{q}, \boldsymbol{\mu}) = \mathbf{0}$. Each column of the orthogonal matrix $\mathbf{N}_c(\mathbf{q}, \boldsymbol{\mu})$ represents a basis vector of the null space of $\mathbf{J}_c(\mathbf{q}, \boldsymbol{\mu})$. These null space bases can also be viewed as the tangent space bases of the constraint manifold as illustrated in Figure 1. We can construct a tangent space velocity of the constraint manifold by a coordinate $\boldsymbol{\alpha}$ as

$$\begin{bmatrix} \dot{\mathbf{q}}_\tau \\ \dot{\boldsymbol{\mu}}_\tau \end{bmatrix} = \mathbf{N}_c(\mathbf{q}, \boldsymbol{\mu})\boldsymbol{\alpha}, \quad (4)$$

Substituting $[\dot{\mathbf{q}} \ \dot{\boldsymbol{\mu}}]^\top$ of (3) by $[\dot{\mathbf{q}}_\tau \ \dot{\boldsymbol{\mu}}_\tau]^\top$ of (4), we have the constraint velocity

$$\dot{\mathbf{c}}(\mathbf{q}, \boldsymbol{\mu}, \dot{\mathbf{q}}, \dot{\boldsymbol{\mu}}) = \mathbf{J}_c(\mathbf{q}, \boldsymbol{\mu})\mathbf{N}_c(\mathbf{q}, \boldsymbol{\mu})\boldsymbol{\alpha} = \mathbf{0}. \quad (5)$$

Equation (5) implies that the constraints do not change regardless of the choice of $\boldsymbol{\alpha}$. Based on this concept, the ATACOM method can be summarized as follows: Starting from a feasible point

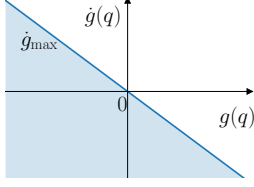


Figure 2: Viability Constraints

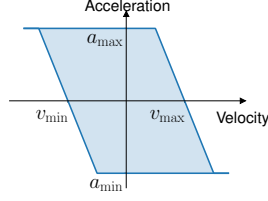


Figure 3: Feasible Acceleration Region

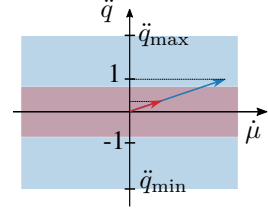


Figure 4: Tangent space bases

$(\mathbf{q}(0), \boldsymbol{\mu}(0)) \in \{(\mathbf{q}, \boldsymbol{\mu}) | \mathbf{c}(\mathbf{q}, \boldsymbol{\mu}) = \mathbf{0}\}$, we choose the tangent space velocity $[\dot{\mathbf{q}}_{\mathcal{T}}(t), \dot{\boldsymbol{\mu}}_{\mathcal{T}}(t)]^T = \mathbf{N}_c(\mathbf{q}(t), \boldsymbol{\mu}(t))\boldsymbol{\alpha}(t)$ and the corresponding action as $\mathbf{a}(t) = \boldsymbol{\Lambda}(\dot{\mathbf{q}}_{\mathcal{T}}(t))$. Thus, the constrained RL problem is converted into an unconstrained RL problem. The resulting trajectory $\mathbf{q}(t)$ satisfies the constraints $\mathbf{c}(\mathbf{q}(t), \boldsymbol{\mu}(t)) = \mathbf{0}$.

2.2 Viability Constraints

For a physical system, it is often required a continuous velocity command. However, directly sampling velocities $\dot{\mathbf{q}}$ does not ensure this continuity. A simple solution is to sample accelerations, apply force to the system or determine the velocity via integration. Furthermore, when considering inequality constraints, it is also desirable that $\dot{g}(\mathbf{q}, \dot{\mathbf{q}}) \leq 0$ when $g(\mathbf{q}) = 0$ to avoid overshooting. We convert the original state constraints (1) to *viability constraints* inspired by the linear viability condition [33]

$$\begin{aligned} \mathbf{f}(\mathbf{q}) + \mathbf{K}_f \dot{\mathbf{f}}(\mathbf{q}, \dot{\mathbf{q}}) &= \mathbf{f}(\mathbf{q}) + \mathbf{K}_f \mathbf{J}_f(\mathbf{q}) \dot{\mathbf{q}} = \mathbf{0}, \\ \mathbf{g}(\mathbf{q}) + \mathbf{K}_g \dot{\mathbf{g}}(\mathbf{q}, \dot{\mathbf{q}}) &= \mathbf{g}(\mathbf{q}) + \mathbf{K}_g \mathbf{J}_g(\mathbf{q}) \dot{\mathbf{q}} \leq \mathbf{0}, \end{aligned} \quad (6)$$

with diagonal matrices $\mathbf{K}_f \in \mathbb{R}^{F \times F}$, $\mathbf{K}_g \in \mathbb{R}^{G \times G}$ having all positive entries. The matrices \mathbf{K}_f and \mathbf{K}_g determine the maximum velocities of the constraints $\dot{\mathbf{f}}$ and $\dot{\mathbf{g}}$ w.r.t to the value of the constraints. The viability constraint of the inequality constraint is illustrated in Figure 2. When $g(\mathbf{q}) < 0$, the upper bound of the constraint velocity is $\dot{g}_{\max} > 0$ which means that it is still possible to get close to the constraint boundary. However, if $g(\mathbf{q}) > 0$, the upper bound of constraint velocity \dot{g}_{\max} should be smaller than zero to pull the violations back.

Analogous to the derivations from equation (2) and (3), we have

$$\mathbf{c}(\mathbf{q}, \dot{\mathbf{q}}, \boldsymbol{\mu}) = \begin{bmatrix} \mathbf{f}(\mathbf{q}) + \mathbf{K}_f \mathbf{J}_f(\mathbf{q}) \dot{\mathbf{q}} \\ \mathbf{g}(\mathbf{q}) + \mathbf{K}_g \mathbf{J}_g(\mathbf{q}) \dot{\mathbf{q}} + \frac{1}{2} \boldsymbol{\mu}^2 \end{bmatrix} = \mathbf{0}, \quad (7)$$

and

$$\dot{\mathbf{c}}(\mathbf{q}, \dot{\mathbf{q}}, \ddot{\mathbf{q}}, \boldsymbol{\mu}, \dot{\boldsymbol{\mu}}) = \underbrace{\begin{bmatrix} \mathbf{K}_f \mathbf{J}_f(\mathbf{q}) & \mathbf{0} \\ \mathbf{K}_g \mathbf{J}_g(\mathbf{q}) & \text{diag}(\boldsymbol{\mu}) \end{bmatrix}}_{\mathbf{J}_c(\mathbf{q}, \boldsymbol{\mu})} \begin{bmatrix} \ddot{\mathbf{q}} \\ \dot{\boldsymbol{\mu}} \end{bmatrix} + \underbrace{\begin{bmatrix} \mathbf{J}_f(\mathbf{q}) \dot{\mathbf{q}} + \mathbf{K}_f \mathbf{b}_f(\mathbf{q}, \dot{\mathbf{q}}) \\ \mathbf{J}_g(\mathbf{q}) \dot{\mathbf{q}} + \mathbf{K}_g \mathbf{b}_g(\mathbf{q}, \dot{\mathbf{q}}) \end{bmatrix}}_{\boldsymbol{\psi}(\mathbf{q}, \dot{\mathbf{q}})} = \mathbf{0}, \quad (8)$$

where $\mathbf{b}_f(\mathbf{q}, \dot{\mathbf{q}}) = \dot{\mathbf{q}}^T \mathbf{H}_f(\mathbf{q}) \dot{\mathbf{q}}$, $\mathbf{b}_g(\mathbf{q}, \dot{\mathbf{q}}) = \dot{\mathbf{q}}^T \mathbf{H}_g(\mathbf{q}) \dot{\mathbf{q}}$ and $\mathbf{H}_f \in \mathbb{R}^{F \times Q \times Q}$, $\mathbf{H}_g \in \mathbb{R}^{G \times Q \times Q}$ are Hessians of $\mathbf{f}(\mathbf{q})$, $\mathbf{g}(\mathbf{q})$, respectively. We can construct the joint acceleration as

$$\begin{bmatrix} \ddot{\mathbf{q}} \\ \dot{\boldsymbol{\mu}} \end{bmatrix} = -\mathbf{J}_c^\dagger(\mathbf{q}, \boldsymbol{\mu}) \boldsymbol{\psi}(\mathbf{q}, \dot{\mathbf{q}}) + \mathbf{N}_c(\mathbf{q}, \boldsymbol{\mu}) \boldsymbol{\alpha}, \quad (9)$$

with the pseudo-inverse $\mathbf{J}_c^\dagger(\mathbf{q}, \boldsymbol{\mu})$ and the null space matrix $\mathbf{N}_c(\mathbf{q}, \boldsymbol{\mu})$ of the Jacobian $\mathbf{J}_c(\mathbf{q}, \boldsymbol{\mu})$, respectively. The first term in equation (9) is the necessary acceleration that maintains the curvature of the constraints manifold (7) and the second term is the tangent space acceleration of the constraints. When starting from the point $[\mathbf{q}(0), \dot{\mathbf{q}}(0), \boldsymbol{\mu}(0)] \in \{(\mathbf{q}, \dot{\mathbf{q}}, \boldsymbol{\mu}) | \mathbf{c}(\mathbf{q}, \dot{\mathbf{q}}, \boldsymbol{\mu}) = \mathbf{0}\}$ and sampling over $\boldsymbol{\alpha}$, the joint acceleration $\ddot{\mathbf{q}}$ and the corresponding action \mathbf{a} satisfy the constraints.

2.3 Viability Acceleration Bound

In robotics as well as other mechanical systems, it is important to consider the velocity constraints of the actuator. Also, the acceleration should be bounded properly to avoid overshooting. We again use the concept of viability to determine the upper and lower bound of the acceleration

$$\begin{aligned} \mathbf{a}_u &= \max(\min(\mathbf{a}_{\max}, -\mathbf{K}_a(\mathbf{q} - \mathbf{v}_{\max})), \mathbf{a}_{\min}), \\ \mathbf{a}_l &= \min(\max(\mathbf{a}_{\min}, -\mathbf{K}_a(\mathbf{q} - \mathbf{v}_{\min})), \mathbf{a}_{\max}), \end{aligned}$$

with the minimum and the maximum joint velocity limits $\mathbf{v}_{\min, \max}$ and the acceleration limits $\mathbf{a}_{\min, \max}$, $K_a > 0$ is a constant. The feasible acceleration region is illustrated in Figure 3. Analogous to the viability constraints, the feasible region of the acceleration is modified depending on the state of joint velocities. This technique effectively prevents overshooting.

2.4 Error Correction and Control Action Selection

For time-continuous systems, the state is obtained at a certain sampling rate and the action is applied for a certain period. This time discretization results in constraint violations at each time step. Therefore, we add an error correction term. We construct a P-controller with a diagonal matrix \mathbf{K}_c for the constraints

$$\begin{bmatrix} \ddot{\mathbf{q}}_E \\ \dot{\boldsymbol{\mu}}_E \end{bmatrix} = -\mathbf{J}_c^\dagger \mathbf{K}_c \mathbf{c}(\mathbf{q}, \dot{\mathbf{q}}, \boldsymbol{\mu}). \quad (10)$$

Combining (9) with (10), we get the joint acceleration applied to the system

$$\begin{bmatrix} \ddot{\mathbf{q}} \\ \dot{\boldsymbol{\mu}} \end{bmatrix} = -\mathbf{J}_c^\dagger(\mathbf{q}, \dot{\mathbf{q}}, \boldsymbol{\mu}) [\mathbf{K}_c \mathbf{c}(\mathbf{q}, \dot{\mathbf{q}}, \boldsymbol{\mu}) + \boldsymbol{\psi}(\mathbf{q}, \dot{\mathbf{q}})] + \mathbf{N}_c(\mathbf{q}, \dot{\mathbf{q}}, \boldsymbol{\mu}) \boldsymbol{\alpha}. \quad (11)$$

The first term on the RHS is the necessary accelerations/velocities to maintain the constraints and the second term on the RHS is the tangent acceleration that can be explored freely. Figure 5 illustrates the vector field of error correction term and null space term of the circle constraint. The gray curves show the sampled trajectories converging to the constraint manifold due to the error correction.

