

Information-Theoretic Stochastic Optimal Control via Incremental Sampling-based Algorithms

Oktaý Arslan Evangelos Theodorou Panagiotis Tsiotras

Abstract. This paper considers optimal control of dynamical systems which are represented by nonlinear stochastic differential equations. It is well-known that the optimal control policy for this problem can be obtained as a function of a value function that satisfies a nonlinear partial differential equation, namely, the Hamilton-Jacobi-Bellman equation. This nonlinear PDE must be solved backwards in time, and this computation is intractable for large scale systems. Under certain assumptions, and after applying a logarithmic transformation, an alternative characterization of the optimal policy can be given in terms of a path integral. Path Integral (PI) based control methods have recently been shown to provide elegant solutions to a broad class of stochastic optimal control problems. One of the implementation challenges with this formalism is the computation of the expectation of a cost functional over the trajectories of the unforced dynamics. Computing such expectation over trajectories that are sampled uniformly may induce numerical instabilities due to the exponentiation of the cost. Therefore, sampling of low-cost trajectories is essential for the practical implementation of PI-based methods. In this paper, we use incremental sampling-based algorithms to sample useful trajectories from the unforced system dynamics, and make a novel connection between Rapidly-exploring Random Trees (RRTs) and information-theoretic stochastic optimal control. We show the results from the numerical implementation of the proposed approach to several examples.

Keywords: path integral, stochastic optimal control, sampling-based algorithms

1 Introduction

In [19,20], the authors showed the connection between Kullback-Leibler (KL) and Path Integral (PI) control with an information-theoretic view of stochastic optimal control. In addition, the authors derived the iterative path integral optimal control without relying on policy parameterizations, as in [17]. We review the work in [19,20] starting with the definitions of free energy and relative entropy and their connections to dynamic programming. In addition, we discuss how the iterative scheme developed in [19] and [20] can be modified to incorporate incremental sampling-based methods such as Rapidly-exploring Random Trees (RRTs) to guide sampling.

Within the mathematical framework of path integral control, the Feynman-Kac lemma plays an essential role, since it creates a connection between Stochastic Differential Equations (SDEs) and backward Partial Differential Equations (PDEs). This fundamental connection between SDEs and backward PDEs has inspired new avenues for the

The authors are with the Daniel Guggenheim School of Aerospace Engineering and the Institute for Robotics and Intelligent Machines, Georgia Institute of Technology, Atlanta, GA.

development of stochastic control algorithms such as Policy Improvement with Path Integrals (PI²) [18] that rely on forward sampling. PI² has been applied to a plethora of motor control tasks from robotic object manipulation and locomotion to general trajectory optimization and gain scheduling [2, 15, 16, 18], but it relies on a suitable parameterization of the optimal control policy. While policy parameterization such as Dynamic Movement Primitives (DMPs) [7] improves sampling by steering trajectories in high-dimensional state spaces towards areas of interest, it does not exploit the feedback structure provided by the path integral control framework. In PI² trajectories are sampled from the initial state of the task, the optimal parameter variations are computed, and the parameters are updated. In the next iteration, trajectories are sampled again from the same initial state and the iterative process continues until convergence. It is clear that in the case of policy parameterization one has to explicitly design the structure of the feedback control policy and then treat the gains as parameters to be optimized.

In this work, we use an alternative approach, which steers state trajectories towards relevant areas of the state space without the requirement of policy parameterization. In addition, the proposed approach improves sampling, while also allowing the use of path integral control in a feedback form.

2 Notation

A *probability space* is a triple (Ω, \mathcal{F}, p) where (Ω, \mathcal{F}) is a measurable space with Ω a non-empty set, which is called the *sample space*, $\mathcal{F} \subseteq 2^\Omega$ a σ -algebra of subsets of Ω , whose elements are called events, and p is a *probability measure* on \mathcal{F} , that is, p is a finite measure on \mathcal{F} with $p(\Omega) = 1$.

A real random variable is a function $X : \Omega \rightarrow \mathbb{R}$ with the property that $\{\omega \in \Omega : X(\omega) \leq x\} \in \mathcal{F}$ for each $x \in \mathbb{R}$. Such a function is said to be \mathcal{F} -measurable. An extended (real) random variable can also take the values $\pm\infty$. If X is a random variable on the probability space (Ω, \mathcal{F}, p) , then its *expectation* is defined by

$$\mathbb{E}_p[X] = \int_{\Omega} X(\omega) dp(\omega), \quad (1)$$

provided that the integral in the right-hand side exists. As usual, and for notational simplicity, in the sequel we will drop the explicit dependence on $\omega \in \Omega$ in (1). In other words, the notation $\mathbb{E}_p[X]$ is another (shorter) notation for the integral $\int X dp$.

3 Stochastic Control Based on Free Energy and Relative Entropy Dualities

Let (Ω, \mathcal{F}) be a measurable space where Ω is a non-empty set and $\mathcal{F} \subseteq 2^\Omega$ is a σ -algebra of subsets of Ω , and let $\mathbf{P}(\Omega)$ be the set of all probability measures defined on (Ω, \mathcal{F}) .

Definition 1 Let $\mathbf{p} \in \mathbf{P}(\Omega)$ be a probability measure, $\mathbf{x} = \mathbf{x}(\omega)$, $\omega \in \Omega$ be a random variable, $t, \rho \in \mathbb{R}$ be real numbers, and let $\mathcal{J}(\mathbf{x}, t)$ be a measurable function. The *Helmholtz free energy* of $\mathcal{J}(\mathbf{x}, t)$ with respect to \mathbf{p} is defined by

$$\mathcal{E}_{\mathbf{p}}(\mathcal{J}(\mathbf{x}, t); \rho) = \log \left(\int \exp(\rho \mathcal{J}(\mathbf{x}, t)) d\mathbf{p} \right) = \log \mathbb{E}_{\mathbf{p}}[\exp(\rho \mathcal{J}(\mathbf{x}, t))]. \quad (2)$$

Definition 2 Let $\mathbf{p}, \mathbf{q} \in \mathbf{P}(\Omega)$ be two probability measures. The *relative entropy* of \mathbf{p} with respect to \mathbf{q} is defined as¹:

$$\mathbb{KL}(\mathbf{q} \parallel \mathbf{p}) = \begin{cases} \int \log \left(\frac{d\mathbf{q}}{d\mathbf{p}} \right) d\mathbf{q} & \text{if } \mathbf{q} \ll \mathbf{p} \text{ and } \log \left(\frac{d\mathbf{q}}{d\mathbf{p}} \right) \in L^1, \\ +\infty & \text{otherwise.} \end{cases} \quad (3)$$

We will also consider the function $\xi(\mathbf{x}, t)$, defined by

$$\xi(\mathbf{x}, t) = \frac{1}{\rho} \mathcal{E}_{\mathbf{p}}(\mathcal{J}(\mathbf{x}, t); \rho) = \frac{1}{\rho} \log \mathbb{E}_{\mathbf{p}}[\exp(\rho \mathcal{J}(\mathbf{x}, t))]. \quad (4)$$

To derive the basic relationship between free energy and relative entropy [4], we express the expectation $\mathbb{E}_{\mathbf{p}}$ taken under the probability measure \mathbf{p} as a function of the expectation $\mathbb{E}_{\mathbf{q}}$ taken under the probability measure \mathbf{q} . More precisely, we have:

$$\mathbb{E}_{\mathbf{p}}[\exp(\rho \mathcal{J}(\mathbf{x}, t))] = \int \exp(\rho \mathcal{J}(\mathbf{x}, t)) \frac{d\mathbf{p}}{d\mathbf{q}} d\mathbf{q}.$$

