# Model-based Chance-Constrained Reinforcement Learning via Separated Proportional-Integral Lagrangian

Baiyu Peng, Jingliang Duan, Jianyu Chen, Shengbo Eben Li, Genjin Xie, Congsheng Zhang, Yang Guan, Yao Mu, Enxin Sun

*Abstract*—**Safety is essential for reinforcement learning (RL) applied in the real world. Adding chance constraints (or probabilistic constraints) is a suitable way to enhance RL safety under uncertainty. Existing chance-constrained RL methods like the penalty methods and the Lagrangian methods either exhibit periodic oscillations or learn an over-conservative or unsafe policy. In this paper, we address these shortcomings by proposing a separated proportional-integral Lagrangian (SPIL) algorithm. We first review the constrained policy optimization process from a feedback control perspective, which regards the penalty weight as the control input and the safe probability as the control output. Based on this, the penalty method is formulated as a proportional controller, and the Lagrangian method is formulated as an integral controller. We then unify them and present a proportional-integral Lagrangian method to get both their merits, with an integral separation technique to limit the integral value in a reasonable range. To accelerate training, the gradient of safe probability is computed in a model-based manner. We demonstrate our method can reduce the oscillations and conservatism of RL policy in a car-following simulation. To prove its practicality, we also apply our method to a real-world mobile robot navigation task, where our robot successfully avoids a moving obstacle with highly uncertain or even aggressive behaviors.**

*Index Terms*—**Safe reinforcement learning, constrained control, robot navigation, neural networks.**

## I. INTRODUCTION

RECENT advances in deep reinforcement learning (RL) have demonstrated state-of-the-art performance in a variety of tasks, including video games [1]–[3], autonomous driving [4], [5] and robotics [6], [7]. However, most RL successes still remain in virtual environments or simulation platforms. For safety-critical real-world tasks, RL is not yet fully mature or ready to serve as an "off-the-shelf" solution. One of the reasons is the lack of safety constraints [8]. Consequently, how to handle safety constraints has become a popular and essential topic in RL community.

Baiyu Peng, Jingliang Duan, Shengbo Eben Li, Genjin Xie, Congsheng Zhang, Yang Guan, Yao Mu and Enxin Sun are with the State Key Laboratory of Automotive Safety and Energy, School of Vehicle and Mobility, Tsinghua University, Beijing, China (e-mail: lishbo@tsinghua.edu.cn).

Jianyu Chen is with the Institute for Interdisciplinary Information Sciences, Tsinghua University, Beijing, China and Shanghai Qi Zhi Institute, Shanghai, China.

The safety constraints used in safe RL mainly fall into three categories: expected constraints, worst-case constraints, and chance constraints. Especially, the popular Constrained Markov Decision Process (CMDP) framework [9] is a special case of the expected constraints, which constrains the expected cumulative cost to be below a predetermined boundary. Many well-known safe RL algorithms build on this framework, including Constrained Policy Optimization (CPO) [10], Projection-based Constrained Policy Optimization (PCPO) [11], and Reward Constrained Policy Optimization (RCPO) [12]. However, these methods only guarantee constraint satisfaction in expectation, which is inadequate for safety-critical engineering applications. In this case, the probability of the constraint violation is about 50% (roughly speaking) [13]. The second type of constraint is the worst-case constraint, which guarantees constraint satisfaction under any uncertain conditions. Nevertheless, the worst-case constraint tends to be overly conservative, and it only supports systems with bounded noise [8]. The third form is the chance constraint, where the constraint holds with a predefined probability. The chance constraint clearly limits the occurring probability of the unsafe event, which is selected according to the different demands of users and the tasks. Therefore, the chance constraint is quite suitable for various real-world applications. In this paper, we will focus on the chance-constrained RL problems, i.e., how to learn an optimal policy satisfying the chance constraints.

Next, we will briefly introduce the existing chance-constrained RL studies. In 2005, Geibel and Wysotzki use an indicator function to estimate the safe probability by sampling, and add a large penalty term in the reward function if the safe probability is low [14]. Then, the reshaped reward is optimized by an actor-critic method. Giuseppi and Pietrabissa (2020) view the reward and safe probability as two objectives, and propose a corresponding multi-objective RL method [15]. Paternain et al. (2019) derive a lower bound of the safe probability, which is employed to construct a surrogate constraint since it has an additive structure and easier to tackle [16]. Then, the transformed problem is solved by the Lagrangian method, which introduces a Lagrangian multiplier to balance policy performance and constraints satisfaction. To reduce the conservatism introduced by constraining a lower bound, Peng et al. (2020) only use the lower bound to obtain an update direction, but still evaluate the feasibility of the policy using original chance constraints [17]. In addition, they employ a penalty method with increasing weight to enforce

constraint satisfaction.

The previous chance-constrained RL mainly relies on the penalty method or the Lagrangian method. However, they actually face several challenges such as poor policy performance, constraint violations, and unstable learning process [18]. The penalty methods require a well-designed penalty weight to balance the reward and constraint, which unfortunately is non-trivial and hard to tune. As shown in Fig. 1(a), a large penalty is prone to rapid oscillations and frequent constraint violations, while a small penalty always violates the constraint seriously [12]. As for the Lagrangian method, it usually suffers from the Lagrange multiplier overshooting (see Fig. 1(b)), which will lead to an overly conservative policy [16], [19]. Besides, due to the delay between the policy optimization and the Lagrange multiplier adaptation, the Lagrangian multiplier usually oscillates periodically during the training process, which further results in policy oscillations [20], [21]. In addition, most existing methods do not support model-based optimization since they all rely on a non-differentiable indicator function to estimate the safe probability. Therefore, they can only use model-free methods to optimize the safe probability, which is generally slower than model-based methods.
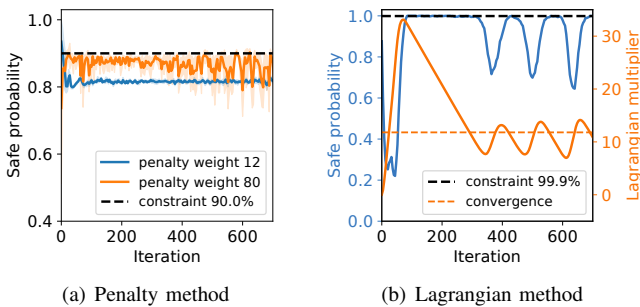


Fig. 1. Examples of learning curves for penalty and Lagrangian methods. (a) Penalty method exhibits oscillations and violates the constraint. (b) Lagrangian method exhibits Lagrange multiplier overshooting and oscillations, and further harms the policy learning.

To overcome the problems mentioned above, we propose a model-based Separated Proportional-Integral Lagrangian (SPIL) method for chance-constrained RL, which can fulfill the safety requirements with a steady and fast learning process. The contributions of this paper are summarized as follows,

1) This paper interprets the chance-constrained policy optimization process from a feedback control perspective, which regards the penalty weight as the control input and the safe probability as the control output. Based on this, we unify the existing constraint optimization methods into the PID control methodology, in which the penalty method is formulated as a proportional controller and the Lagrangian method is formulated as an integral controller. Then, we develop the proportional-integral Lagrangian (PIL) framework by combining the proportional and integral modules to get their both merits.

2) To prevent policy from over-conservatism caused by multiplier overshooting, we draw inspiration from PID control and introduce an integral separation technique, which separates the integrator out when the integral value exceeds a predetermined threshold. Then, we embed this technology into the PIL framework to propose the SPIL method for chance-constrained RL. Simulations demonstrate our SPIL achieves better policy performance compared with existing penalty and Lagrangian methods [14], [16], [19].