The control action can be determined by $\mathbf{a} = \Lambda(\ddot{\mathbf{q}})$ at different levels. For example, we can use the inverse dynamics model to calculate the joint torque when the robot is controlled via torque command. We can also apply the integration method to determine the desired positions/velocities, then use a sufficient accurate tracking controller (e.g., PID controller + Feedforward Term) to track the desired trajectory. However, the tracking errors could potentially cause hazardous constraint violations. In this paper, we control the joint torque calculated by a perfect dynamic model in simulation to simplify the analysis and to exclude the constraint violations caused by the tracking error of the controller. We present the block diagram of the controlling framework in Appendix A.

2.5 Null Space Convention

The orthogonal null space matrix \mathbf{N}_c can be determined through SVD or QR decomposition. However, the representation of the null space bases is not unique. It is difficult to preserve the consistency

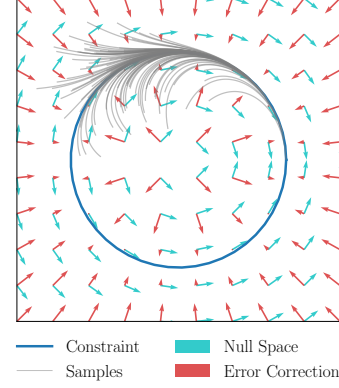


Figure 5: Vector Field of a Circle Constraint. The constraint $q_1^2 + q_2^2 - 1 = 0$ is the blue circle. The cyan arrows show the $\mathbf{N}_c \boldsymbol{\alpha}$ with $\boldsymbol{\alpha} = 1$. The red arrows demonstrate the error correction term $-\mathbf{J}_c^\dagger \mathbf{K}_c \mathbf{c}$. The gray lines show 100 trajectories from different initial points. All trajectories converge to the constraint manifold.

Algorithm 1: ATACOM

Input: Constraint: $\mathbf{f}, \mathbf{g}, \mathbf{J}_f, \mathbf{J}_g, \mathbf{b}_f, \mathbf{b}_g$. Scale parameter: $\mathbf{K}_c, \mathbf{K}_f, \mathbf{K}_g$. Time step ΔT .

- 1 **for each episode do**
- 2 Initial feasible state \mathbf{s}_0 , slack variable $\boldsymbol{\mu}_0$.
- 3 **for each time step k do**
- 4 Sample policy action $\boldsymbol{\alpha}_k \sim \pi(\cdot | \mathbf{s}_k)$.
- 5 Observe the $\mathbf{q}_k, \dot{\mathbf{q}}_k$ from \mathbf{s}_k .
- 6 Compute $\mathbf{J}_{c,k} = \mathbf{J}_c(\mathbf{q}_k, \dot{\mathbf{q}}_k, \boldsymbol{\mu}_k)$, $\boldsymbol{\psi}_k = \boldsymbol{\psi}(\mathbf{q}_k, \dot{\mathbf{q}}_k)$, $\mathbf{c}_k = \mathbf{c}(\mathbf{q}_k, \dot{\mathbf{q}}_k, \boldsymbol{\mu}_k)$.
- 7 Compute the RCEF of tangent space basis of $\mathbf{N}_{c,k}^R$.
- 8 Compute the tangent space acceleration $[\ddot{\mathbf{q}}_k \ \dot{\boldsymbol{\mu}}_k]^\top \leftarrow -\mathbf{J}_{c,k}^\dagger [\mathbf{K}_c \mathbf{c}_k + \boldsymbol{\psi}_k] + \mathbf{N}_{c,k}^R \boldsymbol{\alpha}_k$.
- 9 Clip the joint acceleration $\ddot{\mathbf{q}}_k \leftarrow \text{clip}(\ddot{\mathbf{q}}_k, \mathbf{a}_l, \mathbf{a}_u)$.
- 10 Integrate the slack variable $\boldsymbol{\mu}_{k+1} = \boldsymbol{\mu}_k + \dot{\boldsymbol{\mu}}_k \Delta T$.
- 11 Apply the control action $\mathbf{a}_k = \Lambda(\ddot{\mathbf{q}}_k)$ to the environment.
- 12 Observe the next state \mathbf{s}_{k+1} and reward r_k from the environment.
- 13 Provide the transition tuple $(\mathbf{s}_k, \boldsymbol{\alpha}_k, \mathbf{s}_{k+1}, r_k)$ to the RL algorithm.

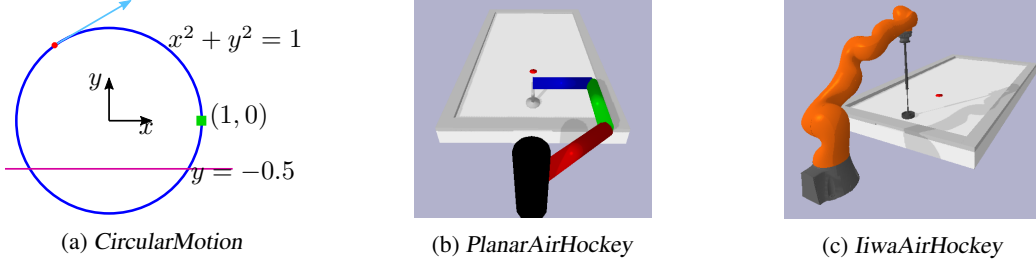


Figure 6: Experiment Environments

of the null space bases computed by the numerical decomposition method [34, 35]. To solve this issue, we propose a convention to ensure the uniqueness of the null space bases.

Each column of the null space matrix N_c is a unit vector indicating a direction of $[\ddot{\mathbf{q}} \ \dot{\boldsymbol{\mu}}]^T$. However, this unit vector could sometimes contribute majorly to the part of the slack variable and the entries for the joint accelerations could be very small. As a result, the joint acceleration obtained from $\boldsymbol{\alpha} \in [\boldsymbol{\alpha}_{\min}, \boldsymbol{\alpha}_{\max}]$ can only cover a small region of the acceleration. As illustrated in Figure 4, the red arrow is a unit basis vector of the tangent space, and the reachable joint acceleration by a universal scaling factor could only cover part of the feasible joint acceleration, as the red area shown in Figure 4. To alleviate the previously mentioned issue, we compute the Reduced Column Echelon Form (RCEF) of the null space matrix $N_c^R = \text{RCEF}(N_c)$. Given that the RCEF of a matrix is unique, we obtain unique bases of the null space. In addition, for RCEF, each row containing a leading 1 has zeros in all its other entries. Generally speaking, there exist N independent joints whose acceleration can be solely determined by $\boldsymbol{\alpha}$, where N is the dimensions of the null space. Also, we can define the feasible range of $\boldsymbol{\alpha}$ as $\alpha_i \in [\ddot{q}_{i,\min}, \ddot{q}_{i,\max}]$. Through this convention, the joint acceleration is able to cover the full feasible range. Null space bases and feasible region are shown as the blue vector and the blue area in Figure 4.

3 Experiments and Evaluation

To illustrate the properties of our approach, we demonstrate three different experiments in this section. We first demonstrate a toy task, *CircularMotion*, shown in Figure 6a. In this task, we consider state equality, inequality, and velocity constraints. Secondly, we show a robotic environment with only inequality constraints, *PlanarAirHockey* shown in Figure 6b. A 3 DoF planar robot playing the hitting task in the air hockey scenario while keeping the end-effector inside the table boundary and the robot’s joint positions and velocities within its limits. Finally, we demonstrate, *IiwaAirHockey* in Figure 6c, a 7-DoF KUKA IIWA robot learning the hitting task in the simulator. In addition to the constraints of the 3-dimensional task, we add an equality constraint to ensure that the robot end-effector stays on the table surface. More details can be found in the Appendix B and D.

CircularMotion. In this task, shown in Figure 6a, the red point tries to move along a unit circle in 2D space while keeping the velocity of each direction below the velocity limits and maintaining the position above a certain height. The objective is to reach the target point (green square) located in $(1, 0)$. The control action is the acceleration $\mathbf{a} = [\ddot{x} \ \ddot{y}]^T$.

We compare ATACOM with two other approaches for the task. (i) *TerminatedCircularMotion* where the episode terminates when the maximum constraint violations up to a threshold. (ii) *ErrorCorrectionCircularMotion*, where the error correction term in (10) is added before the action is applied to the environment. We test five model-free RL algorithms (SAC, DDPG, TD3, TRPO, and PPO implemented in *Mushroom-RL* [36]) for each approach.

Figure 7 shows the learning curve and constraint violations of all test RL algorithms for ATACOM. Every algorithm is able to improve the learning performance and SAC outperforms the others methods, which matches our expectations. Figure 7b and 7c show the maximum constraint function and maximum joint velocity constraints at each time step. It can be shown that the maximum constraint violations during the whole learning process remain small. The velocity limit violations are zero after 30 epochs which means the learned policies try to fully exploit the constraints.

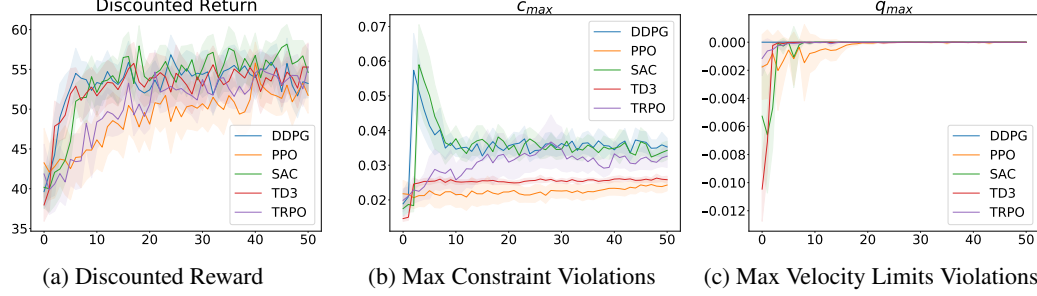


Figure 7: ATACOM for the *CircularMotion*. 7a shows the discounted return at each epochs. 7b and 7c shows the maximum constraint violations and maximum joint velocity limits violations.

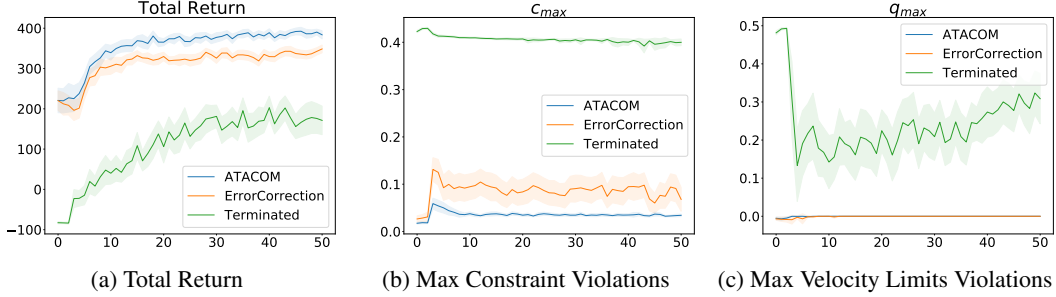


Figure 8: Comparison between ATACOM, *Terminated-*, and *ErrorCorrectionCircularMotion*.