By taking the logarithm of both sides of the previous equation and by making use of Jensen's inequality [4], it can be shown that:

$$\log \mathbb{E}_{\mathbf{p}}[\exp(\rho \mathcal{J}(\mathbf{x}, t))] \geq \int \rho \mathcal{J}(\mathbf{x}, t) d\mathbf{q} - \mathbb{KL}(\mathbf{q} \parallel \mathbf{p}). \quad (5)$$

Let $\rho < 0$. By multiplying both sides of (5) with $-1/|\rho|$, one obtains:

$$\boxed{\xi(\mathbf{x}, t) = -\frac{1}{|\rho|} \mathcal{E}_{\mathbf{p}}(\mathcal{J}(\mathbf{x}, t); \rho) \leq \mathbb{E}_{\mathbf{q}}[\mathcal{J}(\mathbf{x}, t)] + \frac{1}{|\rho|} \mathbb{KL}(\mathbf{q} \parallel \mathbf{p})} \quad (6)$$

where $\mathbb{E}_{\mathbf{q}}[\mathcal{J}(\mathbf{x}, t)] = \int \mathcal{J}(\mathbf{x}, t) d\mathbf{q}$. The inequality (6) provides us with a duality relationship between relative entropy and free energy. Essentially, one could define the following minimization problem:

$$-\frac{1}{|\rho|} \mathcal{E}_{\mathbf{p}}(\mathcal{J}(\mathbf{x}, t); \rho) = \inf_{\mathbf{q} \in \mathbf{P}(\Omega)} \left(\mathbb{E}_{\mathbf{q}}[\mathcal{J}(\mathbf{x}, t)] + \frac{1}{|\rho|} \mathbb{KL}(\mathbf{q} \parallel \mathbf{p}) \right). \quad (7)$$

It can be shown that the infimum in (7) is attained at \mathbf{q}^* , where

$$d\mathbf{q}^* = \frac{\exp(-|\rho| \mathcal{J}(\mathbf{x}, t))}{\int \exp(-|\rho| \mathcal{J}(\mathbf{x}, t)) d\mathbf{p}} d\mathbf{p}. \quad (8)$$

¹ Given two probability measures \mathbf{p} and \mathbf{q} , we say that \mathbf{q} is *absolute continuous* with \mathbf{p} and write $\mathbf{q} \ll \mathbf{p}$ if $\mathbf{q} = 0 \Rightarrow \mathbf{p} = 0$, see page 161 of [13].

A rather intuitive way of writing (6) is to express it in the following form:

$$\underbrace{-|\rho|^{-1}\mathcal{E}_p(\mathcal{J}(\mathbf{x}, t); \rho)}_{\text{Helmholtz Free Energy}} \leq \underbrace{\text{State Cost} + |\rho|^{-1}\text{Information Cost}}_{\text{Non-Equilibrium Free Energy}} \quad (9)$$

where “State Cost” and “Information Cost” are defined as $\mathbb{E}_q[\mathcal{J}(\mathbf{x}, t)]$ and $\mathbb{KL}(q\|\mathbf{p})$, respectively.

In the next sections, we derive the form of (7) for the case when \mathbf{x} is the state of a nonlinear stochastic differential equation affine in noise and control.

3.1 Application of the Legendre Transformation to Stochastic Differential Equations

We consider the general uncontrolled and controlled stochastic dynamics affine in noise as follows:

$$d\mathbf{x} = \mathbf{A}(\mathbf{x}) dt + \mathbf{C}(\mathbf{x}) d\mathbf{w}^{(0)}, \quad (10)$$

$$d\mathbf{x} = \mathbf{F}(\mathbf{x}, \mathbf{u}) dt + \mathbf{C}(\mathbf{x}) d\mathbf{w}^{(1)}, \quad (11)$$

where $\mathbf{x} \in \mathbb{R}^n$ denotes the state of the system, $\mathbf{u} \in \mathbb{R}^m$ denotes the control input, $\mathbf{C}(\mathbf{x}) \in \mathbb{R}^{n \times m}$ is the diffusion matrix, $\mathbf{F}(\mathbf{x}, \mathbf{u}) \in \mathbb{R}^n$ is the drift dynamics, and $\mathbf{w}^{(0),(1)} \in \mathbb{R}^m$ are Wiener processes (Brownian motion). The upper-scripts (0) and (1) are used to distinguish the two noise processes in the uncontrolled and controlled dynamics, respectively. The drift term $\mathbf{A}(\mathbf{x}) \in \mathbb{R}^n$ is defined by $\mathbf{A}(\mathbf{x}) = \mathbf{F}(\mathbf{x}, 0)$. The diffusion matrix may be partitioned as $\mathbf{C}(\mathbf{x}) = [\mathbf{0} \ \mathbf{C}_c^\top(\mathbf{x})]^\top$ where $\mathbf{0} \in \mathbb{R}^{(n-m) \times m}$ and $\mathbf{C}_c(\mathbf{x}) \in \mathbb{R}^{m \times m}$ is invertible. Similarly, the drift term in the controlled dynamics may be partitioned as $\mathbf{F}(\mathbf{x}, \mathbf{u}) = [\mathbf{F}_1^\top(\mathbf{x}, \mathbf{u}) \ \mathbf{F}_2^\top(\mathbf{x}, \mathbf{u})]^\top$ where $\mathbf{F}_1(\mathbf{x}, \mathbf{u}) \in \mathbb{R}^{m \times (n-m)}$ and $\mathbf{F}_2(\mathbf{x}, \mathbf{u}) \in \mathbb{R}^{m \times m}$; and the drift term in the uncontrolled dynamics may be partitioned as $\mathbf{A}(\mathbf{x}) = [\mathbf{A}_1^\top(\mathbf{x}) \ \mathbf{A}_2^\top(\mathbf{x})]^\top$ where $\mathbf{A}_1(\mathbf{x}) \in \mathbb{R}^{m \times (n-m)}$ and $\mathbf{A}_2(\mathbf{x}) \in \mathbb{R}^{m \times m}$. The class of systems whose matrices can be partitioned as such contains rigid body, and multi body dynamics as well as kinematic models such as the ones considered in this work. Henceforth, for simplicity, we will assume that $m = n$. The case when $m < n$ can be treated similarly; see for instance [22]. Let $\Sigma(\mathbf{x}) = \mathbf{C}(\mathbf{x})\mathbf{C}^\top(\mathbf{x}) \in \mathbb{R}^{m \times m}$ and also define the following quantity:

$$\delta\mathbf{F}(\mathbf{x}, \mathbf{u}) = \mathbf{F}(\mathbf{x}, \mathbf{u}) - \mathbf{A}(\mathbf{x}) = \mathbf{F}(\mathbf{x}, \mathbf{u}) - \mathbf{F}(\mathbf{x}, 0), \quad \forall \mathbf{x}, \mathbf{u}.$$

To the system (11) we also associated the state cost

$$\mathcal{J}(\mathbf{x}(\cdot), t) = \Phi(\mathbf{x}(t_f)) + \int_t^{t_f} q(\mathbf{x}(\tau), \tau) d\tau. \quad (12)$$

With a slight abuse of notation we will also use $\mathcal{J}(\mathbf{x}, t)$ to denote the value of $\mathcal{J}(\mathbf{x}(\cdot), t)$ along the trajectory $\mathbf{x}(\cdot)$ starting from $\mathbf{x} = \mathbf{x}(t)$ at time t . Expectations evaluated on trajectories generated by the uncontrolled dynamics and controlled dynamics will be represented by $\mathbb{E}_p[\cdot]$ and $\mathbb{E}_q[\cdot]$, respectively. The following fact can be found in [22].