3) To accelerate the training process in a model-based way, we adopt an approximated model-based gradient of the safe probability to participate in the policy optimization. The approximated gradient is proved to approach to the true gradient under mild conditions. Compared with a popular model-free method, the learning speed is improved by at least five times.

4) Finally, a real-world mobile robot navigation experiment proves the effectiveness of model-based SPIL algorithm in practical engineering problems.

This article is further organized as follows. The chance-constrained RL problem is formulated in Section II. The model-based SPIL method is proposed in Section III. In Section IV, our method is verified and compared in a car-following simulation. In Section V, the mobile robot navigation experiment is presented to prove the practicability of the algorithm. Section VI concludes this paper.

## II. PRELIMINARIES

### A. Problem Description

In our model-based RL setting, we assume a dynamic model is already available, either learned by collected data or derived by prior knowledge. To indicate the gap between the model and reality, an uncertain term is included in this model. This assumption is reasonable in many physical engineering problems like autonomous driving, where there are many established dynamic models. Therefore, we can directly learn a policy through the given uncertainty model. To ensure the policy feasibility for real environments, a chance constraint needs to be introduced in the model-based learning process. In other words, given a reasonable uncertain term, if the policy is safe with a high probability under model uncertainty, we can usually ensure its safety in practical environments. This setting is similar to that in stochastic control [22].

The dynamic model and the chance constraint are mathematically described as:

$$s_{t+1} = f(s_t, a_t, \xi_t), \quad \xi_t \sim p(\xi_t),$$

$$\Pr\left\{ \bigcap_{t=1}^{N} [h(s_t) < 0] \right\} \geq 1 - \delta \tag{1}$$

where $t$ is the current step, $s_t \in \mathcal{S}$ is the state, $a_t \in \mathcal{A}$ is the action, $f(\cdot, \cdot, \cdot)$ is the environmental dynamic, $\xi_t \in \mathbb{R}^n$ is the model uncertainty following an independent and identical distribution $p(\xi_t)$. $h(\cdot)$ is the safety function defining a safe state region. Here the chance constraint takes a joint form, which is initially brought from stochastic systems control [22]. Intuitively, it requires the probability of being safe over the finite horizon $N$ to be at least $1 - \delta$. For simplicity, we only consider one constraint, but our method can readily generalize to multiple constraints.

The chance-constrained problem is defined as maximizing a objective function $J$, i.e., the expected discounted cumulative reward, while keeps a high safe probability $p_s$:

$$\max_\pi J(\pi) = \mathbb{E}_{s_0,\xi}\left\{\sum_{t=0}^{\infty}\gamma^t r(s_t, a_t)\right\}$$

$$\text{s.t. } p_s(\pi) = \Pr\left\{\bigcap_{t=1}^{N}[h(s_t) < 0]\right\} \geq 1 - \delta, \quad (2)$$

$$s_{t+1} = f(s_t, a_t, \xi_t), \quad \xi_t \sim p(\xi_t)$$

where $r(\cdot, \cdot)$ is the reward function, $\gamma \in (0, 1)$ is the discount factor, $\mathbb{E}_{s_0,\xi}(\cdot)$ is the expectation w.r.t. the initial state $s_0$ and uncertainty $\xi_{0:\infty}$. $\pi$ is a deterministic policy, i.e., a mapping from state space $\mathcal{S}$ to action space $\mathcal{A}$. In practice, policy is usually a parameterized neural network with parameters $\theta$, denoted as $\pi(s_t; \theta)$ or $\pi_\theta$.

### B. Penalty and Lagrangian Methods

To find the optimal control policy for problem (2), the penalty and Lagrangian methods are widely employed in existing studies [14]–[17]. The penalty method adds a quadratic penalty term in the objective function to force the satisfaction of the constraint:

$$\max_\pi J(\pi) - \frac{1}{2}K_P\left((1 - \delta - p_s(\pi))^+\right)^2 \quad (3)$$

where $K_p > 0$ is the penalty weight, $p_s(\pi)$ is the joint safe probability (see (2)), and $(\cdot)^+$ means $\max(\cdot, 0)$. This unconstrained problem is usually solved by gradient ascent:

$$\theta^k \leftarrow \theta^{k-1} + \alpha_\theta(\nabla_\theta J^{k-1} + K_P(1 - \delta - p_s^{k-1})^+\nabla_\theta p_s^{k-1}) \quad (4)$$

where $k$ means the $k$-th iteration, $\alpha_\theta > 0$ is the learning rate, $J^k$ and $p_s^k$ are short for $J(\pi_{\theta^k})$ and $p_s(\pi_{\theta^k})$ respectively. For practical applications, it is usually difficult to select an appropriate weight $K_P$ to balance reward and constraint well. A large penalty is prone to rapid oscillations and frequent constraint violations, while a small penalty always violates the constraint seriously.

For the Lagrangian method, it first transforms the chance-constrained problem (2) into an dual problem by introduction of the Lagrange multiplier $\lambda_I$ [23]:

$$\max_{\lambda_I \geq 0}\min_\pi \mathcal{L}(\pi, \lambda_I) = -J(\pi) + \lambda_I(1 - \delta - p_s(\pi)) \quad (5)$$

Then, problem (5) can be solved by dual ascent, i.e., alternatively update the Lagrange multiplier and primal variables:

$$\lambda_I^k \leftarrow (\lambda_I^{k-1} + K_I(1 - \delta - p_s^{k-1}))^+, \quad (6)$$

$$\theta^k \leftarrow \theta^{k-1} + \alpha_\theta(\nabla_\theta J^{k-1} + \lambda_I^k \nabla_\theta p_s^{k-1}) \quad (7)$$

where $K_I > 0$ is the learning rate for $\lambda_I$.

As mentioned in Introduction, the Lagrangian method usually faces the Lagrange multiplier overshooting and multiplier oscillation challenges, resulting in poor policy performance and unstable learning process.

## III. METHODOLOGY

In this section, we first reshape the penalty method and Lagrangian method in a feedback control view. Then we unify them to formulate a proportional-integral Lagrangian method to improve the policy performance without losing the safety requirement. Finally, we introduce an integral separation technique and a model-based gradient to make the whole method practical and efficient.

### A. Feedback Control View of Chance-Constrained Policy Optimization

The key insight of the proposed method comes from a deep and novel understanding of the penalty and Lagrangian methods from a control perspective. From (4) and (7), the update rule of existing methods can be expressed in a similar form

$$\theta^k \leftarrow \theta^{k-1} + \alpha_\theta(\nabla_\theta J^{k-1} + \lambda^k \nabla_\theta p_s^{k-1}) \quad (8)$$

where $\lambda^k$ is actually a balancing weight. The core difference between the two methods lies in the selection of $\lambda$. For the penalty method, $\lambda^k = K_P(1 - \delta - p_s^{k-1})^+$, which indicates the constraint violation at the $k$-th iteration. For the Lagrangian method, $\lambda^k = \lambda_I^k$ in (6), which can be regarded as the sum of constraint violation over previous $k$ iterations. This insight inspires us to review RL from a control perspective.