In Figure 8 we compare ATACOM, *TerminatedCircularMotion*, and *ErrorCorrectionCircularMotion*. We select the best learning algorithms for each approach (SAC for all cases). Compared to the baselines, our method focus on a lower-dimensional exploration space that avoid the constraint violations while the others do not. We can conclude from the Figures that ATACOM has not only lower constraint violations but also better learning performance among the three approaches.

PlanarAirHockey. In this experiment, we apply a 3-joints planar robot for the air-hockey hitting task, as illustrated in Figure 6b. The end-effector of the robot is kept on the table surface and the objective is to hit the puck to the opponent’s goal. In this environment, we only consider inequality constraints, i.e., the robot end-effector should stay inside the table’s region, and the joint positions and velocities should not exceed its limits. The control action is the joint torque obtained by the inverse dynamics model. In this experiment, we assume the dynamics model is perfectly known to eliminate the constraint violation due to the tracking error or the model mismatch.

In this task, we compared ATACOM with the *SafeLayer* method proposed by Dalal et al.,[10] and the *Unconstrained* air-hockey environments. Since the *SafeLayer* method at the beginning requires a free exploration process to learn the constraint function, we only compare the learning performance and the constraint violations after this process. For the unconstrained environment, the robot is completely free to explore, and the episodes only terminate when the maximum episode step is reached. In this experiment, we only compare the best DDPG result after the parameter sweep, as the available implementation of *SafeLayer* only supports DDPG. Additional experiment of *PlanarDefend* can be found in Appendix C.2

The result is shown in Figure 9. We can see that ATACOM have the best learning performance and the minimum constraint violations among the three methods. *SafeLayer* did not learn the constraint function of joint velocities properly. Furthermore, the learned constraints appear to be too restrictive to learn a good policy. Compared to the method of *Unconstrained* approach, although ATACOM has the same dimension as the *Unconstrained*, ATACOM explores only in the feasible region while the *Unconstrained* approach explores the whole state-action space. This consideration explains why ATACOM outperforms the baselines in terms of learning performances.

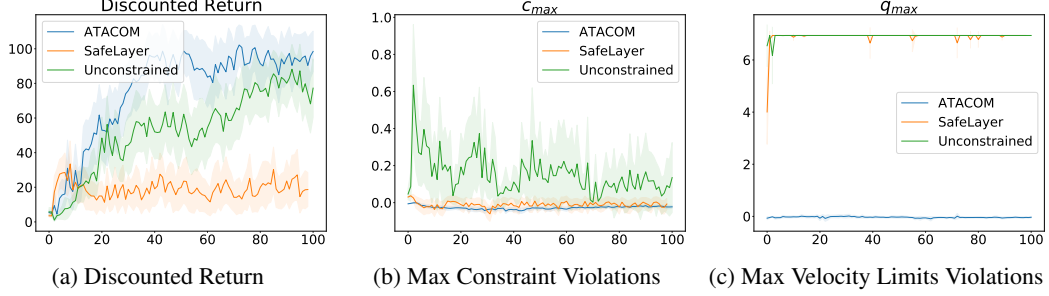


Figure 9: Comparison between ATACOM, *SafeLayer* [10], *Unconstrained AirHockey* in DDPG.

***IiwaAirHockey*.** In the third experiment, we demonstrate the same air-hockey hitting task with a KUKA LBR IIWA14 Robot in the Pybullet simulator. In this task, we add an equality constraint to ensure the end-effector stays on the table surface. We also add inequality constraints to avoid collision of the end-effector and joint limit constraints as mention in the *PlanarAirHockey* task. In addition, we also add inequality constraints to avoid the collision between the 4/6-th link and the table. We enforce the joint velocity limits in the simulation as the real-world’s KUKA controller does. We compare the impact of different simulation step sizes. The step size refers to the sampling frequency in the real world.

At each simulation step, the torque is computed by the previous agent’s action until the new control action is received. The error correction term is added at each time step. We choose the simulation step size as 0.02s, 0.004s, 0.002s, and 0.001s and keep agent control frequency to be 50Hz.

Figure 10 demonstrates the discounted return at the final epoch, the maximum constraint violation c_{max} , and the average constraint violations c_{avg} throughout the learning process. For a sufficiently small step size, such as 0.004s, 0.002s, 0.001s, the learning agent is able to learn the hitting policy. When the step size is too big, e.g., 0.02s, the error correction term dominates the control action and the agent has a poor learning performance. From the result in this experiment, we demonstrate that ATACOM is able to solve high dimensional tasks. The simulation result also provides us the guidance for the real-world application: The higher frequency of sampling and error correction, the smaller constraint violations will occur. In addition, we compare ATACOM with Riemannian Motion Policies [37] in Appendix C.4

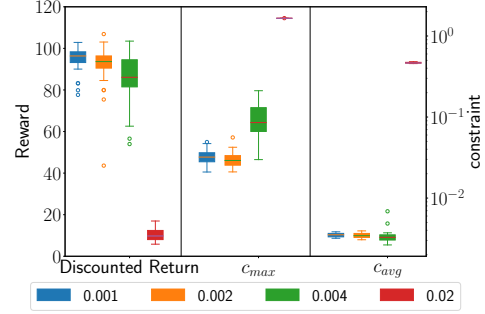


Figure 10: Box Plot of *IiwaAirHockey* with different choice of time step size.

4 Conclusion

In this article, we present ATACOM, a safe exploration method for Constrained RL based on the knowledge of the model and the mathematical formulations of constraints. ATACOM explores the tangent space of the constraint manifold. This exploration technique allows us to utilize any type of model-free RL method while maintaining the constraint violations below a small threshold. From the experiments, we have shown that ATACOM not only has small constraints violations but also better learning performance w.r.t. the other baselines. These performance gains occur because ATACOM only focuses on the safe region (from inequality constraint) and subspace (from equality constraint) of the whole state-action space.

However, our method still has some limitations. Our method requires a sufficiently accurate model or a good tracking controller. This assumption does not hold in most real-world applications since model errors, disturbances, and sensor noise could potentially cause unexpected constraint violations. To deploy this method in real-world robots, we will focus on the model mismatch problem and may require a backup policy to avoid too stringent constraint design. Furthermore, our current approach only focuses on the constraint with only controllable state $c(q) = 0$, even if preliminary results (Appendix E) suggests an extension into constraints with the uncontrollable state.

Acknowledgments

This project was supported by the CSTT fund from Huawei Tech R&D (UK). The support provided by China Scholarship Council (No. 201908080039) is acknowledged.

References

- [1] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the Game of Go with Deep Neural Networks and Tree Search. *Nature*, 529(7587):484–489, 2016.
- [2] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, et al. Grandmaster Level in StarCraft II Using Multi-Agent Reinforcement Learning. *Nature*, 575(7782):350–354, 2019.
- [3] Y. Duan, X. Chen, C. X. B. Edu, J. Schulman, P. Abbeel, and P. B. Edu. Benchmarking Deep Reinforcement Learning for Continuous Control. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2016.
- [4] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous Control with Deep Reinforcement Learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2016.
- [5] J. García and F. Fernández. A Comprehensive Survey on Safe Reinforcement Learning. *Journal of Machine Learning Research*, 16:1437–1480, 2015.
- [6] J. Garcia and F. Fernandez. Safe Exploration of State and Action Spaces in Reinforcement Learning. *Journal of Artificial Intelligence Research*, 45:515–564, 2012. ISSN 10769757.
- [7] F. Berkenkamp, M. Turchetta, A. P. Schoellig, and A. Krause. Safe Model-based Reinforcement Learning with Stability Guarantees. In *Proceedings of the Conference on Neural Information Processing Systems (NIPS)*, 2017.
- [8] M. Alshiekh, R. Bloem, R. Ehlers, B. Könighofer, S. Niekum, and U. Topcu. Safe Reinforcement Learning via Shielding. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- [9] A. Hans, D. Schneeß, A. M. Schäfer, and S. Udfluft. Safe Exploration for Reinforcement Learning. In *ESANN*, 2008. ISBN 2930307080.
- [10] G. Dalal, K. Dvijotham, M. Vecerik, T. Hester, C. Paduraru, and Y. Tassa. Safe Exploration in Continuous Action Spaces. *arXiv preprint arXiv:1801.08757*, 2018.
- [11] Y. Chow, O. Nachum, E. Duenez-Guzman, and M. Ghavamzadeh. A Lyapunov-based Approach to Safe Reinforcement Learning. In *Proceedings of the Conference on Neural Information Processing Systems (NIPS)*, 2018.
- [12] 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 *AAAI Conference on Artificial Intelligence*, pages 3387–3395. AAAI Press, 2019.
- [13] P. Liu, D. Tateo, H. Bou-Ammar, and J. Peters. Efficient and Reactive Planning for High Speed Robot Air Hockey. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021.
- [14] E. Altman. Constrained Markov Decision Processes with Total Cost Criteria: Lagrangian Approach and Dual Linear Program. *Mathematical methods of operations research*, 48(3): 387–417, 1998.
- [15] C. Tessler, D. J. Mankowitz, and S. Mannor. Reward Constrained Policy Optimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- [16] Y. Liu, J. Ding, and X. Liu. IPO: Interior-point Policy Optimization under Constraints. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.

- [17] A. Stooke, J. Achiam, and P. Abbeel. Responsive Safety in Reinforcement Learning by PID Lagrangian Methods. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- [18] D. Ding, X. Wei, Z. Yang, Z. Wang, and M. R. Jovanovic. Provably Efficient Safe Exploration via Primal-Dual Policy Optimization. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 130, 2021.
- [19] J. Achiam, D. Held, A. Tamar, and P. Abbeel. Constrained Policy Optimization. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.
- [20] H. B. Ammar, R. Tutunov, and E. Eaton. Safe policy search for lifelong reinforcement learning with sublinear regret. In *International Conference on Machine Learning*. PMLR, 2015.
- [21] A. I. Cowen-Rivers, D. Palenicek, V. Moens, M. Abdullah, A. Sootla, J. Wang, and H. Ammar. Samba: Safe model-based & active reinforcement learning, 2020.
- [22] Y. Chow, O. Nachum, A. Faust, E. Duenez-Guzman, and M. Ghavamzadeh. Lyapunov-based Safe Policy Optimization for Continuous Control. In *Reinforcement Learning for Real Life (RL4RealLife) Workshop in the 36th International Conference on Machine Learning*, 2019.
- [23] A. Wachi, Y. Sui, Y. Yue, and M. Ono. Safe Exploration and Optimization of Constrained MDPs Using Gaussian Processes. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- [24] T. Koller, F. Berkenkamp, M. Turchetta, and A. Krause. Learning-Based Model Predictive Control for Safe Exploration. In *Proceedings of the IEEE Conference on Decision and Control*, 2018.
- [25] L. Hewing, K. P. Wabersich, M. Menner, and M. N. Zeilinger. Learning-Based Model Predictive Control: Toward Safe Learning in Control. *Annual Review of Control, Robotics, and Autonomous Systems*, 3(1):269–296, 2020. ISSN 2573-5144.
- [26] N. D. Ratliff, J. Issac, D. Kappler, S. Birchfield, and D. Fox. Riemannian motion policies. *arXiv preprint arXiv:1801.02854*, 2018.
- [27] J. Urain, A. Li, P. Liu, C. D’Eramo, and J. Peters. Composable energy policies for reactive motion generation and reinforcement learning. In *Robotics: Science and Systems*, 2021.
- [28] A. K. Akametalu, J. F. Fisac, J. H. Gillula, S. Kaynama, M. N. Zeilinger, and C. J. Tomlin. Reachability-based safe learning with gaussian processes. In *53rd IEEE Conference on Decision and Control*, pages 1424–1431. IEEE, 2014.
- [29] 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*, 6(2):3663–3670, 2021.
- [30] H. Krasowski, X. Wang, and M. Althoff. Safe reinforcement learning for autonomous lane changing using set-based prediction. In *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, pages 1–7. IEEE, 2020.
- [31] R. P. Singh and P. W. Likins. Singular Value Decomposition for Constrained Dynamical Systems. *Journal of Applied Mechanics, Transactions ASME*, 52(4):943–948, 1985.
- [32] S. S. Kim and M. J. Vanderploeg. QR Decomposition for State Space Representation of Constrained Mechanical Dynamic Systems. *Journal of Mechanisms, Transmissions, and Automation in Design*, 108:183–188, 1986.
- [33] M. Faroni, M. Beschi, N. Pedrocchi, and A. Visioli. Viability and Feasibility of Constrained Kinematic Control of Manipulators. *Robotics*, 7(3), jul 2018. ISSN 22186581.
- [34] L. W. Tu. *An introduction to manifolds*. Springer., 2011.
- [35] L. N. Trefethen and D. Bau III. *Numerical linear algebra*, volume 50. Siam, 1997.