Proposition 1. Given the measures \mathbf{p}, \mathbf{q} induced by the trajectories of (10) and (11), respectively, the *Radon-Nikodym derivative* of \mathbf{q} with respect to \mathbf{p} is defined by

$$\frac{d\mathbf{q}}{d\mathbf{p}} = \exp \left(\int_t^{t_f} \delta \mathbf{F}^\top(\mathbf{x}(\tau), \mathbf{u}(\tau)) \mathbf{C}^{-1}(\mathbf{x}(\tau)) d\mathbf{w}^{(1)}(\tau) \right) + \exp \left(\int_t^{t_f} \frac{1}{2} \delta \mathbf{F}^\top(\mathbf{x}(\tau), \mathbf{u}(\tau)) \boldsymbol{\Sigma}^{-1}(\mathbf{x}(\tau)) \delta \mathbf{F}(\mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau \right). \quad (13)$$

Given equation (13), the relative entropy term in (6) takes the form:

$$\begin{aligned} \frac{1}{|\rho|} \mathbb{KL}(\mathbf{q} \parallel \mathbf{p}) &= \mathbb{E}_{\mathbf{q}} \left[\frac{1}{|\rho|} \int_t^{t_f} \delta \mathbf{F}^\top(\mathbf{x}(\tau), \mathbf{u}(\tau)) \mathbf{C}^{-1}(\mathbf{x}(\tau)) d\mathbf{w}^{(1)}(\tau) \right] + \\ &\quad \mathbb{E}_{\mathbf{q}} \left[\frac{1}{|\rho|} \int_t^{t_f} \frac{1}{2} \delta \mathbf{F}^\top(\mathbf{x}(\tau), \mathbf{u}(\tau)) \boldsymbol{\Sigma}^{-1}(\mathbf{x}(\tau)) \delta \mathbf{F}(\mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau \right] \\ &= \mathbb{E}_{\mathbf{q}} \left[\frac{1}{2|\rho|} \int_t^{t_f} \delta \mathbf{F}^\top(\mathbf{x}(\tau), \mathbf{u}(\tau)) \boldsymbol{\Sigma}^{-1}(\mathbf{x}(\tau)) \delta \mathbf{F}(\mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau \right], \end{aligned}$$

where the first term in the previous expression vanishes since the expectations term $\mathbb{E}_{\mathbf{q}} [\delta \mathbf{F}^\top(\mathbf{x}(\tau), \mathbf{u}(\tau)) \mathbf{C}^{-1}(\mathbf{x}(\tau)) d\mathbf{w}^{(1)}(\tau)]$ becomes

$$\mathbb{E}_{\mathbf{q}} [\delta \mathbf{F}^\top(\mathbf{x}(\tau), \mathbf{u}(\tau)) \mathbf{C}^{-1}(\mathbf{x}(\tau))] \mathbb{E}_{\mathbf{q}} [d\mathbf{w}^{(1)}(\tau)] = 0, \quad \forall \tau, t \leq \tau \leq t_f \quad (14)$$

Substituting the previous expression of the Kullback-Leibler divergence into (6) one obtains

$$\begin{aligned} -\frac{1}{|\rho|} \mathcal{E}_{\mathbf{p}}(\mathcal{J}(\mathbf{x}, t); \rho) &\leq \mathbb{E}_{\mathbf{q}} [\mathcal{J}(\mathbf{x}, t)] + \\ &\quad \mathbb{E}_{\mathbf{q}} \left[\frac{1}{2|\rho|} \int_t^{t_f} \delta \mathbf{F}^\top(\mathbf{x}(\tau), \mathbf{u}(\tau)) \boldsymbol{\Sigma}^{-1}(\mathbf{x}(\tau)) \delta \mathbf{F}(\mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau \right]. \end{aligned}$$

The previous equation can be written in the form (9) with state cost term defined as

$$\mathbb{E}_{\mathbf{q}} [\mathcal{J}(\mathbf{x}, t)], \quad (15)$$

and information cost defined as

$$\mathbb{E}_{\mathbf{q}} \left[\frac{1}{2|\rho|} \int_t^{t_f} \delta \mathbf{F}^\top(\mathbf{x}(\tau), \mathbf{u}(\tau)) \boldsymbol{\Sigma}^{-1}(\mathbf{x}(\tau)) \delta \mathbf{F}(\mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau \right]. \quad (16)$$

Next, we further specialize the class of systems where (9) is applied to, and discuss its connections to stochastic optimal control as in [4, 19, 20]. To this end, let us consider the special case of (10) and (11) with uncontrolled and controlled stochastic dynamics of the following form, respectively:

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}) dt + \frac{1}{\sqrt{|\rho|}} \mathbf{B}(\mathbf{x}) d\mathbf{w}^{(0)}, \quad (17)$$

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}) dt + \mathbf{B}(\mathbf{x}) \left(\mathbf{u} dt + \frac{1}{\sqrt{|\rho|}} d\mathbf{w}^{(1)} \right), \quad (18)$$

where $\mathbf{x} \in \mathbb{R}^n$ denotes the state of the system, $\mathbf{B}(\mathbf{x}) \in \mathbb{R}^{n \times m}$ is the control/diffusion matrix, $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^n$ is the passive dynamics, $\mathbf{u} \in \mathbb{R}^m$ is the control vector and $\mathbf{w}^{(0),(1)}$ are m -dimensional Wiener noise processes.

For the dynamics in (17) and (18) the form of the Radon-Nikodym derivative in (13) can be computed as follows. Noticing that $\delta \mathbf{F}(\mathbf{x}, \mathbf{u}) = \mathbf{B}(\mathbf{x})\mathbf{u}$, $\mathbf{C}(\mathbf{x}) = \mathbf{B}(\mathbf{x})/\sqrt{|\rho|}$ and $\Sigma(\mathbf{x}) = \mathbf{B}(\mathbf{x})\mathbf{B}^\top(\mathbf{x})/|\rho|$, and substituting these expressions in (13) yields

$$\frac{dq}{dp} = \exp(|\rho|\eta(\mathbf{u}, t)) \quad \text{and} \quad \frac{dp}{dq} = \exp(-|\rho|\eta(\mathbf{u}, t)), \quad (19)$$

where $\eta(\mathbf{u}, t)$ is given by:

$$\eta(\mathbf{u}, t) = \frac{1}{2} \int_t^{t_f} \mathbf{u}^\top(\tau) \mathbf{u}(\tau) d\tau + \frac{1}{\sqrt{|\rho|}} \int_t^{t_f} \mathbf{u}^\top(\tau) d\mathbf{w}^{(1)}. \quad (20)$$

Substitution of (19) and (20) into inequality (6) yields the following result:

$$-\frac{1}{|\rho|} \log \mathbb{E}_p [\exp(-|\rho|\mathcal{J}(\mathbf{x}, t))] \leq \mathbb{E}_q \left[\mathcal{J}(\mathbf{x}, t) + \frac{1}{|\rho|} \eta(\mathbf{u}, t) \right]. \quad (21)$$

The expectation on the right side of the inequality in (21) is further simplified as follows:

$$\underbrace{-\frac{1}{|\rho|} \log \mathbb{E}_p [\exp(-|\rho|\mathcal{J}(\mathbf{x}, t))]}_{\xi(\mathbf{x}, t)} \leq \underbrace{\mathbb{E}_q \left[\mathcal{J}(\mathbf{x}, t) + \frac{1}{2} \int_t^{t_f} \mathbf{u}(\tau)^\top \mathbf{u}(\tau) d\tau \right]}_{\text{Total Cost}}. \quad (22)$$