As shown in Fig. 2, one can view the policy optimization as a feedback control process, where $\theta^k$ is the state, $\lambda^k$ is the control signal, policy update (8) is the system dynamics (state transition equation), $p_s^k$ is the system output, $1 - \delta$ is the desired output and $1 - \delta - p_s^{k-1}$ is the tracking error (or constraint violation). Then, the essence of this system is the design of the controller, i.e., given the tracking error $1 - \delta - p_s$, how can we decide the control signal $\lambda$? The penalty method actually adopts $\lambda^k$ proportional to the violation. The Lagrangian method instead computes $\lambda^k$ as the sum of previous constraint violations. In such an insight, the penalty method becomes exactly a proportional (P) controller, and the Lagrangian method becomes an integral (I) controller. Subsequently, one can easily understand the merits and flaws of these two methods by analogy. For pure "P" control, small $K_P$ leads to steady-state error, while large $K_P$ exhibits oscillation. This matches the phenomenons we observe in the penalty method. Similarly, the problems of overshooting and oscillation in Lagrangian method can also be explained by properties of pure "I" control.
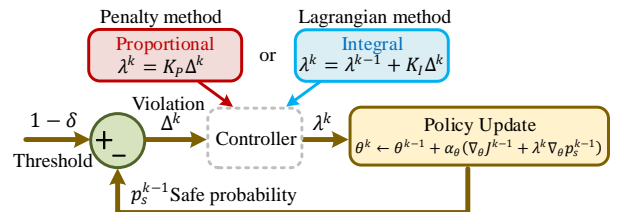
Fig. 2. Feedback control view of chance-constrained policy optimization.

With this insight, we naturally propose to combine the penalty method and Lagrangian method by computing $\lambda^k$ as

a weighted sum of proportional and integral values, which leads to the proportional-integral Lagrangian method (PIL). The update process of PIL at $k-$iteration is given as:

$$\Delta^k \leftarrow 1 - \delta - p_s^{k-1}, \tag{9}$$

$$I^k \leftarrow (I^{k-1} + \Delta^k)^+, \tag{10}$$

$$\lambda^k \leftarrow (K_P \Delta^k + K_I I^k)^+, \tag{11}$$

$$\theta^k \leftarrow \theta^{k-1} + \frac{\alpha_\theta}{1 + \lambda^k}(\nabla_\theta J^{k-1} + \lambda^k \nabla_\theta p_s^{k-1}) \tag{12}$$

where $\Delta^k$ and $I^k$ are proportional and integral values, respectively, with $K_P$ and $K_I$ as coefficients. The proportional term $K_P \Delta^k$ serves as an immediate feedback of the constraint violation. The integral term $K_I I^k$ eliminates the steady-state error at convergence. In such a framework, the penalty method and the Lagrangian method can be regarded as two special cases of PIL with $K_P > 0, K_I = 0$ and $K_P = 0, K_I > 0$, respectively. Note that, to maintain a relatively consistent step size and make policy update more stable, the gradient for $\theta^k$ in (12) is re-scaled by $\frac{1}{1+\lambda^k}$.

This mechanism is expected to realize a steady learning process, just like how the proportional-integral controller works. However, there are still two key issues for practical applications:

1) Integral value $I$ easily gets overly large, leading to policy over-conservatism.
2) $\nabla_\theta p_s^k$ is hard to compute, especially in a model-based paradigm.

The following subsections will explain and solve these two problems.

### B. Integral Separation Technique

The integral value $I^k$ increases according to the constraint violation $\Delta^k$. However, when the initial policy is relatively unsafe, $\Delta^k$ can be quite large, which will cause the overshooting of $I^k$ and $\lambda^k$. With a large $\lambda^k$ in (12), the policy tends to be overly conservative since the weight of $\nabla_\theta p_s^{k-1}$ is overly large. Even worse, since the maximal safe probability is 1, the overshooting and conservatism problems can hardly recover by themselves. For e.g., suppose $1 - \delta = 0.999, p_s^k = 1$, and $\lambda^k$ is already overshooting, the integral term $I^k$ can only fall slowly with the speed of $\Delta^k = -0.001$. Therefore, the policy in such a case will deteriorate in a long time. This issue is also not well recognized and resolved in previous similar works like [21].

To deal with the overshooting problem of $I^k$ and $\lambda^k$, we draw inspiration from PID control [24] and introduce an integral separation technique. As shown in Fig. 3, the integrator will be activated only when $\Delta^k$ is less than a certain value. If $\Delta^k$ is too large, the integrator will be blocked to limit the increase of $I^k$. Specifically, (10) is modified to:

$$I^k \leftarrow (I^{k-1} + K_S \Delta^k)^+,$$

$$K_S = \begin{cases} 0 & \varepsilon_1 < \Delta^k \\ \beta & \varepsilon_2 < \Delta^k \le \varepsilon_1 \\ 1 & \Delta^k \le \varepsilon_2 \end{cases} \tag{13}$$

where $K_S$ is the separation function, $1 > \beta > 0$ and $\varepsilon_1 > \varepsilon_2 > 0$ are predetermined parameters. The piece-wise function $K_S$ separates the integrator out or slows it down if the error is relatively large. Such a recipe restrains the occurrence of overshooting and over-conservatism. Our simulation indicates it greatly improves the performance for a large safety threshold $1 - \delta$, such as $1 - \delta = 99.9\%$. We refer to the combination of PIL and the separation technique as the separated proportional-integral Lagrangian (SPIL) method.
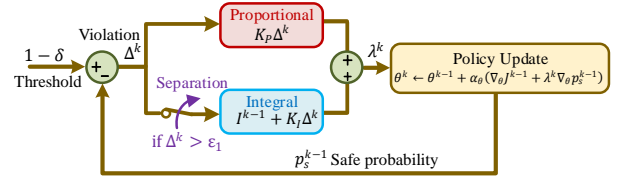


Fig. 3. The framework of the proposed SPIL method.

### C. Model-based Gradient of Safe Probability

Next, we will figure out how to estimate $\nabla_\theta p_s$ in (12) in an efficient way. Generally, there is not any analytic solution for a joint probability and its gradient [22], [25], i.e., they are nearly intractable. Therefore, previous researchers usually introduce an indicator function to estimate the probability through sampling [16], [17]. Due to the discontinuity of the indicator function, these methods are mostly model-free, which are generally believed to be slower than their model-based counterparts [26], [27].

Inspired by recent advances in stochastic optimization [25], we introduce a model-based alternative to $\nabla_\theta p_s$, which enables us to estimate $\nabla_\theta p_s$ efficiently. We first define an indicator-like function $\phi(z, \tau)$:

$$\phi(z, \tau) = \frac{1 + b_1 \tau}{1 + b_2 \tau \exp(-\frac{z}{\tau})}, \tag{14}$$

$$0 < b_2 < \frac{b_1}{1 + b_1}, 0 < \tau < 1$$

where $z$ and $\tau$ are scalar variables of the function, $b_1, b_2$ are the parameters. The expected production of $\phi(z, \tau)$ over $N$ horizon is defined as:

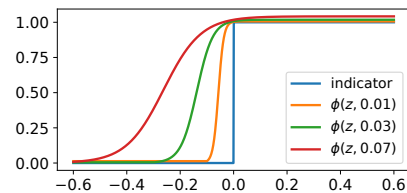$$\Phi(\pi, \tau) = \mathbb{E}_{s_0, \xi} \left\{ \prod_{t=1}^{N} \phi(-h(s_t), \tau) \right\} \tag{15}$$



Fig. 4. Comparison of indicator function and $\phi(z, \tau)$ with different $\tau$.