- [36] C. D’Eramo, D. Tateo, A. Bonarini, M. Restelli, and J. Peters. Mushroomrl: Simplifying reinforcement learning research. *Journal of Machine Learning Research*, 22(131):1–5, 2021.
- [37] C.-A. Cheng, M. Mukadam, J. Issac, S. Birchfield, D. Fox, B. Boots, and N. Ratliff. Rmpflow: A computational graph for automatic motion policy generation. In *International Workshop on the Algorithmic Foundations of Robotics*, pages 441–457. Springer, 2018.

This appendix is structured as follows: we introduce the block diagram of the ATACOM in Appendix A. We describe the experiment environment in Appendix B. Then, we show the comparison of ATACOM, *ErrorCorrection*, *Terminated* approach with different RL method in the *CircularMotion* task in Appendix C.1. The comparison of ATACOM with different algorithms in *PlanarDefend* and *IiwaAirHockey* are shown in Appendix C.2 and Appendix C.3, respectively. The parameters of the environments, the learning algorithms, and the results of the parameter sweep are shown in Appendix D. Finally, we append an extension of ATACOM including the uncontrollable state in Appendix E. The results of the *CollisionAvoidance* environment are illustrated in Appendix E.1.

A Block Diagram of ATACOM

Here we will describe the controlling diagram of the ATACOM using viability constraints. We copy the overall control acceleration in Equation (11) here

$$\begin{bmatrix} \ddot{q} \\ \dot{\mu} \end{bmatrix} = \underbrace{-J_c^\dagger(q, \dot{q}, \mu) [K_c c(q, \dot{q}, \mu) + \psi(q, \dot{q})]}_{[\ddot{q}_{mm} \quad \dot{\mu}_{mm}]^\top} + \underbrace{N_c(q, \dot{q}, \mu) \alpha}_{[\ddot{q}_{null} \quad \dot{\mu}_{null}]^\top}.$$

The first term on the RHS tries to maintain the state on the constraint manifold. The second term tries to explore the tangent space of the constraint manifold based on the agent policy.

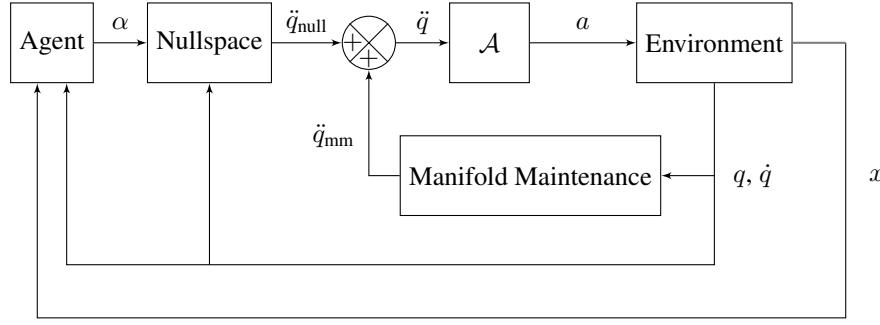


Figure 11: Block Diagram of the proposed method

As shown in Figure 11, the control action is determined by the inverted dynamics $a = \mathcal{A}(\ddot{q})$. In practice, we can apply the integration method to get desired positions and velocities as the input to an accurate tracking controller (e.g., PID controller) to determine the control action. However, this tracking error could potentially bring additional issues, like the interplay between the error in the constraints and tracking performance. In the experiment of this paper, we assume the model of dynamics is perfectly known to exclude the errors caused by the factors that are not considered in this paper.

B Environment Description

B.1 CircularMotion

In this environment, shown in Figure 6a, a point is moving on a circle with the following constraints

$$\begin{aligned} f : x^2 + y^2 - 1 &= 0, & g : -y - 0.5 &< 0, \\ |\dot{x}| - 1 &< 0, & |\dot{y}| - 1 &< 0. \end{aligned}$$

The first constraint ensures that the point moves on the circle. The second one ensures that the point is moving above the height -0.5 . The last two constraints limit the velocity of each component. The control action is the acceleration along each axis. The objective is to minimize the distance to the goal $(1, 0)$, shown as the green square

$$r(x, y) = \exp(-\sqrt{(x-1)^2 + (y-0)^2})$$

B.2 PlanarAirHockey

In this experiment, we solve the air hockey task with a 3-joints planar robot, as illustrated in Figure 6b. In this environment, we consider only inequality constraints, i.e., the robot end-effector should stay inside the table’s region, and the joint positions and velocities should not exceed its limits. The constraints are the following

$$\begin{aligned} g_1 : -x_{ee} + x_{table,l} < 0, & \quad g_2 : -y_{ee} + y_{table,l} < 0, & \quad g_3 : y_{ee} - y_{table,u} < 0, \\ g_{4,5,6} : q_i^2 - q_{i,l}^2 < 0 & \quad |\dot{q}_i| - \dot{q}_{i,l} < 0 & \quad \forall i \in \{1, 2, 3\} \end{aligned}$$

where (x_{ee}, y_{ee}) is the position of the robot end-effector, $x_{table,l}, y_{table,l}, y_{table,u}$ are the boundaries of the air-hockey table, q_i, \dot{q}_i refers to position and velocity of the i -th joint, and $q_{i,l}, \dot{q}_{i,l}$ are the position and velocity limit for the joint i . The control action is the torque applied on each joint. In this task, the forward kinematics, Jacobian matrix, Hessian matrix, and the inverse dynamics are calculate with the help of the *Pinocchio* rigid bodies dynamics library.

We initialize the puck and the robot in a given configuration. The objective is to hit the puck to the opponent’s goal as fast as possible. We define the reward as

$$r = \begin{cases} \exp[-8 \cdot \|\mathbf{p}_{ee} - \mathbf{p}_{puck}\| \cdot \text{clip}(\cos \theta, 0, 1)] - \lambda, & \text{if has not hit,} \\ 1 + r_{hit} + 0.1v_{x,hit} - \lambda, & \text{if has hit,} \\ 150, & \text{if goal is scored.} \end{cases} \quad (12)$$

with $\cos \theta = \langle \frac{\mathbf{p}_{puck} - \mathbf{p}_{ee}}{\|\mathbf{p}_{puck} - \mathbf{p}_{ee}\|}, \frac{\mathbf{p}_{goal} - \mathbf{p}_{puck}}{\|\mathbf{p}_{goal} - \mathbf{p}_{puck}\|} \rangle$, $\mathbf{p}_{ee}, \mathbf{p}_{puck}, \mathbf{p}_{goal}$ is the 2D position of the end-effector, the puck and the goal, r_{hit} is the reward when hitting occurs, $v_{x,hit}$ is the velocity of the puck along x -direction (longitude direction of the table) at the hitting moment, $\lambda = 0.001 \cdot \|\mathbf{a}\|$ is the penalty of the action. If the robot has not hit the puck, the reward encourages to get the end-effector close to the puck along the direction from the puck to the goal. If the hitting occurs, the reward is only influenced by the last reward before hitting and the velocity along the longitude direction, as further actions will not affect the final behavior. When the agent scores a goal, we provide a bonus reward and terminate the episode.

B.3 IiwaAirHockey

In this environment, we learn the robot air hockey task with a 7 DoF KUKA IIWA. To increase the reachability of the robot arm, we design a new end-effector. The end-effector is composed of an extension rod, a gas spring, a universal joint, and a mallet. The total length is 0.5m. To ensure that the mallet stays perpendicular to the table surface, we add a separate controller on the 7th joint of the IIWA robot, which forces the axis of the universal joint to be parallel to the table surface. We disabled the collision of the robot with the table to verify that the constraint is actively guaranteed. Therefore, the universal joint is not able to adapt its joint positions passively. We add a position controller in the simulator to ensure the mallet makes proper contact with the table. The control action is a six-dimensional vector representing the torque applied at each joint. As the real-world’s KUKA controller enforces the joint velocity limits by default, we force the joint velocity constraints in the simulation.

In addition to the inequality constraints described in the *PlanarAirHockey*, we add an equality constraint to ensure the end-effector is moving on the table surface, two inequality constraints on the 4th link and the 6th link in order to prevent the collision. The final constraint set is

$$\begin{aligned} f : \quad z_{ee} - z_{table} &= 0, & g_1 : \quad -z_4 + \hat{z}_{4,l} < 0, & g_2 : \quad -z_6 + \hat{z}_{6,l} < 0, \\ g_3 : \quad -x_{ee} + x_{table,l} &< 0, & g_4 : \quad -y_{ee} + y_{table,l} < 0, & g_5 : y_{ee} - y_{table,u} < 0, \\ g_{6,7,\dots,11} : \quad q_i^2 - q_{i,l}^2 &< 0, & \forall i \in \{1, 2, \dots, 6\} \end{aligned}$$

where f ensures the end-effector’s height z_{ee} is the same as table’s height z_{table} . g_1, g_2 constraint the height of the 4-th link z_4 and the 6-th link z_6 above its limit $\hat{z}_{4,l}$ and $\hat{z}_{6,l}$, respectively. g_3, g_4, g_5 ensure the end-effector is moving inside the table’s range. $g_{6,7,\dots,11}$ are the joint position limits constraints. In the *IiwaHit* task, the reward is same as *PlanarHit* specified in (12).

C Additional Experiments

C.1 CircularMotion

In this section, we compare ATACOM with the *ErrorCorrection* and the *Terminated* in the *CircularMotion* environment. The *ErrorCorrection* approach adds the acceleration calculated by (10) at each time step and the *Terminated* approach terminate the episode when the constraint violations bigger than a threshold. The details of the experiment parameter can be found in Appendix D.

Figure 12 illustrated comparison of different approaches in each RL algorithms. We select the best performed parameter for each algorithm after the parameter sweep. Here we plot the total return instead of discounted return, as the difference will be more evident. It is clear that the ATACOM has better performance on each algorithms and less constraint violations.

C.2 PlanarDefend

In this section, we demonstrate the the task *PlanarDefend*. In this task, the puck and the robot is initialized at a certain position and with the same initial velocity. The constraints and environment parameters are same as the hitting task. The ultimate objective is to stop the puck's at the line $x = -0.6$. The reward is designed as following

$$r = \begin{cases} \exp(-3\|\mathbf{p}_{des} - \mathbf{p}_{ee}\|), & \text{if no short sides of the table has been hit} \\ & \text{and the puck has not been hit,} \\ 1 + \exp(-5|x_{puck} + 0.6|) + 5 \exp(-5\|\mathbf{v}_{puck}\|), & \text{if no short sides of the table has been hit} \\ & \text{and the puck has been hit,} \\ 0, & \text{if any short sides of the table has been hit,} \\ -50, & \text{if the puck is in the defender goal.} \end{cases}$$

Figure 13 shows the comparison of the five RL algorithms. In this task, the policy is very hard to obtain. Small changes of the end-effector's position will result in complete different movements on the puck. The constraint violations remains small as expected.