The right-hand side term in the above inequality corresponds to the cost function of a stochastic optimal control problem that is bounded from below by the free energy. Surprisingly, inequality (22) was derived without relying on any principle of optimality. Inequality (22) essentially defines a minimization process in which the right-hand side part of the inequality is minimized with respect to $\eta(\mathbf{u}, t)$ and therefore with respect to the corresponding control \mathbf{u} . At the minimum, when $\mathbf{u} = \mathbf{u}^*$, the right-hand side of inequality in (22) attains its optimal value $\xi(\mathbf{x}, t)$. Under the optimal control \mathbf{u}^* , and according to (8), the corresponding optimal distribution takes the form

$$dq^* = \frac{\exp(-|\rho|\Phi(\mathbf{x}(t_f))) \exp\left(-|\rho| \int_t^{t_f} q(\mathbf{x}(\tau), \tau) d\tau\right)}{\int \exp(-|\rho|\Phi(\mathbf{x}(t_f))) \exp\left(-|\rho| \int_t^{t_f} q(\mathbf{x}(\tau), \tau) d\tau\right) dp} dp. \quad (23)$$

The work [19, 20] inspired by early mathematical developments in control theory [4, 5], has shown that the value function $\xi(\mathbf{x}, t)$ in (22) satisfies the Hamilton-Jacobi-Bellman equation and it has made the connection with more recent work in machine learning [8, 21] on Kullback-Leibler and path integral control.

3.2 Connection with Dynamic Programming (DP)

An important question that arises is: What is the link between (22) and the principle of optimality in dynamic programming? To address this question, we show that $\xi(\mathbf{x}, t)$ satisfies the Hamilton-Jacobi-Bellman (HJB) equation associated with the optimal control problem (18)-(12) and hence, $\xi(\mathbf{x}, t)$ is the corresponding value function of the following minimization problem

$$\begin{aligned}\xi(\mathbf{x}, t) &= \min_{\substack{\mathbf{u}(\tau) \\ t \leq \tau \leq t_f}} \mathbb{E}_q \left[\Phi(\mathbf{x}(t_f)) + \int_t^{t_f} (q(\mathbf{x}(\tau), \tau) + \frac{1}{2} \mathbf{u}^\top(\tau) \mathbf{u}(\tau)) d\tau \right] \\ &= \min_{\substack{\mathbf{u}(\tau) \\ t \leq \tau \leq t_f}} \mathbb{E}_q \left[\mathcal{J}(\mathbf{x}, \tau) + \frac{1}{2} \int_t^{t_f} \mathbf{u}^\top(\tau) \mathbf{u}(\tau) d\tau \right],\end{aligned}\quad (24)$$

where the expectation is computed over the trajectories of (18). To see this, we introduce $\Psi(\mathbf{x}, t) \triangleq \mathbb{E}_p [\exp(\rho \mathcal{J}(\mathbf{x}, t))]$ and apply the Feynman-Kac lemma [6] to arrive at the backward Chapman-Kolmogorov partial differential equation (PDE)

$$\begin{aligned}-\partial_t \Psi(\mathbf{x}, t) &= -|\rho| q(\mathbf{x}, t) \Psi(\mathbf{x}, t) + \mathbf{f}^\top(\mathbf{x}) \nabla \Psi_{\mathbf{x}}(\mathbf{x}, t) \\ &\quad + \frac{1}{2|\rho|} \text{tr}(\nabla \Psi_{\mathbf{xx}}(\mathbf{x}, t) \mathbf{B}(\mathbf{x}) \mathbf{B}(\mathbf{x})^\top)\end{aligned}\quad (25)$$

with boundary condition $\Psi(\mathbf{x}(t_f), t_f) = \exp(-|\rho| \Phi(\mathbf{x}(t_f)))$, which governs the evolution of $\Psi(\mathbf{x}, t)$ along the trajectories of (18) subject to $\mathbf{x} = \mathbf{x}(t)$. Since $\xi(\mathbf{x}, t) = -\log \Psi(\mathbf{x}, t)/|\rho|$, it follows that

$$\begin{aligned}\partial_t \Psi(\mathbf{x}, t) &= -|\rho| \Psi(\mathbf{x}, t) \partial_t \xi(\mathbf{x}, t), \\ \nabla \Psi_{\mathbf{x}}(\mathbf{x}, t) &= -|\rho| \Psi(\mathbf{x}, t) \nabla \xi_{\mathbf{x}}(\mathbf{x}, t), \\ \nabla \Psi_{\mathbf{xx}}(\mathbf{x}, t) &= |\rho| \Psi(\mathbf{x}, t) \nabla \xi_{\mathbf{xx}}(\mathbf{x}, t) - |\rho|^2 \Psi(\mathbf{x}, t) \nabla \xi_{\mathbf{x}}(\mathbf{x}, t) \nabla \xi_{\mathbf{x}}^\top(\mathbf{x}, t).\end{aligned}$$

In this case, it can be shown that $\xi(\mathbf{x}, t)$ satisfies the nonlinear PDE

$$\begin{aligned}-\partial_t \xi(\mathbf{x}, t) &= q(\mathbf{x}, t) + \nabla \xi_{\mathbf{x}}^\top(\mathbf{x}, t) \mathbf{f}(\mathbf{x}) - \frac{1}{2} \nabla \xi_{\mathbf{x}}^\top(\mathbf{x}, t) \mathbf{B}(\mathbf{x}) \mathbf{B}^\top(\mathbf{x}) \nabla \xi_{\mathbf{x}}(\mathbf{x}, t) \\ &\quad + \frac{1}{2|\rho|} \text{tr}(\nabla \xi_{\mathbf{xx}}(\mathbf{x}, t) \mathbf{B}(\mathbf{x}) \mathbf{B}^\top(\mathbf{x})),\end{aligned}\quad (26)$$

subject to the boundary condition $\xi(\mathbf{x}(t_f), t_f) = \Phi(\mathbf{x}(t_f))$. The nonlinear PDE (26) corresponds to the HJB equation associated with the optimal control problem (24) and hence $\xi(\mathbf{x}, t)$ is the corresponding minimizing value function [14]. It is important to note, however, that the principle of optimality was not used to derive (26).

3.3 Path Integral Control with Initial Sampling Policies

According to (22), in order to find the value function $\xi(\mathbf{x}, t)$, sampling of trajectories under the uncontrolled dynamics is performed, and the left-hand side of (22) is evaluated on these trajectories. However, in high-dimensional spaces, it is desirable to steer

sampling towards specific areas of the state space. To do so, we have to incorporate an initial control policy into the uncontrolled dynamics. Therefore, instead of sampling from the uncontrolled dynamics (17), we sample, instead, based on the stochastic dynamics:

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}) dt + \mathbf{B}(\mathbf{x}) \left(\mathbf{u}_{\text{in}} dt + \frac{1}{\sqrt{|\rho|}} d\mathbf{w}^{(1)} \right), \quad (27)$$

where \mathbf{u}_{in} is an initial control policy. In [19, 20], the authors derived an iterative PI control without relying on previous policy parameterizations. More precisely, when sampling from the dynamics (27) the work in [20] and [19] showed that the value function $\xi(\mathbf{x}, t)$ is expressed as

$$\xi(\mathbf{x}, t) = -\frac{1}{|\rho|} \log \left(\int \exp(-|\rho| S(\mathbf{x}, \mathbf{u}_{\text{in}}(\mathbf{x}, t), t)) d\mathbf{q}_{\text{in}} \right)$$

where the term $S(\mathbf{x}, \mathbf{u}_{\text{in}})$ is defined as

$$\begin{aligned} S(\mathbf{x}, \mathbf{u}_{\text{in}}) = & \underbrace{\Phi(\mathbf{x}(t_f)) + \int_t^{t_f} q(\mathbf{x}(\tau), \tau) d\tau}_{\mathcal{J}(\mathbf{x}, t)} + \\ & \underbrace{\frac{1}{2} \int_t^{t_f} \mathbf{u}_{\text{in}}^T(\tau) \mathbf{u}_{\text{in}}(\tau) d\tau + \frac{1}{\sqrt{|\rho|}} \int_t^{t_f} \mathbf{u}_{\text{in}}^T(\tau) d\mathbf{w}^{(1)}(\tau)}_{\eta(\mathbf{u}_{\text{in}}, t)}, \end{aligned} \quad (28)$$