As shown in Fig. 4, $\phi(z, \tau)$ can be intuitively regarded as a differentiable approximation of indicator function

for constraint violation, and its expected product $\Phi(\pi, \tau)$ approximates joint safe probability. (To see this, one can image $\phi(z, \tau)$ in (15) as the indicator function. Then its expected production is actually the joint safe probability.) The parameter $\tau$ controls how well the indicator function is approximated. Intriguingly, this approximation is not only an intuitive trick and it does have strong theoretical support. Regardless of the nonlinearity and nonconvexity of $h(s_t)$, the gradient of $\Phi(\pi, \tau)$ is proved to converge to the true gradient of joint safe probability $p_s(\pi)$ as $\tau$ approaches 0 under mild assumptions [25]:

$$\lim_{\tau \to 0+} \sup_{\theta \in \Theta} \nabla_\theta \Phi(\pi, \tau) = \nabla_\theta p_s(\pi) \qquad (16)$$

where $\Theta$ is a small ball around the current policy network parameter. For simplicity, we omit mathematical details; interested readers are recommended to refer to [25] for a rigorous explanation.

In practice, one only needs to pick a relatively small fixed $\tau$ and compute $\nabla_\theta \Phi(\pi, \tau)$ to substitute $\nabla_\theta p_s(\pi)$, where the expectation is estimated by sampling average. Therefore, the original policy update rule of SPIL is rewritten as:

$$\theta^k \leftarrow \theta^{k-1} + \frac{\alpha_\theta}{1 + \lambda^k} (\nabla_\theta J^{k-1} + \lambda^k \nabla_\theta \Phi(\pi_{\theta^{k-1}}, \tau)) \quad (17)$$

It should be pointed out that the introduction of $\phi(z, \tau)$ in this subsection is only used for the computation of $\nabla_\theta p_s(\pi)$. The safe probability $p_s$ itself is still estimated through Monte-Carlo sampling, i.e., suppose there are $m$ safe trajectories among all the $M$ trajectories, then the safety probability is $p_s \approx \frac{m}{M}$.

### D. Model-based SPIL Algorithm

Based on aforementioned gradient, we will propose the model-based SPIL algorithm for practical applications. Firstly, we define the state-action value of $(s, a)$ under policy $\pi$ as:

$$Q^\pi(s, a) = \mathbb{E}_{\xi, \pi} \left\{ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \Big| s_0 = s, a_0 = a \right\} \quad (18)$$

Thus, the expected cumulative reward $J$ in (2) can be expressed as a $N$-step form:

$$J(\pi) = \mathbb{E}_{s_0, \xi} \left\{ \sum_{t=0}^{N-1} \gamma^t r(s_t, a_t) + \gamma^N Q^\pi(s_N, s_N) \right\} \quad (19)$$

For large and continuous state spaces, both value function and policy are parameterized:

$$Q(s, a) \cong Q(s, a; w), \quad a \cong \pi(s; \theta) \quad (20)$$

The parameterized state-action value function with parameter $w$ is usually named the "critic", and the parameterized policy with parameter $\theta$ is named the "actor" [27].

The parameterized critic is trained by minimizing the average square error:

$$J_Q^k = \mathbb{E}_{s_0, \xi} \left\{ \frac{1}{2} \left( Q_{\text{target}} - Q(s_0, a_0; w^k) \right)^2 \right\} \quad (21)$$

where $Q_{\text{target}} = \sum_{t=0}^{N-1} \gamma^t r(s_t, a_t) + \gamma^N Q(s_N, a_N; w^k)$ is the $N$-step target. The rollout length $N$ is identical to the horizon of chance constraint.

The parameterized actor is trained via gradient ascent in (17). In particular, we first compute $J$ and $\Phi$ through model rollout. Then $\nabla_\theta J$ and $\nabla_\theta p_s$ are computed via backpropagation though time with the dynamic model [27]. In practice, this process can be easily finished by any autograd package. The pseudo-code of the proposed algorithm is summarized in Algorithm 1.

---

**Algorithm 1** SPIL algorithm

---

Initialize $\pi_{\theta^0}$, $Q_{w^0}$, $I^0 = 0$, $s_0 \in \mathcal{S}$, $k = 1$
**repeat**
    Rollout $M$ trajectories by $N$ steps using policy $\pi^{k-1}$
    Estimate safe probability
        $p_s^{k-1} \leftarrow \frac{m}{M}$
    Update $\lambda$ via SPIL rules
        $\Delta^k \leftarrow 1 - \delta - p_s^{k-1}$
        $I^k \leftarrow (I^{k-1} + K_S \Delta^k)^+$
        $\lambda^k \leftarrow (K_P \Delta^k + K_I I^k)^+$
    Update critic:
        $\omega^k \leftarrow \omega^{k-1} + \alpha_\omega \nabla_\omega J_Q^{k-1}$
    Update actor:
        $\theta^k \leftarrow \theta^{k-1} + \frac{\alpha_\theta}{1 + \lambda^k} \left( \nabla_\theta J^{k-1} + \lambda^k \nabla_\theta \Phi(\pi_{\theta^{k-1}}, \tau) \right)$
    $k \leftarrow k + 1$
**until** $|Q^k - Q^{k-1}| \leq \zeta$ and $|\pi^k - \pi^{k-1}| \leq \zeta$

---

## IV. SIMULATION VALIDATION

### A. Example Description

In this section, the proposed SPIL is verified and compared in a car-following simulation in Fig. 6, where the ego car expects to drive fast and closely with the front car to reduce wind drag [28], while keeping a minimum distance between the two cars with a high probability. Concretely, we assume the ego car and front car follow a simple kinematics model, and the velocity of the front car is varying with uncertainty $\xi$.

The dynamics of the car-following example is given as:

$$\begin{aligned} s_{t+1} &= As_t + Ba_t + D\xi_t, \\ s &= [v_e, \quad v_f, \quad \epsilon]^\top, \\ A &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -T & T & 1 \end{bmatrix}, \\ B &= [T, 0, 0]^\top, \quad D = [0, T, 0]^\top \end{aligned} \quad (22)$$

where $v_e$ (m/s) is the velocity of ego car, $v_f$ (m/s) is velocity of front car, and $\epsilon$ (m) is the distance between the two cars. The action $a \in (-4, 3)\text{m/s}^2$ is the acceleration of ego car. The uncertainty $\xi_t \sim \mathcal{N}(0, 0.7)$ is truncated in the interval $(-7, 7)$. $T = 0.1s$ is the simulation time step. With a chance constraint on the minimum distance, the policy optimization problem is defined as:

$$\max_\pi \sum_{t=0}^{\infty} \gamma^t (0.2 v_{e,t} - 0.1\epsilon_t - 0.02 a_t^2)$$

$$\text{s.t.} \quad \Pr \left\{ \bigcap_{t=1}^{N} (\epsilon_t > 2) \right\} \geq 1 - \delta, \quad (23)$$