C.3 IiwaAirHockey

In this subsection, we test the *IiwaAirHockey* task with five RL algorithms (DDPG, TD3, SAC, TRPO, PPO). In Figure 14 we present the results of the best parameter of each algorithms after the parameter sweep. These results show that all algorithms are able to improve the policy while maintaining the constraint violation low.

C.4 Comparison with Riemannian Motion Policies

In the following section, we compare our method with Riemannian Motion Policies (RMP) [37] in the *IiwaAirHockey* environment. We follow the dynamics of the collision avoidance as described in [26]. The collision avoidance policy is used to avoid the collision between the End-Effector and the table's boundary, the wrist link and the table, as well as the elbow link and the table. In addition, we also apply the joint limits avoidance policy described in [37]. Unfortunately, RMP does not consider equality constraints. Therefore, to maintain the end-effector on the table surface, we deploy a simple PD-controller to maintain the height of the end-effector z_{ee} at the table's height z_{table} as:

$$\ddot{z}_{ee} = -P(z_{ee} - z_{table}) - D\dot{z}_{ee}$$

with P, D the controller gain. Finally, we add a learning policy on top of the composed RMP policy. The action space of the learning agent is the joint acceleration in each dimension.

Figure 15 shows the comparison between RMP and ATACOM. ATACOM can learn faster since it explores only on the constraint manifold, instead of on the whole joint space as in RMP. ATACOM keeps the constraint violations small during the whole learning process, differently from the RMP approach. These higher violations are due to the fact that the potential fields of RMP model soft constraints and not hard ones. While the final performance of RMP is slightly better than the one achieved by ATACOM, the learned trajectories violate the constraints, allowing for faster movements that are not possible on real robots.

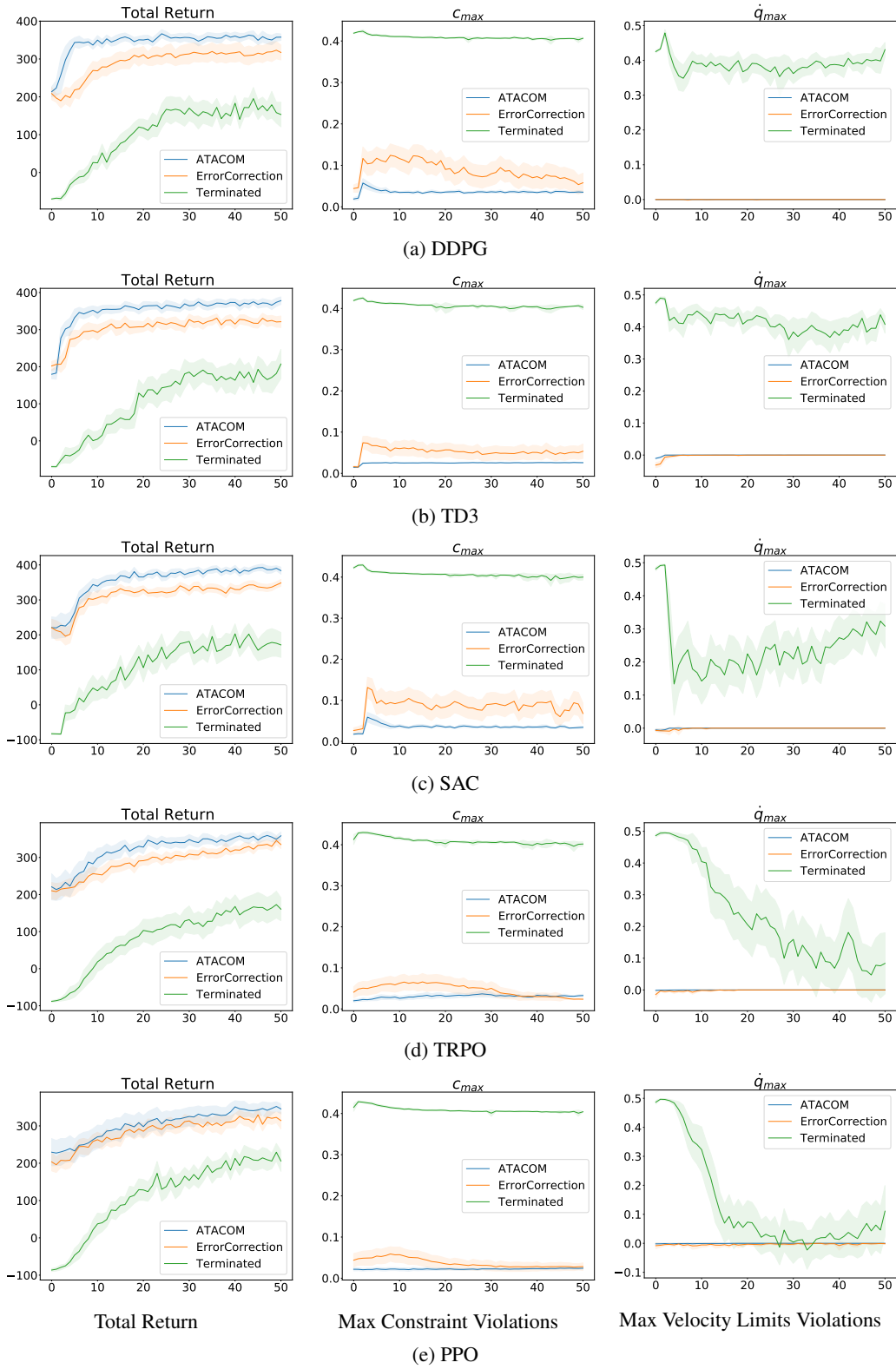


Figure 12: Comparison of ATACOM, ErrorCorrection and Terminated in CircularMotion.

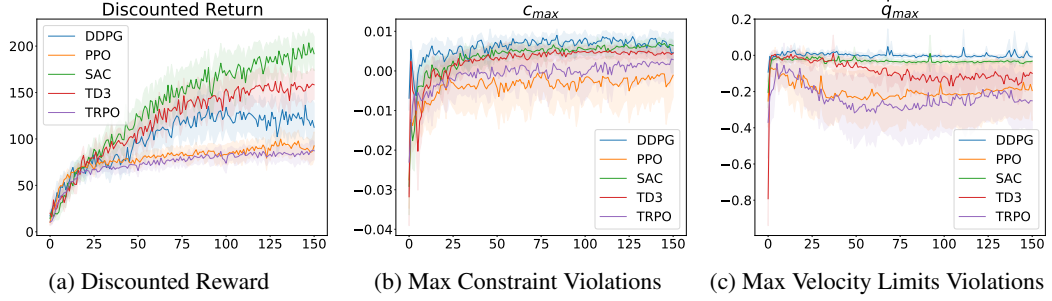


Figure 13: Comparison of RL algorithms in *PlanarDefend*.

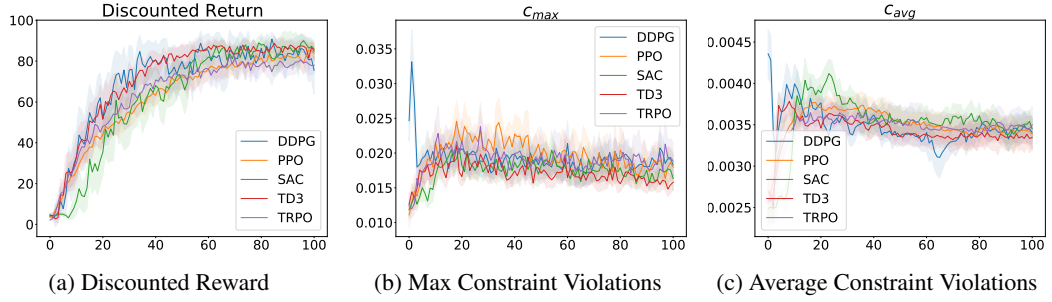


Figure 14: Comparison of RL algorithms in *IiwaAirHockey*.

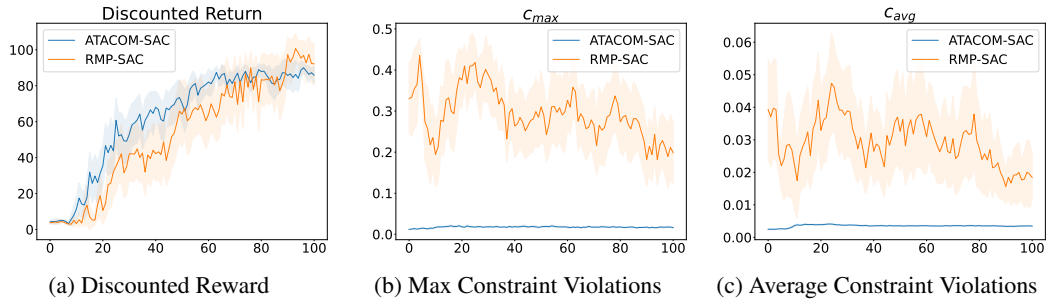


Figure 15: Comparison of RMP with ATACOM in *IiwaAirHockey*.

D Parameters Sweep and Results

In this section, we listed all the parameters and the results of the parameter sweep of each experiment. We launched 25 independent runs for each sweeping parameter.

D.1 CircularMotion

We compare ATACOM with the *ErrorCorrection* and the *Terminated* approaches. The *ErrorCorrection* approach applies only the error correction (10) at each time step and the *Terminated* approach terminate the environment when the constraint violations bigger than a threshold. The hyperparameters for the experiments are shown in Table 1 and Table 2.