where the term $\eta(\mathbf{u}_{\text{in}}, t)$ appears due to sampling based on the dynamics (27), while the term $\mathcal{J}(\mathbf{x}, t)$ is the state-dependent part of the total cost function in (22). The path integral control is now expressed as [19]

$$\mathbf{u}_{\text{PI}}(\mathbf{x}, t) dt = \mathbf{u}_{\text{in}}(\mathbf{x}, t) dt + \delta \mathbf{u}(\mathbf{x}, t), \quad (29)$$

where the term $\delta \mathbf{u}(\mathbf{x}, t)$ is defined by

$$\delta \mathbf{u}(\mathbf{x}, t) = \frac{1}{\sqrt{|\rho|}} \mathbb{E}_{\mathbf{q}^*} [d\mathbf{w}^{(1)}] = \frac{1}{\sqrt{|\rho|}} \int d\mathbf{w}^{(1)} d\mathbf{q}^*, \quad (30)$$

and where the expectation is taken under the optimal probability

$$d\mathbf{q}^* = \frac{\exp(-|\rho| S(\mathbf{x}, \mathbf{u}_{\text{in}}))}{\int \exp(-|\rho| S(\mathbf{x}, \mathbf{u}_{\text{in}})) d\mathbf{q}_{\text{in}}} d\mathbf{q}_{\text{in}}. \quad (31)$$

During implementation, equation (32) is approximated as

$$\delta \mathbf{u}(\mathbf{x}, t) = \frac{1}{\sqrt{|\rho|}} \sum_{k=1}^{\#\text{traj}} p_k d\mathbf{w}^{(1)}(\omega_k) \quad \text{with} \quad p_k = \frac{\exp(-|\rho| S(\mathbf{x}_k, \mathbf{u}_{\text{in}}))}{\sum_{\ell=1}^{\#\text{traj}} \exp(-|\rho| S(\mathbf{x}_\ell, \mathbf{u}_{\text{in}}))} \quad (32)$$

The initial policy \mathbf{u}_{in} can be a suboptimal control law, a hand-tuned PD, PID control, or feedforward control. In this paper, we consider a feedforward control given by the RRT algorithm as the initial control policy. In this case, the RRT-based optimal path integral control takes the form

$$\mathbf{u}_{\text{PI}}(\mathbf{x}, t) dt = \mathbf{u}_{\text{RRT}}(t) dt + \delta \mathbf{u}(\mathbf{x}, t). \quad (33)$$

In the next section, we discuss how to use the RRT algorithm to compute the initial control policy \mathbf{u}_{RRT} .

4 Trajectory Sampling via Sampling-based Algorithms

As shown in the previous sections, sampling of useful trajectories from the unforced dynamics can be a tedious task. This issue can be addressed by first computing a “good enough” initial trajectory and then sampling local trajectories in the neighborhood of this trajectory. In the proposed approach, we use a probabilistic algorithm to compute an initial trajectory quickly. Probabilistic methods have proven to be very efficient for the solution of motion planning problems with dynamic constraints in high dimensional search spaces. Among them, Rapidly-exploring Random Trees (RRTs) [3, 11, 12] are among the most popular for solving single query motion planning problems. The main body of the RRT algorithm is given in Algorithm ??.

In the proposed approach, we leverage the speed and exploration capabilities of the RRT algorithm to compute an initial policy quickly by making a minor modification of the RRT primitive procedures. Since both final time and final state are given, the search space is formed by adding an additional time dimension T to the state space \mathcal{X} . Our search space, goal set and free space are thus defined as $\mathcal{Z} = \mathcal{X} \times T$, $\mathcal{Z}_{\text{goal}} = \mathcal{X}_{\text{goal}} \times T_{\text{goal}}$, and $\mathcal{Z}_{\text{free}} = \mathcal{Z} \setminus \mathcal{Z}_{\text{goal}}$, respectively. The RRT algorithm is then run to find a trajectory starting from an initial point $z_{\text{init}} = (x_{\text{init}}, t_{\text{init}})$ to the goal set $\mathcal{Z}_{\text{goal}}$ while avoiding the obstacles in \mathcal{X} . The primitive procedures used by the RRT algorithm are given below:

Sampling: **Sample** : $\mathbb{N} \rightarrow \mathcal{Z}_{\text{free}}$ returns independent, identically distributed (i.i.d) samples from $\mathcal{Z}_{\text{free}}$.

Nearest neighbor: **Nearest** returns a point from a given finite set V , which is the point closest to a given point \mathbf{z} in terms of a given distance function.

Steering: Given two points \mathbf{z}_1 and \mathbf{z}_2 in $\mathcal{Z}_{\text{free}}$, **Steer** extends \mathbf{z}_1 towards \mathbf{z}_2 by sampling trajectories from the unforced dynamics of the system. Specifically, the procedure samples a set of trajectories emanating from \mathbf{z}_1 and returns the closest end point of this set of trajectories with respect to a given distance function.

Collision checking: Given a trajectory σ , the Boolean function **ObstacleFree**(σ) checks whether σ belongs to $\mathcal{Z}_{\text{free}}$ or not. It returns **True** if the trajectory is a subset of $\mathcal{Z}_{\text{free}}$, i.e., $\sigma \subset \mathcal{Z}_{\text{free}}$, and **False** otherwise.

Graph extension: **Extend** is a function that extends the nearest vertex of the graph \mathcal{G} toward the randomly sampled point z_{rand} . Since time always flows in forward direction, we make sure that **Extend** computes valid connections, i.e., it returns false if the time value of z_{rand} is less than that of the nearest vertex in the graph. The **Extend** procedure of the RRT algorithm is shown in Algorithm ??.

The body of the path-integral based RRT algorithm is shown in Algorithm ?? . It runs in a receding horizon fashion, that is, it computes a “good enough” control input and executes the first portion of the control signal at each time step. The algorithm starts by initializing the current time and state with the initial values in Lines 2-3. The algorithm then computes an initial policy in Line 5 by using the RRT algorithm. The steering procedure in the RRT algorithm is slightly modified in order to sample dynamically feasible trajectories. The steering procedure first samples a fixed number of trajectories from the unforced dynamics and then chooses the one that has the closest terminal state towards the desired point. Once a trajectory that reaches the goal set has been computed, the corresponding trajectory σ_{RRT} , along with control the signal \mathbf{u}_{RRT} , are extracted from the computed data structure in Line 6. Then, the algorithm proceeds by locally sampling trajectories around $(\sigma_{\text{RRT}}, \mathbf{u}_{\text{RRT}})$ and computes the variation in the control $\delta \mathbf{u}(\mathbf{x}, t)$ according to (30) by using information of local trajectories. Since we have M number of local trajectories, the expectation in (30) is numerically approximated by using the expression in (32). For each local trajectory σ_k , a cost value is computed as $S(\sigma_k, \mathbf{u}_{\text{RRT}})$ and its desirability value is computed by exponentiating the corresponding cost value, i.e., $d_k = \exp(-|\rho|S(\sigma_k, \mathbf{u}_{\text{RRT}}))$. Then, the variation term in control $\delta \mathbf{u}(\mathbf{x}, t)$ is computed by taking the weighted average of all noise profiles which create the local trajectories and the weight of each trajectory is computed as the normalized desirability value, i.e., $p_k = d_k / \sum_{\ell=1}^N d_{\ell}$. The iteration of the algorithm is completed by executing the first τ times of the computed control signal and the algorithm keeps repeating the same steps until the final time is reached.