$$s_{t+1} = As_t + Ba_t + D\xi_t$$

(a) Cumulative reward under 90.0% threshold

(b) Cumulative reward under 99.9% threshold

(c) Safe probability under 90.0% threshold

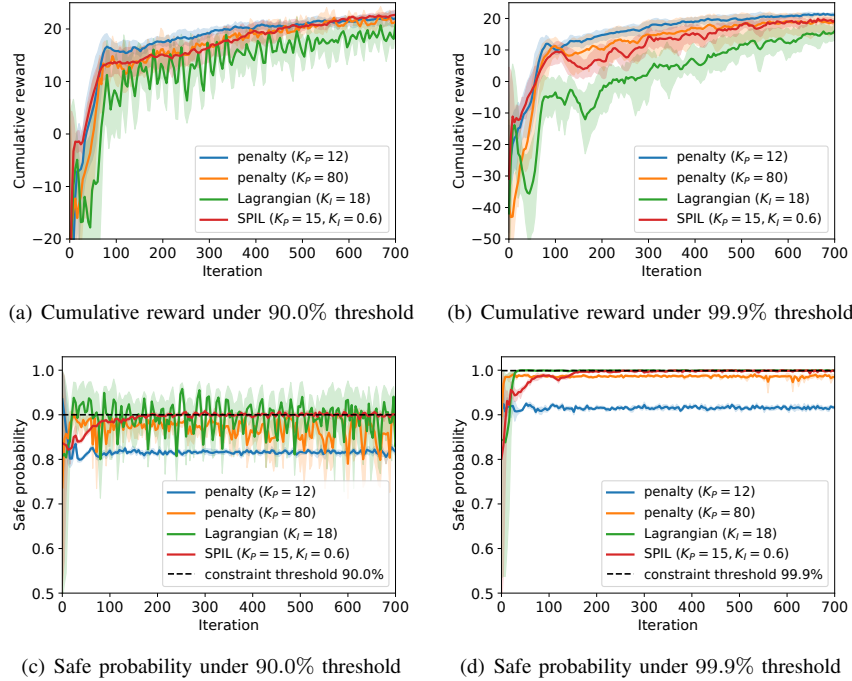(d) Safe probability under 99.9% threshold

Fig. 5. Comparison of performance among SPIL (separated Proportional-Integral Lagrangian), the penalty method and the Lagrangian method. The solid lines correspond to the mean and the shaded regions correspond to 95% confidence interval over 5 runs.
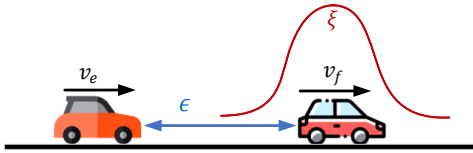


Fig. 6. Car-following scenario.

TABLE I
SPIL HYPER-PARAMETERS FOR SIMULATIONS

| Parameters | Symbol & Value |
|---|---|
| trajectories number | $M = 4096$ |
| constraint horizon | $N = 40$ |
| discount factor | $\gamma = 0.99$ |
| learning rate of policy network | $\alpha_\theta = 3e-4$ |
| learning rate of value network | $\alpha_\omega = 2e-4$ |
| parameters of $K_S$ | $(\beta, \varepsilon_1, \varepsilon_2) = (0.3, 0.2, 0.05)$ |
| parameters of $\phi(\cdot)$ | $(\tau, b_1, b_2) = (1e-3, 1, 0.45)$ |

where $v_{e,t}$ denotes the velocity of the ego car at step $t$.

### B. Algorithm Details

Three algorithms are employed to find the nearly optimal car-following policy, including SPIL (ours), the penalty method (amounts to proportional-only PIL), and the Lagrangian method (amounts to integral-only PIL). Note that all the algorithms are trained in the model-based manner. The coefficients of SPIL are selected as $K_P = 15$, $K_I = 0.6$ since it achieves the best results. The penalty method is sensitive to the penalty weight selection, so we adopt two weights $K_P = 12$ and 80. Actually, both of them cannot totally ensure the constraint satisfaction. For the Lagrangian method, we set $K_I = 18$ because it achieved the best performance in the pre-test compared to other values. The cumulative reward and safe probability in horizon $N$ are compared under two chance constraint thresholds: 90.0% and 99.9%, i.e., $\delta = 0.1$ and $\delta = 0.001$.

Both actor and critic are approximated by fully-connected neural networks. Each network has two hidden layers using rectified linear unit (ReLU) as activation functions, with 64 units per layer. The optimizer for the networks is Adam [29]. The main hyper-parameters are listed in Table I.

### C. Results

*1) Overall Performance:* The learning curves of all methods under two thresholds are illustrated in Fig. 5. We emphasize that any comparison should consider both safety and reward-winning. Generally, the proposed SPIL algorithm not only succeeds to satisfy the chance constraint without periodic oscillations, but also achieves the best cumulative reward among methods that meet the safety threshold.

For safety, the proposed SPIL satisfies the chance constraint in both thresholds as shown in Fig. 5(c) and Fig. 5(d). While the penalty methods with two weights both fail to meet the thresholds. Additionally, the large weight of $K_P = 80$ also leads to oscillation. The Lagrangian method basically satisfies the constraint, but suffers from periodic oscillations under 90.0% threshold. In a feedback control view, the SPIL combines the advantages of integral and proportion control, thus having a stable learning process with no steady-state errors.

In terms of reward-wining plotted in Fig. 5(a) and Fig. 5(b), SPIL achieves more reward than the Lagrangian method but less than the penalty method with $K_P = 12$. This is because

that the penalty method actually wins high performance at the cost of constraint violation.

*2) Ablation Study:* To demonstrate the necessity of integral separation proposed in Section III-B, we compare the results of our method with and without integral separation with five unsafe initial policies. As shown in Fig. 7, since the initial safe probability is low, the integral value $I$ increases rapidly at first. With a large $I$ and $\lambda$, the policy quickly becomes 100% safe. Unfortunately, since $\Delta = 0.999 - 1 = -0.001$ is too small in (10), the decline of $I$ and $\lambda$ is quite slow, leading to an overly conservative policy with poor reward. On the contrary, once the integral separation is equipped, the above problem is immediately solved. Note that the results in Fig. 5 and Fig. 7 are not comparable since the latter are conducted under manually chosen unsafe initial policies.



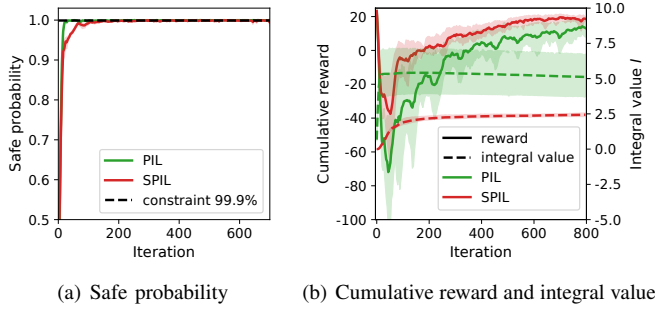(a) Safe probability     (b) Cumulative reward and integral value

Fig. 7. Comparison of performance of SPIL with and without integral separation technique (PIL) under 99.9% threshold. The solid and dotted lines correspond to the mean and the shaded regions correspond to 95% confidence interval over 5 runs. In (b), solid lines represent reward, dotted line represent integral value. The values for SPIL are in red and the values for PIL are in green.