Environment Parameter	ATACOM	ErrorCorrection	Terminated
episode duration	5s	5s	5s
discount factor	0.99	0.99	0.99
simulation step size	0.01s	0.01s	0.01s
acceleration limit a_{\max}	10	10	10
K_c in (10)	diag([100, 100])	diag([100, 100])	-
K_f for the equality constraint	diag([0.1])	diag([0.1])	-
K_g for the inequality constraint	diag([2])	diag([2])	-
K_a for the joint accelerations	diag([20, 20])	diag([20, 20])	-
error correction frequency	100	100	-
termination tolerance	-	-	0.4
action space	1D tangent space	2D acceleration	2D acceleration

Table 1: Parameters for *CircularMotion* Environment

	DDPG	TD3	SAC	TRPO	PPO
Sweeping parameter					
actor/critic learning rate	{ $1e^{-3}$, $5e^{-4}$, $1e^{-4}$ }			-	-
clipping coefficient	-	-	-	{0.1, 0.05, 0.01}	-
maximum kl	-	-	-	-	{0.2, 0.1, 0.05}
Default parameter					
epochs	50	50	50	50	50
steps per epoch	5000	5000	5000	5000	5000
steps per fit	1	1	1	2000	2000
episodes per test	25	25	25	25	25
actor/critic network size	[80 80]	[80 80]	[32 32]	[32 32]	[32 32]
batch size	64	64	64	64	64
initial policy covariance	0.2	1.0	-	0.25	0.25
initial replay size	5000	5000	5000	-	-
max replay size	200000	200000	200000	-	-
soft updates coefficient	$1e^{-3}$	$1e^{-3}$	$1e^{-3}$	-	-
warm-up transitions	-	-	10000	-	-
learning rate alpha	-	-	0.0003	-	-
target entropy	-	-	-6	-	-
GAE* update coefficient	-	-	-	0.95	0.95
entropy regularisation	-	-	-	5e-5	5e-5
line searches per fit	-	-	-	10	-
CG** steps per fit	-	-	-	10	-
CG damping	-	-	-	$1e^{-2}$	-
CG tolerance	-	-	-	$1e^{-10}$	-

Table 2: Training Parameter for Algorithms. GAE*: Generalized Advantage Estimation. CG**: Conjugate Gradient

D.2 PlanarAirHockey

For this task, we have demonstrated two tasks, *PlanarHit* and *PlanarDefend*. The Parameter choice of the environment and RL algorithms are shown in Table 3 and Table 4

Environment Parameter	<i>PlanarHit</i>	<i>PlanarDefend</i>
episode duration	2s	3s
discount factor		0.99
simulation step size		1 / 240s
acceleration limit \mathbf{a}_{\max}		[10, 10, 10]
velocity limit \mathbf{v}_{\max}		[2.3562, 2.3562, 2.3562]
\mathbf{K}_c in (10)		diag([240])
\mathbf{K}_f for the equality constraint		-
\mathbf{K}_g for the inequality constraint		diag([0.5, 0.5, 0.5, 1, 1, 1])
\mathbf{K}_a for the joint accelerations		diag($2 \cdot \mathbf{a}_{\max} / \mathbf{v}_{\max}$)
error correction frequency		240
control frequency		60
maximum simulated joint velocity		$1.5 \cdot \mathbf{v}_{\max}$

Table 3: Parameters for *PlanarAirHockey* Environment

	DDPG	TD3	SAC	TRPO	PPO
Sweeping parameter					
actor/critic learning rate	{ $1e^{-3}$, $5e^{-4}$, $1e^{-4}$ }			-	-
clipping coefficient	-	-	-	{0.1, 0.05, 0.01}	-
maximum kl	-	-	-	-	{0.2, 0.1, 0.05}
Default parameter (<i>PlanarHit/PlanarDefend</i>)					
epochs				50 / 150	
steps per epoch				5000 / 12000	
steps per fit	1	1	1	600	600
episodes per test				25	
actor/critic network size	[80 80]	[80 80]	[64 64]	[64 64]	[64 64]
batch size				64	
initial policy covariance	0.2	1.0	-	0.25	0.25
initial replay size	5000	5000	5000	-	-
max replay size	200000	200000	200000	-	-
soft updates coefficient	$1e^{-3}$	$1e^{-3}$	$1e^{-3}$	-	-
warm-up transitions	-	-	10000	-	-
learning rate alpha	-	-	0.0003	-	-
target entropy	-	-	-6	-	-
GAE* update coefficient	-	-	-	0.95	0.95
entropy regularisation	-	-	-	5e-5	5e-5
line searches per fit	-	-	-	10	-
CG** steps per fit	-	-	-	10	-
CG damping	-	-	-	$1e^{-2}$	-
CG tolerance	-	-	-	$1e^{-10}$	-

Table 4: Training Parameter for *PlanarHit/PlanarDefend*. GAE*: Generalized Advantage Estimation. CG**: Conjugate Gradient

D.3 IiwaAirHockey

Environment Parameter	<i>IiwaAirHockey</i>
Sweeping Parameter (SAC)	
simulation step size	[1/50s, 1/250s, 1/500s, 1/1000s]
Default Parameter	
episode duration	2s
discount factor	0.99
acceleration limit \mathbf{a}_{\max}	[10, 10, 10, 10, 10, 10]
velocity limit \mathbf{v}_{\max}	[1.4835, 1.4835, 1.7453, 1.3090, 2.2689, 2.3562]
\mathbf{K}_c in (10)	diag([500])
\mathbf{K}_f for the equality constraint	diag([0.1])
\mathbf{K}_g for the inequality constraint	diag([0.5, 0.5, 0.5, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0])
\mathbf{K}_a for the joint accelerations	-
control frequency	50
maximum simulated joint velocity	\mathbf{v}_{\max}

Table 5: Parameters for *IiwaAirHockey* Environment

	DDPG	TD3	SAC	TRPO	PPO
Sweeping parameter					
actor/critic learning rate	$\{1e^{-3}, 5e^{-4}, 1e^{-4}\}$			-	-
clipping coefficient	-	-	-	{0.1, 0.05, 0.01}	-
maximum kl	-	-	-	-	{0.2, 0.1, 0.05}
Default parameter (<i>PlanrHit/PlanrDefend</i>)					
epochs				50 / 150	
steps per epoch				5000 / 12000	
steps per fit	1	1	1	600	600
episodes per test				25	
actor/critic network size	[80 80]	[80 80]	[64 64]	[64 64]	[64 64]
batch size				64	
initial policy covariance	0.2	1.0	-	0.25	0.25
initial replay size	5000	5000	5000	-	-
max replay size	200000	200000	200000	-	-
soft updates coefficient	$1e^{-3}$	$1e^{-3}$	$1e^{-3}$	-	-
warm-up transitions	-	-	10000	-	-
learning rate alpha	-	-	0.0003	-	-
target entropy	-	-	-6	-	-
GAE* update coefficient	-	-	-	0.95	0.95
entropy regularisation	-	-	-	5e-5	5e-5
line searches per fit	-	-	-	10	-
CG** steps per fit	-	-	-	10	-
CG damping	-	-	-	$1e^{-2}$	-
CG tolerance	-	-	-	$1e^{-10}$	-

Table 6: Training Parameter for *IiwaAirHockey*. GAE*: Generalized Advantage Estimation. CG**: Conjugate Gradient

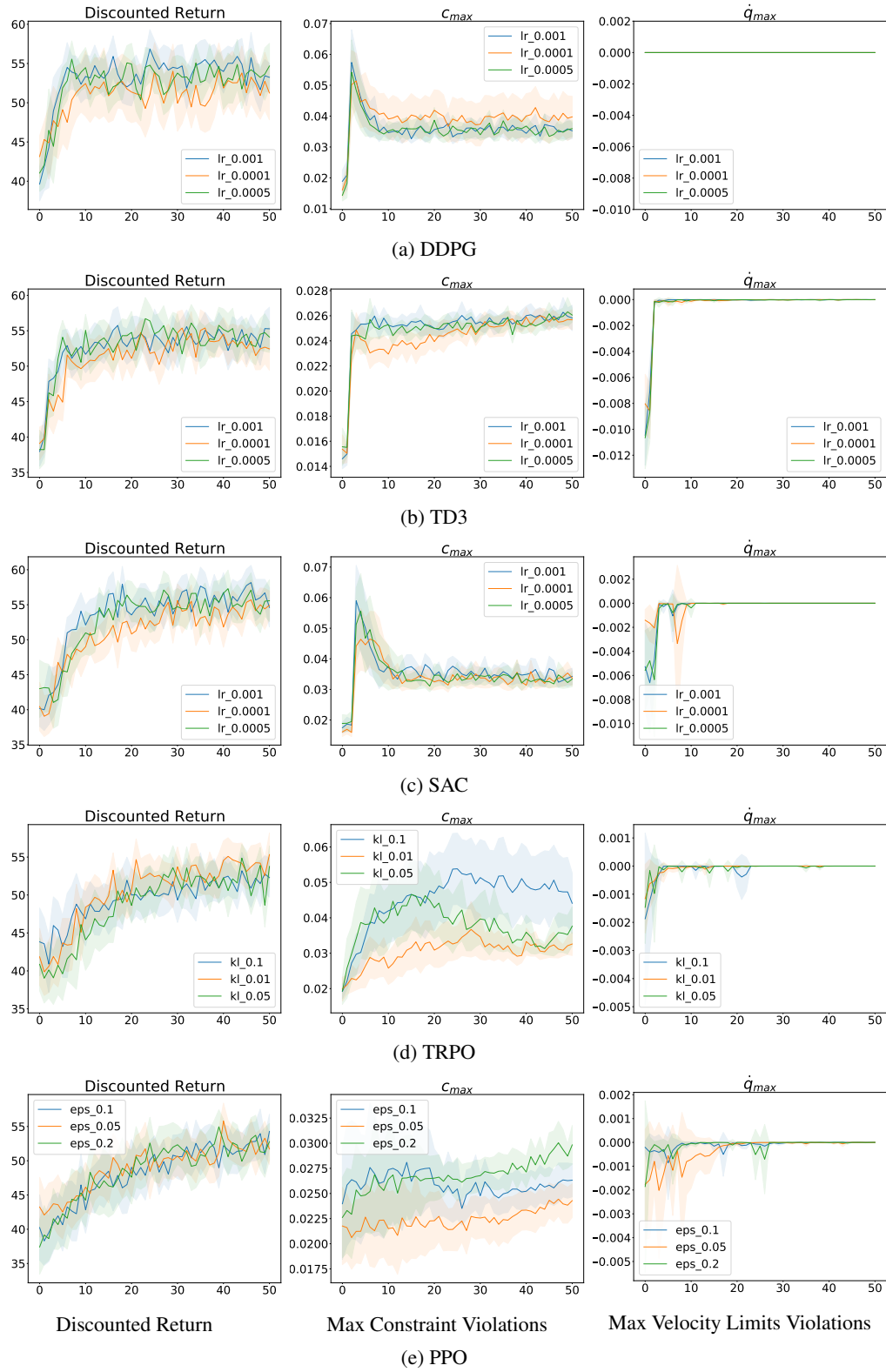


Figure 16: Parameter sweep of ATACOM in *CircularMotion*.