5 Numerical Simulations

In this section, we present a series of simulated experiments using a kinematic car model. We are interested in controlling a vehicle, whose motion is described by the following kinematic equations:

$$\dot{x} = v \cos \theta, \quad \dot{y} = v \sin \theta, \quad \dot{\theta} = w/r \quad (34)$$

where x, y are the Cartesian coordinates of a reference point of the vehicle, v is its speed, w is the control input and r is a positive constant. We assume that the admissible control inputs, are restricted by $w \in [-1, 1]$. We would like to find an optimal policy for the heading rate w to move the vehicle from a given initial configuration $(x_i, y_i, \theta_i)^\top$ to a final configuration $(x_f, y_f, \theta_f)^\top$ within some fixed final time t_f .

Let $\mathbf{x}_1 = x$, $\mathbf{x}_2 = y$, $\mathbf{x}_3 = \theta$ be the states and $\mathbf{u} = w$ be the control input of the system. Then (34) can be rewritten as

$$\dot{\mathbf{x}}_1 = v \cos \mathbf{x}_3, \quad \dot{\mathbf{x}}_2 = v \sin \mathbf{x}_3, \quad \dot{\mathbf{x}}_3 = \mathbf{u}/r. \quad (35)$$

Assuming the system is subjected to noise of intensity α in the control channel, (35) can be written in the standard form

$$\begin{pmatrix} d\mathbf{x}_1 \\ d\mathbf{x}_2 \\ d\mathbf{x}_3 \end{pmatrix} = \begin{pmatrix} v \cos \mathbf{x}_3 \\ v \sin \mathbf{x}_3 \\ 0 \end{pmatrix} dt + \begin{pmatrix} 0 \\ 0 \\ 1/r \end{pmatrix} (\mathbf{u} dt + \alpha d\mathbf{w}), \quad (36)$$

where \mathbf{f} , \mathbf{B} and ρ in (27) are defined as follows

$$\mathbf{f}(\mathbf{x}) = \begin{pmatrix} v \cos \mathbf{x}_3 \\ v \sin \mathbf{x}_3 \\ 0 \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} 0 \\ 0 \\ 1/r \end{pmatrix}, \quad \rho = -\frac{1}{\alpha^2}.$$

The following parameters were used in the numerical simulations: $\mathbf{x}_0 = (-9 \ 0 \ 0)^\top$, $t_0 = 0$, $\mathbf{x}_f = (9 \ 0 \ 0)^\top$, $t_f = 10$, $dt = 0.1$, $v = 2.0$.

5.1 Example 1: Single-slit Obstacle

The objective in this problem is to find trajectories for the vehicle in a square environment with a box-like obstacle having a single slit. The trajectories computed by the PI-RRT algorithm at different stages are shown in Figure 1. The initial state is plotted as a yellow square and the goal region is shown in blue with magenta border (right-most). The computed path by the RRT algorithm following the unforced dynamics is shown in yellow. The locally sampled trajectories which are bundled around the yellow trajectory are shown in different colors. The trajectory of the vehicle due to execution of the control policy for some finite time horizon is shown in magenta.

To understand how the intensity of the noise level affects the patterns of the trajectories of the system, we run the algorithm and analyzed the situation for three different cases, $\alpha = 0.25, 0.5$ and 1.0 corresponding to low, medium and high intensity noise

levels in the control channel. As shown in Figure 1 (a)-(c), the PI-RRT algorithm computes trajectories that pass through the slit most of the time when there is low intensity noise in the control channel. As a first step, the PI-RRT algorithm computes a baseline trajectory using the RRT algorithm. The vertices and the edges of the tree computed by the RRT algorithm are shown in green and blue colors, respectively. During the simulations, it was observed that this baseline trajectory does not necessarily pass through the slit. The RRT algorithm sometimes returns a baseline trajectory that passes close by the upper or the lower sections of the obstacle due to both the noise which is observed in the dynamics and the randomized nature of the algorithm itself. The PI-RRT algorithm then samples a bundle of trajectories around the baseline trajectory in order to compute the variation term for the new control input. The new control input is computed by summing up the baseline control policy returned by the RRT algorithm and the variation term, which is the weighted average of the contribution of each locally sampled trajectory. These weights are computed by using the cost information of each locally sampled trajectory. We observed that the distribution of the trajectories, which pass close to the upper or lower corners or through the slit, changes as the intensity of the noise increases. For higher intensity of the noise, the PI-RRT algorithm computes trajectories which do not pass through the slit but rather pass close to the upper or lower corners. This change in the distribution of trajectories is shown in Figure 1 (d)-(f) for medium intensity noise and in Figure 1 (g)-(i) for high intensity noise.

5.2 Example 2: Double-slit Obstacle

Next, we consider a more challenging motion planning problem. In this case, there are two slits on the obstacle block and the length of the slits is longer than in the previous example. The longer length of the slits results in a higher probability of collision while traversing through the slit, which makes the motion planning problem more challenging.

A study was performed in order to compare the performance of the PI-RRT algorithm with the RRT algorithm. No variation term in the control input was computed for the RRT algorithm, and it was simply executed in a receding horizon fashion. All algorithms were run for 6000 iterations to find a baseline trajectory. The results over 100 trials are shown in Figures 2, 3 and 4. The trajectories that result in collision are plotted in Figure 2 (a), (d) for the low noise level, Figure 3 (a), (d) for the medium noise level, and Figure 4 (a), (d) for the high noise level for the RRT and PI-RRT algorithms, respectively. Also, the distribution of collision-free trajectories is plotted in Figure 2 (c), (f) for the low noise level, Figure 3 (c), (f) for the medium noise level, and Figure 4 (c), (f) for the high noise level for the RRT and PI-RRT algorithms, respectively. The distribution of trajectories and the number of trajectories which result in a collision are summarized in Table I. Under the ‘Success’ column, the rows of the table contain the number of collision-free trajectories which pass through the bottom corner, bottom slit, top slit and top corner of the block. As shown in Table 1, the PI-RRT computes safer control policies which reduce the risk of having a collision. On the other hand, both the RRT and the PI-RRT compute trajectories that are almost equally distributed over both slits.

In summary, it was observed that the behaviors of both algorithms are similar for the case with high noise level. As the noise level decreases, most of the failed cases,