The previous works typically adopt model-free methods to optimize the policy due to the discontinuity of the indicator function. In Section III-C, we propose a model-based gradient of safe probability to enable model-based optimization, which is believed to accelerate the training process. To verify this, we use a popular model-free algorithm PPO [30] to compute the gradient of the safe probability, following [16]. The learning curves of safe probability $p_s$ and cumulative reward $J$ are plotted in Fig. 8, where an iteration amounts to a batch of 4096 state-action pairs for both model-based and model-free versions. Our model-based SPIL reaches the required safe probability within about 250 iterations, at least five times faster than its model-free counterpart. This improvement comes from

the fact that the model-free method has first to learn an accurate cost value function before it optimizes the policy, while the model-based method makes use of the model to directly obtain a relatively precise gradient. It should be pointed out that this improvement is especially helpful for online training in the real world, where the policy should be safe as early as possible.
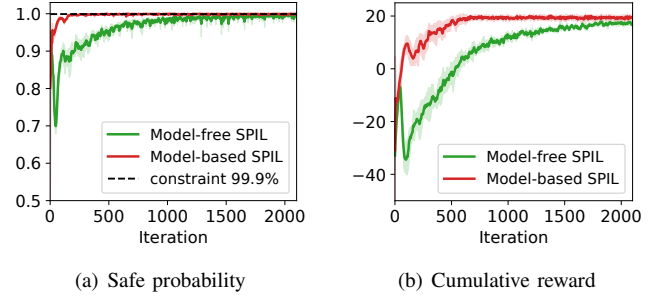


(a) Safe probability     (b) Cumulative reward

Fig. 8. Comparison of performance of SPIL with model-based and model-free optimization under 99.9% constraint threshold. The solid lines correspond to the mean and the shaded regions correspond to 95% confidence interval over 5 runs.

*3) Sensitivity Analysis:* To demonstrate the practicality of SPIL, we show that its high performance is relatively insensitive to hyper-parameter choice. We test the algorithm across different values of $K_P$, $K_I$ and $\beta$, while keeping all other parameters fixed. The results over 5 runs under 90.0% threshold are summarized in Table II, with the best parameters shown in bold. Even the worst case only leads to 1.8% degradation in safe probability and 6.5% degradation in cumulative reward.

## V. EXPERIMENTAL VALIDATION

### A. Experiment Description

To demonstrate the effectiveness of the proposed method for real-world safety-critical application, we apply it in a mobile robot navigation task. As illustrated in Fig. 9, the robot aims to follow the reference path (exactly the positive $x$-axis) without colliding with a moving obstacle. However, it does not know the behaviour or trajectory of the moving obstacle, which may be highly stochastic. Besides, the moving obstacle will not actively avoid the robot.

The robot locates its position and heading angle through a lidar, and it also has sensors to estimate its current

TABLE II
ALGORITHM PERFORMANCE WITH DIFFERENT PARAMETERS

| Value of $K_P$ ($K_I = 0.6$, $\beta = 0.3$) | 3.75 | 7.5 | **15** | 30 | 60 |
|---|---|---|---|---|---|
| Cumulative reward | $22.27 \pm 0.56$ | $21.86 \pm 0.67$ | $\mathbf{23.01 \pm 0.96}$ | $21.91 \pm 1.60$ | $21.49 \pm 1.61$ |
| Safe probability | $89.5 \pm 0.1\%$ | $90.2 \pm 0.3\%$ | $\mathbf{90.0 \pm 1.8\%}$ | $91.6 \pm 1.8\%$ | $88.2 \pm 6.6\%$ |
| Value of $K_I$ ($K_P = 15$, $\beta = 0.3$) | 0.15 | 0.3 | **0.6** | 1.2 | 2.4 |
| Cumulative reward | $21.60 \pm 0.66$ | $22.18 \pm 0.68$ | $\mathbf{23.01 \pm 0.96}$ | $21.77 \pm 0.43$ | $21.87 \pm 0.50$ |
| Safe probability | $89.9 \pm 1.0\%$ | $89.0 \pm 0.4\%$ | $\mathbf{90.0 \pm 1.8\%}$ | $90.5 \pm 0.2\%$ | $89.5 \pm 0.3\%$ |
| Value of $\beta$ ($K_P = 15$, $K_I = 0.6$) | 0.1 | 0.2 | **0.3** | 0.4 | 0.5 |
| Cumulative reward | $21.88 \pm 0.90$ | $21.25 \pm 0.44$ | $\mathbf{23.01 \pm 0.96}$ | $21.86 \pm 1.44$ | $21.74 \pm 0.68$ |
| Safe probability | $90.0 \pm 0.7\%$ | $90.8 \pm 1.0\%$ | $\mathbf{90.0 \pm 1.8\%}$ | $89.9 \pm 1.1\%$ | $90.1 \pm 0.4\%$ |

velocity and angular velocity. In this experiment, we let the obstacle share its current motion information through socket communications.
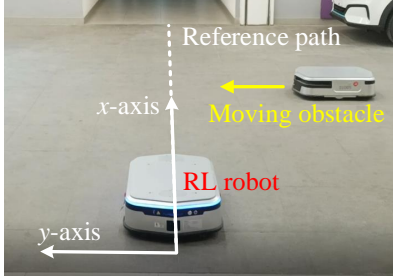


Fig. 9. Mobile robot navigation task.

### B. Experiment Details

*1) Dynamic Model:* We first present the model used for training. The wheeled robot adopts a two-wheel differential drive architecture, and takes the desired velocity $v^d$ and desired angular velocity $w^d$ as control commands. Nonetheless, its bottom-level control mechanism and response are unknown and varied in different conditions. Therefore, additional random terms $\xi^v$ and $\xi^\omega$ are introduced to describe this uncertainty. To cover the stochastic behaviour of obstacle, the uncertainty of the obstacle is also considered. The motions of both robot and moving obstacle are predicted through a simple kinematic model:

$$s = \begin{bmatrix} P^x \\ P^y \\ \alpha \\ v \\ \omega \end{bmatrix}, a = \begin{bmatrix} v^d \\ w^d \end{bmatrix},$$

$$s_{t+1} = \begin{bmatrix} P_t^x \\ P_t^y \\ \alpha_t \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ v_t^d \\ \omega_t^d \end{bmatrix} + T \begin{bmatrix} v_t \cos \alpha \\ v_t \sin \alpha \\ \omega_t \\ \xi_t^v \\ \xi_t^\omega \end{bmatrix} \quad (24)$$

where $T = 0.4s$ is the time step, $P^x$ and $P^y$ denote the position coordinates, $\alpha$ is the heading angle, $v$ is the velocity, $\omega$ is the angular velocity, $v^d$ is desired velocity, $w^d$ is desired angular velocity, $\xi^v$ is the uncertainty on the velocity, $\xi^\omega$ is the uncertainty on the angular velocity. For the robot, the control inputs $a$ are output by the policy network, which are also bounded by the following input constraints $|v - v^d| \leq 1.8T$, $|\omega - \omega^d| \leq 0.8T$. The velocity and angular velocity uncertainty of the robot is set to $\xi^v \sim \mathcal{N}(0, 0.08)$ and $\xi^\omega \sim \mathcal{N}(0, 0.05)$. For the model of the obstacle, the desired velocity and angular velocity are always set as the same as its current velocity and angular velocity, but with uncertainty $\xi^v \sim \mathcal{N}(0, 0.1)$ and $\xi^\omega \sim \mathcal{N}(0, 0.06)$ to indicate its stochastic behaviors. Although the true future behavior of obstacles is unknown to the ego robot, these uncertainty terms help to improve the robustness of the learned policy in real environments.

We admit the whole model is relatively naive and inaccurate, and the uncertainty term is given by experience instead of

estimation. However, we find it is enough to accomplish this task. A better choice is to update the model and uncertainty online through real-world experimental data. We leave it in our further work.