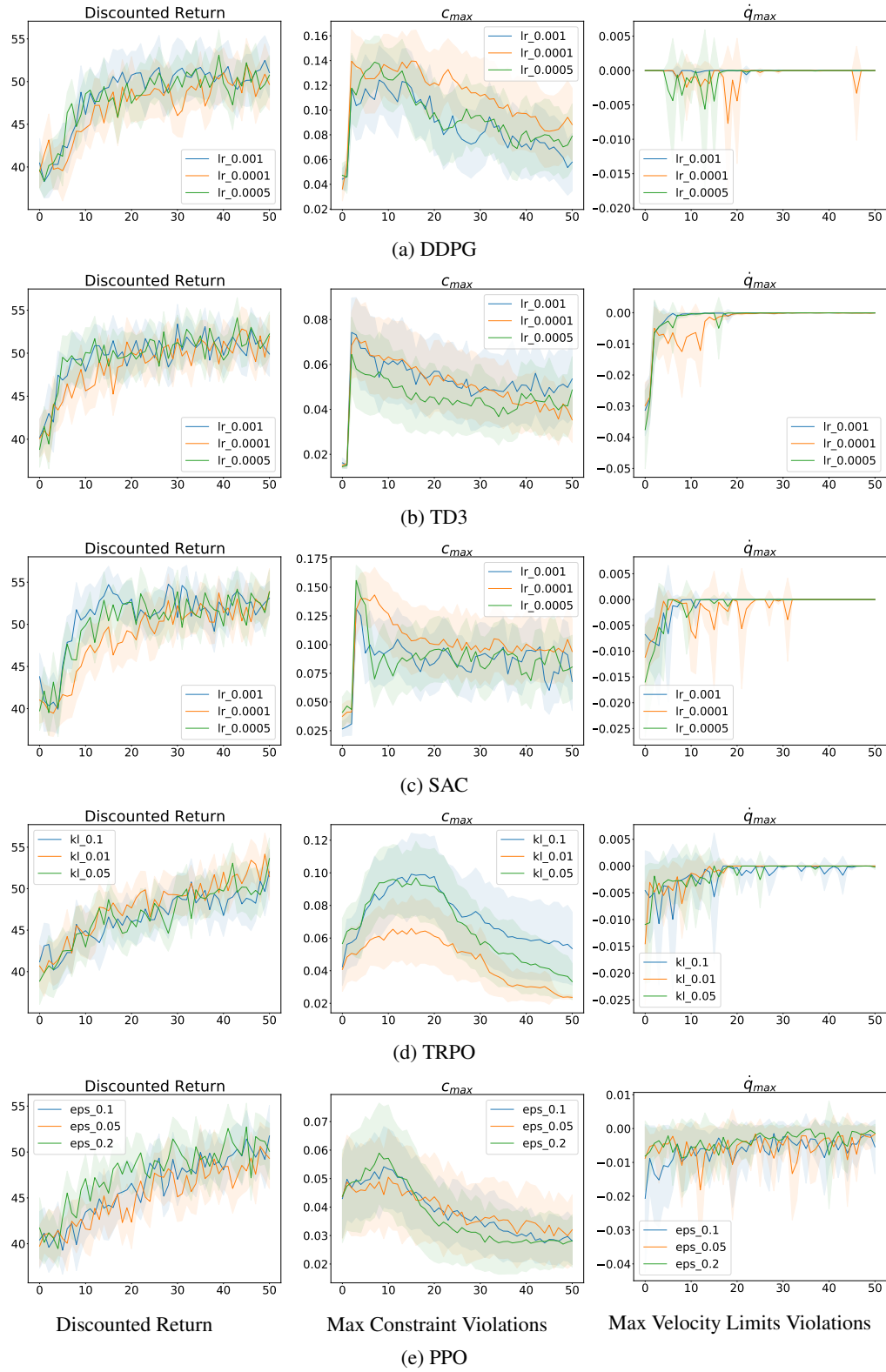


Figure 17: Parameter sweep of *ErrorCorrection* in *CircularMotion*.

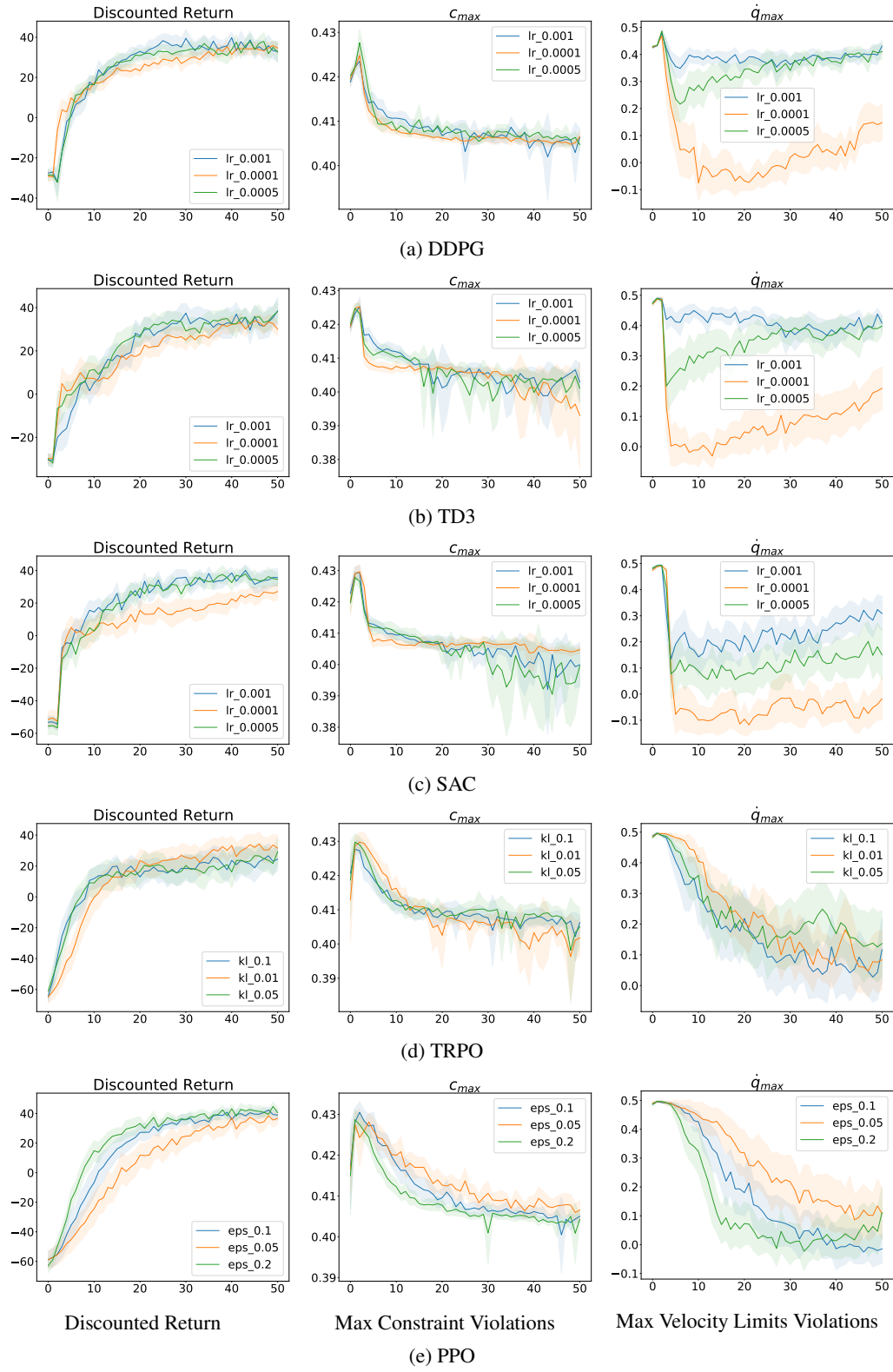


Figure 18: Parameter sweep of *Terminated* in *CircularMotion*.

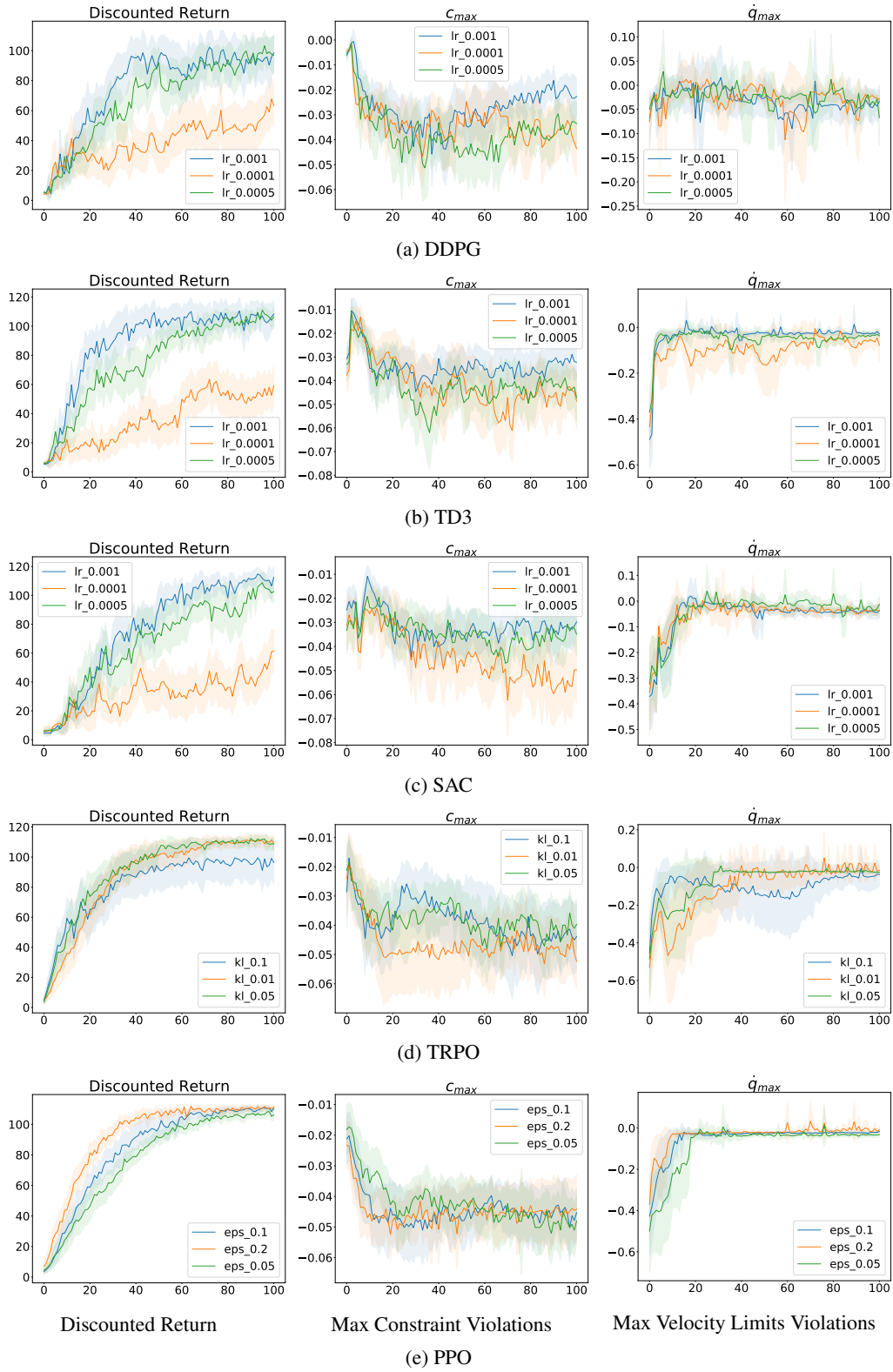


Figure 19: Parameter sweep of ATACOM in *PlanarHit*.