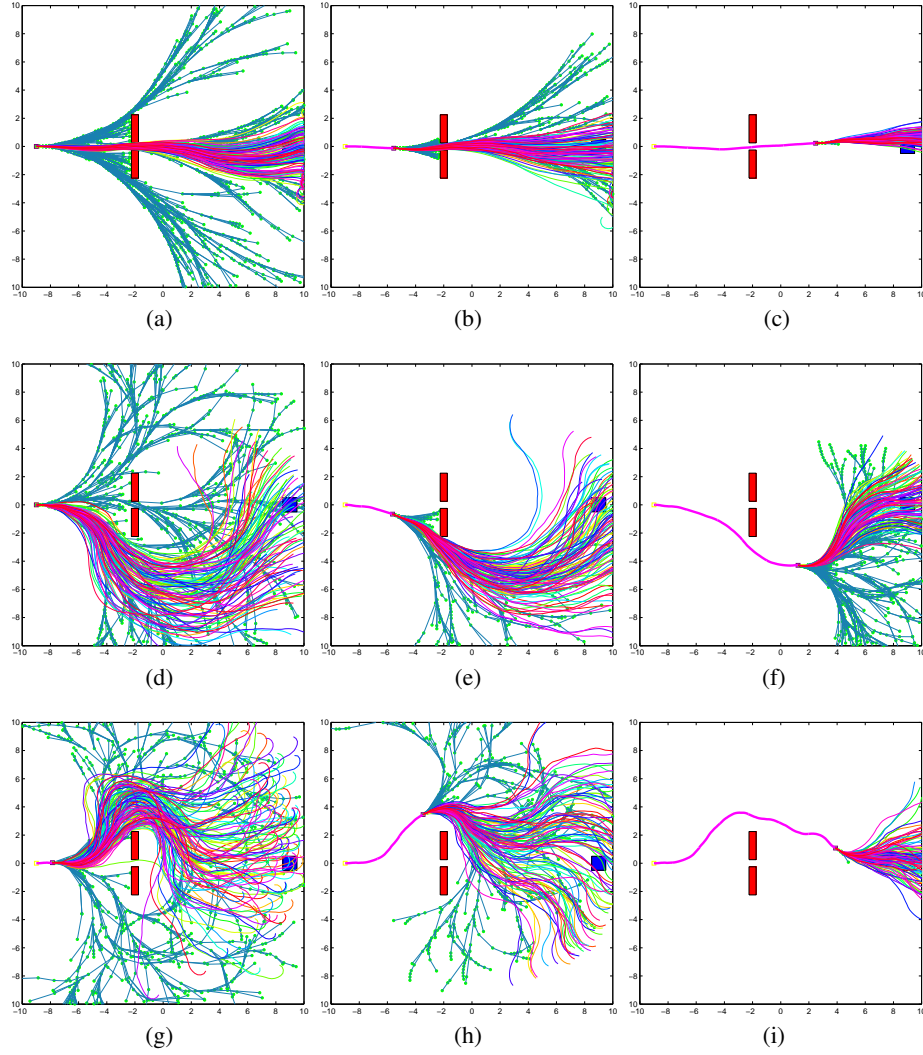


Fig. 1. The trajectories computed by the PI-RRT algorithm for stochastic optimal control of the kinematic car model under different levels of noise injected to the control channel: (a)-(c) is with $\alpha = 0.25$, (d)-(f) is with $\alpha = 0.50$, and (g)-(i) is with $\alpha = 1.0$.

not surprisingly, occur when the algorithms try to compute a path that passes through the slits. Our simulation results demonstrate that the PI-RRT algorithm tends to compute trajectories that have larger clearance from obstacles and hence outperforms the standard RRT algorithm, resulting in a smaller failure rate.

Table 1. Monte-Carlo Results for Double-Slit Obstacle

Algorithm	$\alpha = 0.25$		$\alpha = 0.50$		$\alpha = 1.00$	
	Success	Fail	Success	Fail	Success	Fail
RRT	0 24 20 0	56	23 8 11 27	31	48 0 0 44	8
PI-RRT	0 44 45 0	11	35 9 8 37	11	47 0 0 49	4

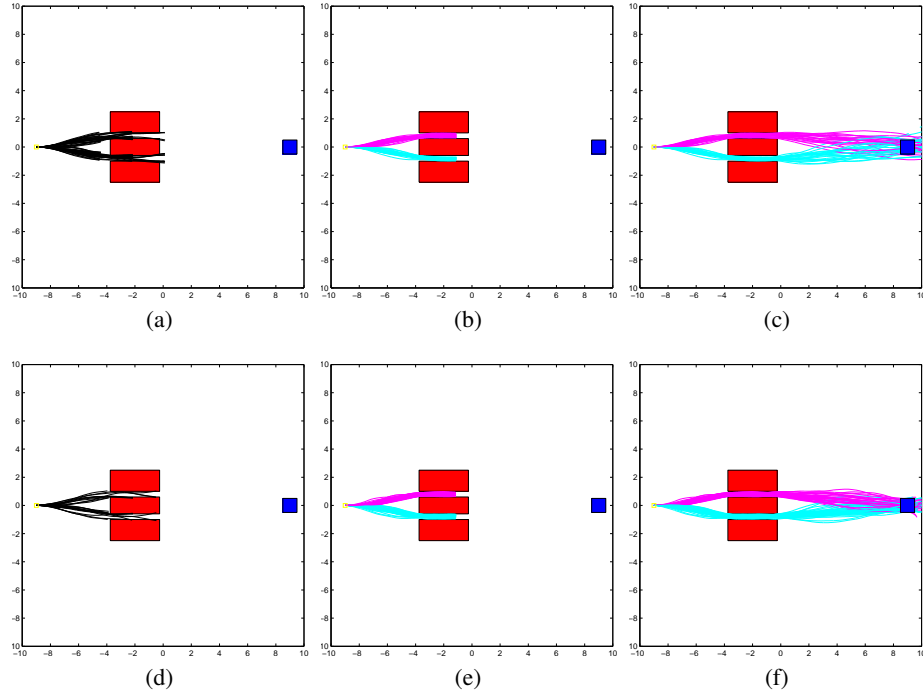


Fig. 2. Distribution of trajectories for kinematic car model under low intensity of noise injected to the control channel ($\alpha = 0.25$) is shown in (a)-(c) for the RRT algorithm, and in (d)-(f) for the PI-RRT algorithm. The trajectories which hit the obstacles are shown in (a), (d). The collision-free trajectories at an intermediate stage are shown in (b), (e), and at the final stage are shown in (c), (f).

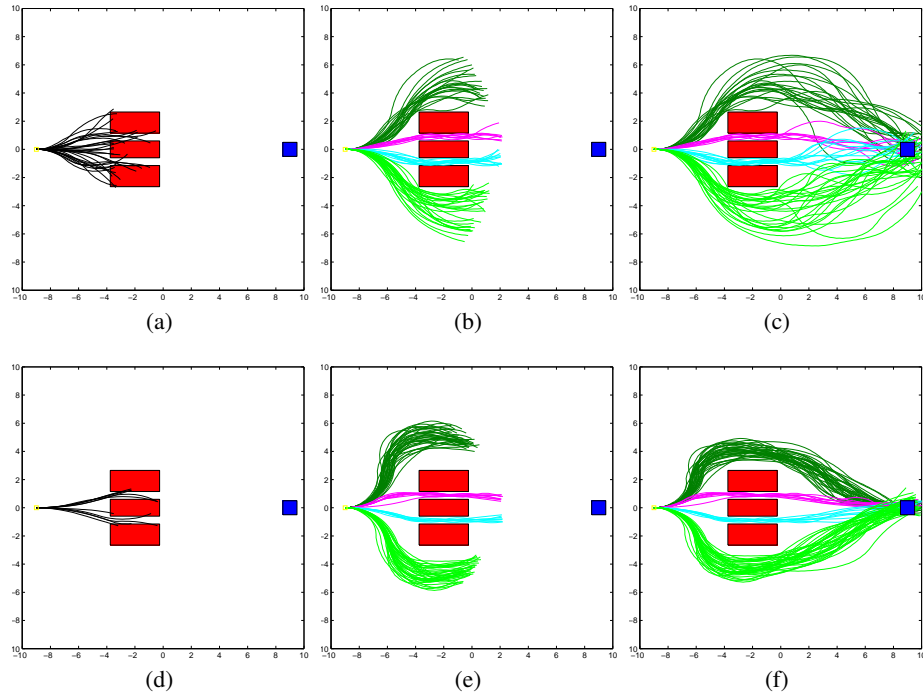


Fig. 3. Distribution of trajectories for kinematic car model under low intensity of noise injected to the control channel ($\alpha = 0.50$) is shown in (a)-(c) for the RRT algorithm, and in (d)-(f) for the PI-RRT algorithm. The trajectories which hit the obstacles are shown in (a), (d). The collision-free trajectories at an intermediate stage are shown in (b), (e), and at the final stage are shown in (c), (f).