*2) Reward and Chance Constraints:* The robot aims to follow a reference path while avoiding collisions with a moving obstacle. The start point for the robot is near (1, 0). The reference path is the positive $x$-axis. Therefore, the reference $y$-position and heading angle are both 0. The reference velocity is set to 0.3 m/s. With additional regularization on the control command, the reward of the task is defined as:

$$r = -1.4(P^y)^2 - \alpha^2 - 16(v - 0.3)^2 - 0.2(v^d)^2 - 0.5(\omega^d)^2 \quad (25)$$

Then we impose the obstacle avoidance constraint. For simplicity, the robot and obstacle are both regarded as circles with a radius of 0.4m, and the distance between their centers is denoted as $\epsilon$. The chance constraint on the minimum distance is:

$$\Pr \left\{ \bigcap_{t=1}^{25} (\epsilon_t > 0.9) \right\} \geq 0.99 \quad (26)$$

*3) Algorithmic Parameters:* The network structure is exactly the same as that of the simulation. The main hyper-parameters are listed in Table III.

TABLE III
SPIL HYPER-PARAMETERS FOR EXPERIMENT

| Parameters | Symbol & Value |
|---|---|
| trajectories number | $M = 4096$ |
| constraint horizon | $N = 25$ |
| learning rate of policy network | $\alpha_\theta = 3e - 2$ |
| proportional coefficient | $K_P = 60$ |
| integral coefficient | $K_I = 0.02$ |
| parameters of $K_S$ | $(\beta, \varepsilon_1, \varepsilon_2) = (0.7, 0.2, 0.1)$ |
| parameters of $\phi(\cdot)$ | $(\tau, b_1, b_2) = (7e - 2, 1, 0.45)$ |

*4) Test Scenarios:* The learned policy will be tested in five scenarios with different obstacle behaviours. In the former three scenarios, the obstacle behaves normally, without a sudden stop or turn. To demonstrate the robustness and high intelligence of the trained robot, the latter two scenarios are under high randomness, where the obstacle drives in a complex trajectory or even deliberately blocks the robot. We stress that, in all experiments, the robot does not know the behavior or trajectory of the obstacle in advance, and all the experiments are conducted with the same network. This means that the method should have high generalization ability and intelligence to pass the five test scenarios.

### C. Results

We trained the policy network using the proposed SPIL algorithm with the dynamic model, and then tested the learned network on the real robot. A video for the results is available at https://youtu.be/oVDB2XqNoCU. We highly recommend readers to watch the video for an intuitive impression. Fig. 7 shows the control results of the learned policy in five scenarios. In Fig. 10(a), when the obstacle was moving at a low speed, the robot actively bypassed it from the left. If we increased
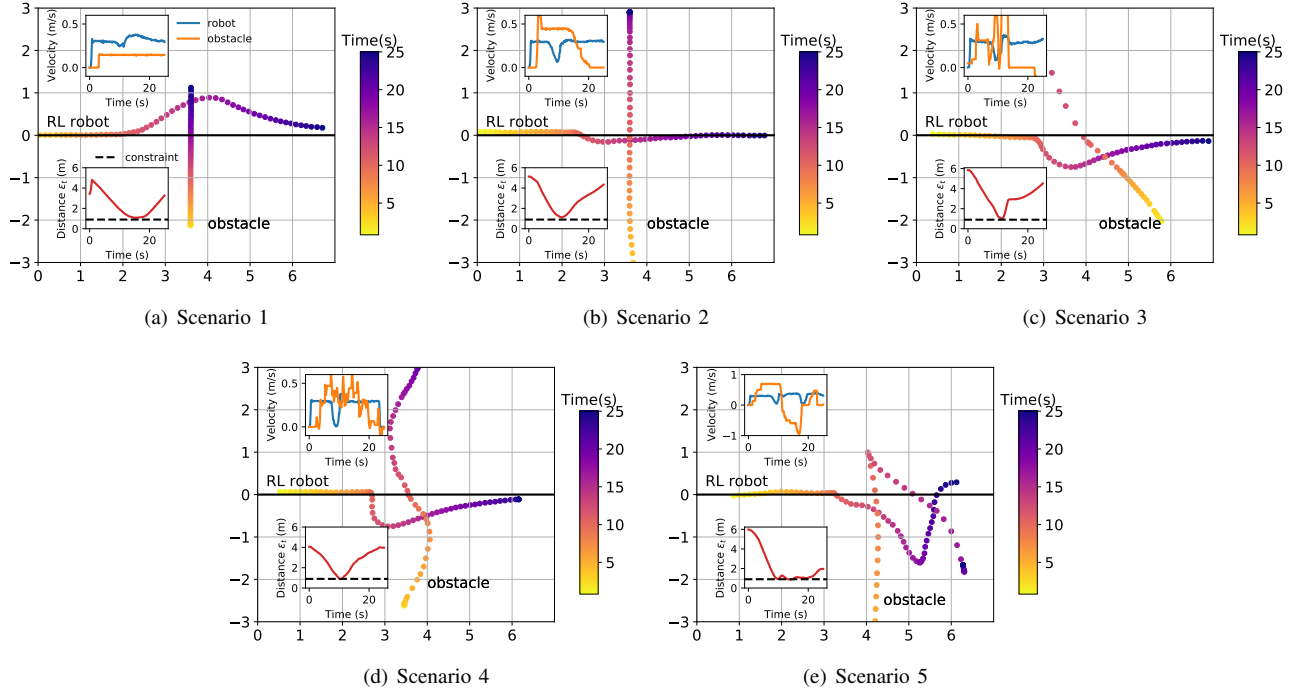
Fig. 10. Trajectories of robot and obstacle in five scenarios. All points are plotted with the same time interval, so the density of points also represents the velocity of the object. The speed of two objects and distance between two objects are also plotted. In most time, our robot keeps a safe distance from the obstacle as specified in the chance constraint.
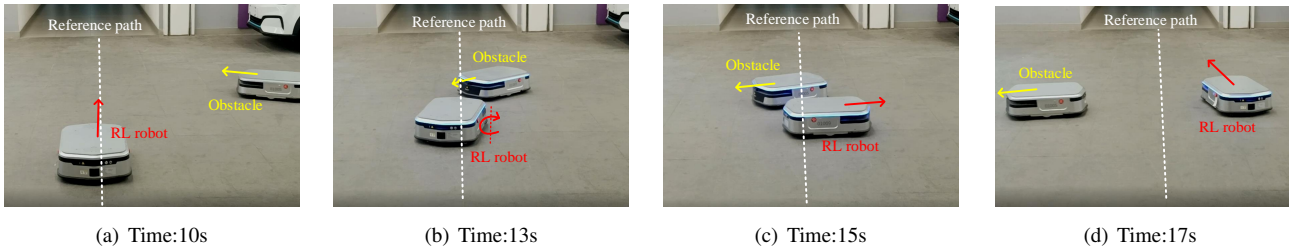


Fig. 11. Snapshots for the scenario 4. (a) The obstacle initially moved towards the positive direction of the reference path. (b) The obstacle suddenly changed its direction and moved towards the robot. Thus the robot urgently turned right to avoid it. (c) The robot passed around the obstacle from the right. (d) The robot returned to the reference path after the obstacle left.

the speed of the obstacle, the robot instead chose to stop and wait until the obstacle moved away, thereby avoiding a collision (See Fig. 10(b)). In Fig. 10(c), the ego robot avoided the moving obstacle from the right side while tracking the reference path. These results indicate that the learned policy can adopt different strategies according to the position and speed of the obstacle to achieve collision avoidance while tracking.

In scenarios 4 (Fig. 10(d)) and 5 (Fig. 10(e)), the behaviour of the obstacle robot was more aggressive. We controlled the obstacle to deliberately collide or block the robot to increase the difficulty of collision avoidance. See Fig. 10(d) and 11 for behavior details of the robot in scenario 4, and Fig. 10(e) and 12 for that in scenario 5. Results show that the learned policy can achieve safe movement even when the obstacle behaves aggressively, which demonstrates the amazing performance of our method.