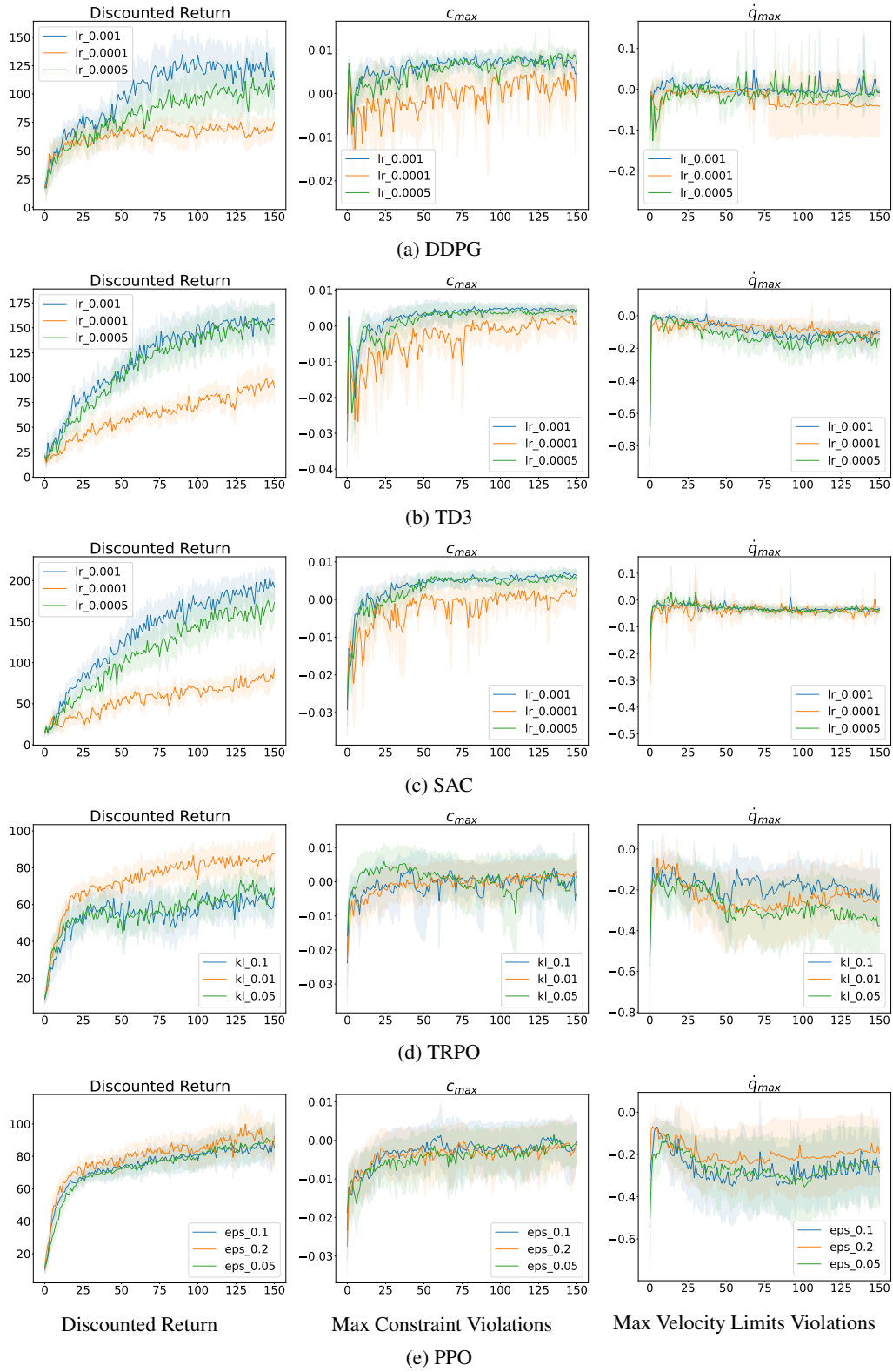


Figure 20: Parameter sweep of ATACOM in *PlanarDefend*.

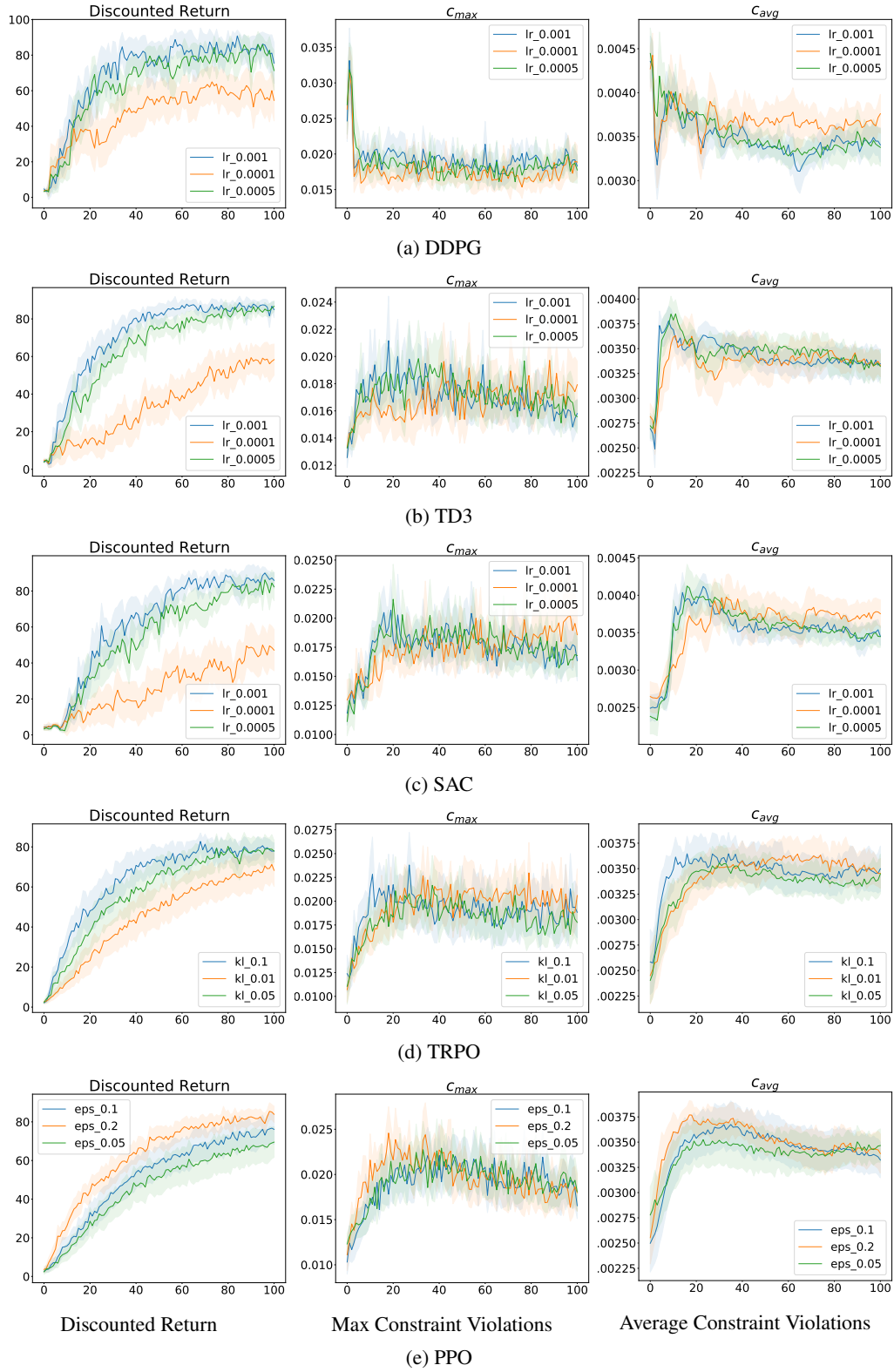


Figure 21: Parameter sweep of ATACOM in *IiwaAirHockey*.

E Constraint with Uncontrollable State

The proposed method ATACOM can be easily extended to the constraints with uncontrollable states. We assume that the velocity of uncontrollable state \dot{x} can be estimated and the acceleration of controllable state $\ddot{x} = 0$. Here we briefly introduce the extension of ATACOM of the inequality constraints with $g(q, x) \leq 0$. The viability constraint set is

$$c(q, \dot{q}, x, \dot{x}, \mu) = g(q, x) + K(J_q(q, x)\dot{q} + J_x(q, x)\dot{x}) + \frac{1}{2}\mu^2 = 0$$

with the partial derivatives $J_q(q, x) = \frac{\partial}{\partial q}g(q, x)$ and $J_x(q, x) = \frac{\partial}{\partial x}g(q, x)$. We use J_q, J_x to simplify the notation. The time derivative of the viability constraint is

$$\dot{c}(q, \dot{q}, \ddot{q}, x, \dot{x}, \ddot{x}, \mu, \dot{\mu}) = J_q\dot{q} + J_x\dot{x} + KJ_q\ddot{q} + K\dot{J}_q\dot{q} + KJ_x\ddot{x} + K\dot{J}_x\dot{x} + \text{diag}(\mu)\dot{\mu}$$

with the time derivative of i -th Jacobian $\dot{J}_{i,q} = \dot{q}^\top H_{i,qq} + \dot{x}^\top H_{i,xq}$ and $\dot{J}_{i,x} = \dot{x}^\top H_{i,xx} + \dot{q}^\top H_{i,qx}$. H_i is the hessian matrix w.r.t i -th constraint.

As mentioned before, we assume \dot{x} is known and $\ddot{x} = 0$. By setting $\dot{c} = 0$, we have

$$\underbrace{[KJ_q \quad \text{diag}(\mu)]}_{J_c} \begin{bmatrix} \ddot{q} \\ \dot{\mu} \end{bmatrix} + \underbrace{J_q\dot{q} + J_x\dot{x} + K\dot{J}_q\dot{q} + K\dot{J}_x\dot{x}}_{\psi} = 0$$

The overall control action is

$$\begin{bmatrix} \ddot{q} \\ \dot{\mu} \end{bmatrix} = -J_c^\dagger(\psi + K_c c) + N_c \alpha$$

with the error correction gain K_c . We validate our method in the *CollisionAvoidance* environment which is described in Appendix E.1.

E.1 CollisionAvoidance

In this experiment, we demonstrate a collision avoidance task with four moving obstacles in a 2d environment shown in Figure 22. The objective is to move the agent (blue circle) to the target (green square) while avoiding the collision with four random moving obstacles (red circle). In this environment, the velocities of the obstacles are known. The agent and obstacles have a radius of 0.3. The control action is the acceleration along x, y direction. The collision avoidance constraint is

$$g_i : 0.6^2 - (q_x - o_{i,x})^2 - (q_y - o_{i,y})^2 < 0, \quad i \in \{1, 2, 3, 4\}.$$

q_x and q_y are agent's positions and $o_{i,x}, o_{i,y}$ are positions of i -th object. The reward function is

$$r(q_x, q_y) = -\frac{1}{c} \|\mathbf{q}_{goal} - \mathbf{q}\|.$$

with a scale constant c .

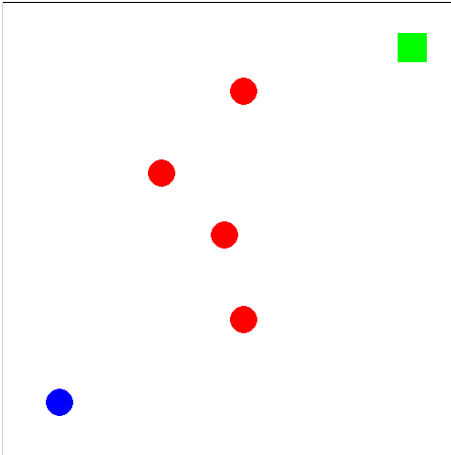


Figure 22: Collision Avoidance Environment

	SAC
default parameter	
actor/critic learning rate	$3e^{-4}$
epochs	100
steps per epoch	10000
steps per fit	11
episodes per test	25
actor/critic network size	[64 64]
batch size	64
initial policy covariance	-
initial replay size	5000
max replay size	200000
soft updates coefficient	$1e^{-3}$
warm-up transitions	10000
learning rate alpha	0.0003
target entropy	-4

Table 7: Training Parameter for Collision Avoidance

E.2 Experiment of Collision Avoidance

In this section, we demonstrate the preliminary result of collision avoidance. We applied the SAC with default parameter setup as shown in Table 7. As shown in Figure 23, the agent learns to reach the goal while the maximum constraint violations throughout the learning process remain low. Small constraint violations during the learning process occur if the agent is surrounded by the obstacles in the corner, and there exists no feasible action to avoid the collision.

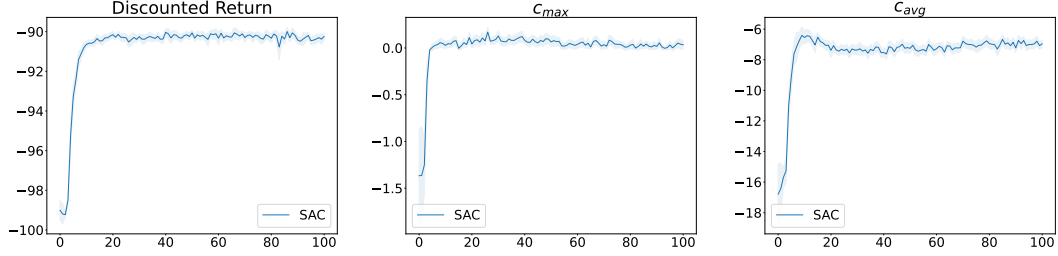


Figure 23: Learning Curve of ATACOM-SAC in *CollisionAvoidance* Environment