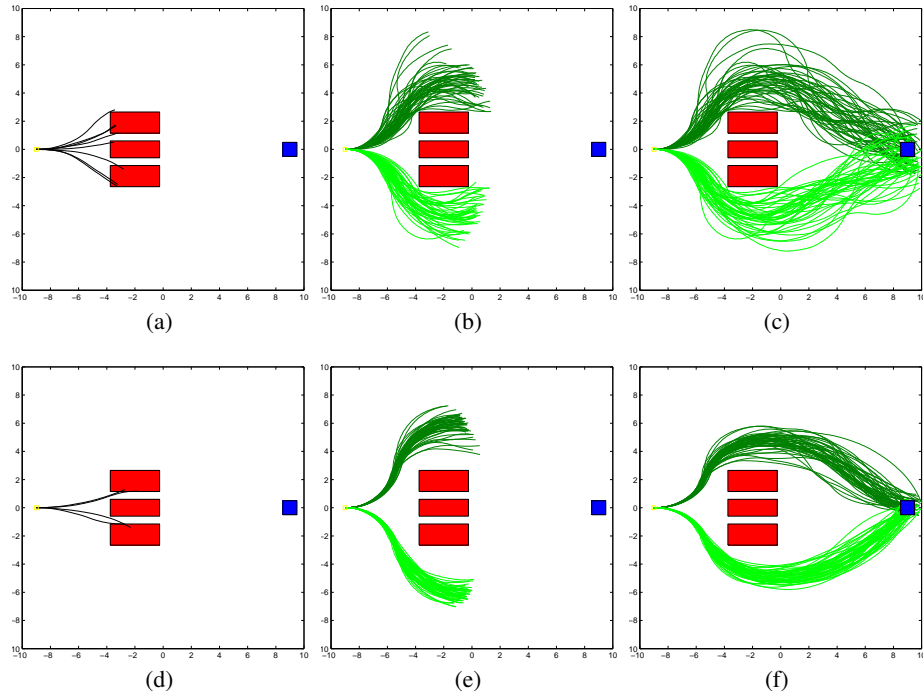


Fig. 4. Distribution of trajectories for kinematic car model under low intensity of noise injected to the control channel ($\alpha = 1.0$) is shown in (a)-(c) for the RRT algorithm, and in (d)-(f) for the PI-RRT algorithm. The trajectories which hit the obstacles are shown in (a), (d). The collision-free trajectories at an intermediate stage are shown in (b), (e), and at the final stage are shown in (c), (f).

6 Conclusion

In this paper, the PI-RRT algorithm is proposed in order to solve a class of stochastic optimal control problems. The proposed approach makes a novel connection between incremental sampling-based algorithms and path integral control. The work in this paper can be extended in several directions. First, a parallel version of the algorithm can be implemented by sampling local trajectories or computing several initial trajectories simultaneously. Second, since there exist many variants of the standard RRT algorithm, one can implement different sampling-based algorithms to compute initial trajectories and incorporate them within the path integral framework. For example, the RRT* [9,10] and the RRT[#] algorithms [1], which are both asymptotically optimal, can be used to compute bundles of good initial trajectories in a single pass; however, such an algorithm would require more elaborate computations for implementing the steering function, e.g., backward integration of a stochastic differential equation. This is part of ongoing work.

References

1. O. Arslan and P. Tsiotras. Use of relaxation methods in sampling-based algorithms for optimal motion planning. In *IEEE International Conference on Robotics and Automation*, pages 2421–2428, Karlsruhe, Germany, May 6–10, 2013.
2. J. Buchli, F. Stulp, E. Theodorou, and S. Schaal. Learning variable impedance control. *The International Journal of Robotics Research*, 30(7):820–833, April 2011.
3. H. Choset, K. Lynch, S. Hutchinson, G. Kantor, W. Burgard, L. Kavraki, and S. Thrun. *Principles of Robot Motion: Theory, Algorithms, and Implementations*. Intelligent Robotics and Autonomous Agents. The MIT Press, May 2005.
4. P. Dai Pra, L. Meneghini, and W. Runggaldier. Connections between stochastic control and dynamic games. *Mathematics of Control, Signals, and Systems*, 9(4):303–326, December 1996.
5. W. H. Fleming and W. M. McEneaney. Risk-sensitive control on an infinite time horizon. *SIAM Journal on Control and Optimization*, 33(6):1881–1915, November 1995.
6. A. Friedman. *Stochastic Differential Equations and Applications*. Dover Books on Mathematics. Dover Publications, December 2006.
7. A. J. Ijspeert, J. Nakanishi, H. Hoffmann, P. Pastor, and S. Schaal. Dynamical movement primitives: Learning attractor models for motor behaviors. *Neural Computation*, 25(2):328–373, 2013.
8. H. J. Kappen, V. Gómez, and M. Opper. Optimal control as a graphical model inference problem. *Machine Learning*, 87(2):159–182, 2012.
9. S. Karaman and E. Frazzoli. Optimal kinodynamic motion planning using incremental sampling-based methods. In *49th IEEE Conference on Decision and Control*, pages 7681–7687, Atlanta, Georgia, Dec. 15–17, 2010.
10. S. Karaman and E. Frazzoli. Sampling-based algorithms for optimal motion planning. *The International Journal of Robotics Research*, 30(7):846–894, 2011.
11. S. M. LaValle. *Planning Algorithms*. Cambridge University Press, 2006.
12. S. M. LaValle and J. J. Kuffner, Jr. Rapidly-exploring random trees: Progress and prospects. In B. R. Donald, K. Lynch, and D. Rus, editors, *New Directions in Algorithmic and Computational Robotics*, pages 293–308. 2001.

13. B.K. Øksendal. *Stochastic differential equations : An introduction with Applications*. Springer, Berlin ; New York, 6th edition, 2003.
14. R. F. Stengel. *Optimal Control and Estimation*. Dover Publications, New York, 1994.
15. F. Stulp, E.A. Theodorou, and S. Schaal. Reinforcement learning with sequences of motion primitives for robust manipulation. *IEEE Transactions on Robotics*, 28(6):1360–1370, 2012.
16. N. Sugimoto and J. Morimoto. Phase-dependent trajectory optimization for CPG-based biped walking using path integral reinforcement learning. In *IEEE-RAS International Conference on Humanoid Robots*, pages 255–260, Bled, Slovenia, Oct. 26–28, 2011.
17. E. Theodorou, J. Buchli, and S. Schaal. A generalized path integral approach to reinforcement learning. *Journal of Machine Learning Research*, (11):3137–3181, 2010.
18. E. Theodorou, J. Buchli, and S. Schaal. Reinforcement learning of motor skills in high dimensions: A path integral approach. In *IEEE International Conference on Robotics and Automation*, pages 2397–2043, Anchorage, Alaska, May 3–8, 2010.
19. E. Theodorou, D. Krishnamurthy, and E. Todorov. From information theoretic dualities to path integral and kullback-leibler control: Continuous and discrete time formulations. In *The Sixteenth Yale Workshop on Adaptive and Learning Systems*, page ???, New Haven, Connecticut, June 5–7, 2013.
20. E. Theodorou, D. Krishnamurthy, and E. Todorov. Time-varying nonlinear policy gradients. In *52nd IEEE Conference on Decision and Control*, pages 7765–7770, Florence, Italy, Dec. 10–13, 2013.
21. E. Todorov. Efficient computation of optimal actions. *Proceedings of the National Academy of Sciences*, 106(28):11478–11483, 2009.
22. J. Yang and J. H. Kushner. A Monte Carlo method for sensitivity analysis and parametric optimization of nonlinear stochastic systems. *SIAM Journal in Control and Optimization*, 29(5):1216–1249, 1991.