## VI. CONCLUSION

We presented a model-based RL algorithm SPIL for chance-constrained policy optimization. Based on a feedback control view, we first reviewed and unified two existing chance-constrained RL methods to formulate a proportional-integral Lagrangian method, and enhanced it with an integral separation technique to prevent policy over-conservatism. To accelerate training, it also adopted a model-based gradient of safe probability for efficient policy optimization. We demonstrated the benefits of SPIL over previous methods in a car-following simulation. To prove its practicality, it was also applied to a real-world robot navigation task, where it successfully tracked the reference path while avoiding a highly stochastic moving obstacle. In the future, we will explore the possibility of online training, where the model and policy are both updated depending on the data from the real world to improve its online performance.

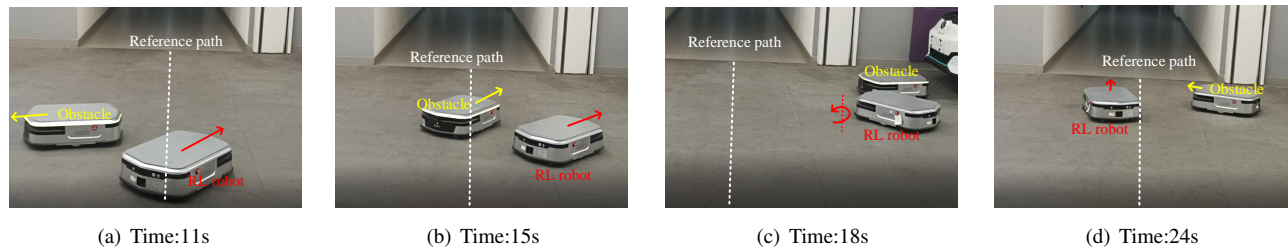(a) Time:11s      (b) Time:15s      (c) Time:18s      (d) Time:24s

Fig. 12. Snapshots for the scenario 5. (a) The robot turned right to avoid the coming obstacle. (b) The obstacle deliberately moved back and blocked the robot from returning to the reference path. Thus the robot had to keep moving in the wrong direction. (c) The robot found an opportunity to turn left when the obstacle was not blocking. (d) The robot successfully returned to the reference path.

## REFERENCES

[1] 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*, vol. 575, no. 7782, pp. 350–354, 2019.

[2] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, pp. 529–533, 2015.

[3] M. Hessel, J. Modayil, H. V. Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. G. Azar, and D. Silver, "Rainbow: Combining improvements in deep reinforcement learning," in *AAAI*, 2018.

[4] J. Duan, S. E. Li, Y. Guan, Q. Sun, and B. Cheng, "Hierarchical reinforcement learning for self-driving decision-making without reliance on labelled driving data," *IET Intelligent Transport Systems*, vol. 14, no. 5, pp. 297–305, 2020.

[5] Y. Ren, J. Duan, S. E. Li, Y. Guan, and Q. Sun, "Improving generalization of reinforcement learning with minimax distributional soft actor-critic," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1–6, IEEE, 2020.

[6] T. Kurutach, I. Clavera, Y. Duan, A. Tamar, and P. Abbeel, "Model-ensemble trust-region policy optimization," in *International Conference on Learning Representations*, 2018.

[7] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *International Conference on Machine Learning*, pp. 1861–1870, PMLR, 2018.

[8] G. Dulac-Arnold, D. Mankowitz, and T. Hester, "Challenges of real-world reinforcement learning," *arXiv preprint arXiv:1904.12901*, 2019.

[9] E. Altman, *Constrained Markov decision processes*, vol. 7. CRC Press, 1999.

[10] J. Achiam, D. Held, A. Tamar, and P. Abbeel, "Constrained policy optimization," in *International Conference on Machine Learning*, pp. 22–31, PMLR, 2017.

[11] T.-Y. Yang, J. Rosca, K. Narasimhan, and P. J. Ramadge, "Projection-based constrained policy optimization," in *International Conference on Learning Representations*, 2019.

[12] C. Tessler, D. J. Mankowitz, and S. Mannor, "Reward constrained policy optimization," in *International Conference on Learning Representations*, 2018.

[13] P. Petsagkourakis, I. O. Sandoval, E. Bradford, F. Galvanin, D. Zhang, and E. A. del Rio-Chanona, "Chance constrained policy optimization for process control and optimization," *arXiv preprint arXiv:2008.00030*, 2020.

[14] P. Geibel and F. Wysotzki, "Risk-sensitive reinforcement learning applied to control under constraints," *Journal of Artificial Intelligence Research*, vol. 24, pp. 81–108, 2005.

[15] A. Giuseppi and A. Pietrabissa, "Chance-constrained control with lexicographic deep reinforcement learning," *IEEE Control Systems Letters*, vol. 4, no. 3, pp. 755–760, 2020.

[16] S. Paternain, M. Calvo-Fullana, L. F. O. Chamon, and A. Ribeiro, "Learning safe policies via primal-dual methods," *2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 6491–6497, 2019.

[17] B. Peng, Y. Mu, Y. Guan, S. Li, Y. Yin, and J. Chen, "Model-based actor-critic with chance constraint for stochastic system," *ArXiv*, vol. abs/2012.10716, 2020.

[18] M. Farina, L. Giulioni, and R. Scattolini, "Stochastic linear model predictive control with chance constraints–a review," *Journal of Process Control*, vol. 44, pp. 53–67, 2016.

[19] Y. Chow, M. Ghavamzadeh, L. Janson, and M. Pavone, "Risk-constrained reinforcement learning with percentile risk criteria," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 6070–6120, 2017.

[20] B. W. Wah, T. Wang, Y. Shang, and Z. Wu, "Improving the performance of weighted lagrange-multiplier methods for nonlinear constrained optimization," *Information Sciences*, vol. 124, no. 1-4, pp. 241–272, 2000.

[21] A. Stooke, J. Achiam, and P. Abbeel, "Responsive safety in reinforcement learning by pid lagrangian methods," in *International Conference on Machine Learning*, pp. 9133–9143, PMLR, 2020.

[22] A. Mesbah, "Stochastic model predictive control: An overview and perspectives for future research," *IEEE Control Systems*, vol. 36, pp. 30–44, 2016.

[23] S. Boyd, S. P. Boyd, and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.

[24] L. Jia and X. Zhao, "An improved particle swarm optimization (PSO) optimized integral separation PID and its application on central position control system," *IEEE Sensors Journal*, vol. 19, pp. 7064–7071, 2019.

[25] A. Geletu, A. Hoffmann, and P. Li, "Analytic approximation and differentiability of joint chance constraints," *Optimization*, vol. 68, pp. 1985 – 2023, 2019.

[26] M. Deisenroth and C. E. Rasmussen, "Pilco: A model-based and data-efficient approach to policy search," in *Proceedings of the 28th International Conference on machine learning (ICML-11)*, pp. 465–472, Citeseer, 2011.

[27] S. E. Li, "Reinforcement Learning and Control." Tsinghua University: Lecture Notes. http://www.idlab-tsinghua.com/thulab/labweb/publications.html, 2020.

[28] F. Gao, S. E. Li, Y. Zheng, and D. Kum, "Robust control of heterogeneous vehicular platoon with uncertain dynamics and communication delay," *IET Intelligent Transport Systems*, vol. 10, no. 7, pp. 503–513, 2016.

[29] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[30] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *ArXiv*, vol. abs/1707.06347, 2017.