



HAL
open science

Revisiting Continuous-Time Reinforcement Learning. A Study of HJB Solvers Based on PINNs and FEMs

Alena Shilova, Thomas Delliaux, Philippe Preux, Bruno Raffin

► To cite this version:

Alena Shilova, Thomas Delliaux, Philippe Preux, Bruno Raffin. Revisiting Continuous-Time Reinforcement Learning. A Study of HJB Solvers Based on PINNs and FEMs. EWRL 2023 Workshop, Sep 2023, Brussels, Belgium. hal-04705270

HAL Id: hal-04705270

<https://hal.science/hal-04705270v1>

Submitted on 22 Sep 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Revisiting Continuous-Time Reinforcement Learning. A Study of HJB Solvers Based on PINNs and FEMs

Alena Shilova

Univ. Lille, Inria, CNRS, Centrale Lille, UMR 9189 - CRISTAL
Lille, France
alena.shilova@inria.fr

Thomas Delliaux

Univ. Lille, Inria, CNRS, Centrale Lille, UMR 9189 - CRISTAL
Lille, France
thomas.delliaux@inria.fr

Philippe Preux

Inria, Univ. Lille, CNRS
Lille, France
philippe.preux@inria.fr

Bruno Raffin

Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LIG
Grenoble, France
bruno.raffin@inria.fr

Abstract

Despite recent advances in Reinforcement Learning (RL), the Markov Decision Processes are not always the best choice to model complex dynamical systems requiring interactions at high frequency. Being able to work with arbitrary time intervals, Continuous Time Reinforcement Learning (CTRL) is more suitable for those problems. Instead of the Bellman equation operating in discrete time, it is the Hamiltonian Jacobi Bellman (HJB) equation that describes value function evolution in CTRL. Even though the value function is a solution of the HJB equation, it may not be its unique solution. To distinguish the value function from other solutions, it is important to look for the viscosity solutions of the HJB equation. The viscosity solutions constitute a special class of solutions that possess uniqueness and stability properties. In this paper, we bring together the formalism of viscosity solutions and practical methods for finding them. We also propose a novel way of training neural networks to obtain viscosity solutions. Finally, we do a comparison of those methods with discrete time RL algorithms to emphasize the benefits of considering the continuous time setting. This paper aims at providing the necessary theoretical basis for working with CTRL and setting a few possible directions for future research.

1 Introduction

Reinforcement learning (RL) is getting more and more attention. Most of state-of-the-art RL methods are designed to work with Markov Decision Processes (MDPs) and, in particular, rely on a discrete time assumption. Even though this assumption is not that restrictive in tasks such as Atari games and go, it is no longer valid in complex problems such as driving cars or finance trading. Such problems are usually described by a dynamical system where discrete time RL (DTRL) often struggles to provide accurate control within short time intervals due to several reasons [8, 32, 23]. As DTRL algorithms, such as Q-learning or policy gradients, typically operate in discrete time steps, it becomes challenging to capture the nuances of continuous state dynamics within small intervals. Even when

the discretization step in time is set to be small, DTRL fails to learn the optimal policy as exploration cannot be done efficiently.

Continuous Time Reinforcement Learning (CTRL) derives from the optimal control theory and thus provides a promising theoretical framework to tackle aforementioned shortcomings of DTRL. CTRL allows us to have control over the discretization of time and it provides a smoother control. The key ingredient in CTRL is the Hamilton-Jacobi-Bellman (HJB) equation [8, 25], which is a Partial Differential Equation (PDE). This PDE describes the evolution of the value function with respect to a state of the system, which in turn evolves continuously with time. In principle, the value function, from which an optimal control policy may be deduced, is a solution of the HJB equation. However, finding the value function still remains a challenging task. The HJB equation may have multiple solutions in a class of continuous functions [25]. Distinguishing the value function from the other solutions of the HJB equation relies on the search for the viscosity solution that corresponds to a certain solution that satisfies uniqueness and stability properties.

However, there is no guarantee in general that the solution found by existing methods is the good solution, *i.e.* the viscosity solution. Moreover, the existing methods for solving the HJB equation are limited and mainly applicable to specific cases such as Linear Quadratic Regulator (LQR) problems or variational problems [9]. There exist approaches that can be applied in general case, such as Finite Difference (FD), Finite Element Methods (FEM) [10], and dynamic programming methods [25] that transform a CTRL problem back to a DTRL problem thanks to the discretization. However, their effectiveness is hindered by the exponentially growing algorithmic complexity with respect to the state space dimensionality. Moreover, they are also affected by discretization error.

Conversely, neural networks (NNs) are emerging as PDE solvers that can cope with the curse of dimensionality [27]. Although Neural Networks have shown a great potential in various domains [19, 15], applying them to solve the HJB equation requires caution because of the non-uniqueness issue. To address this challenge, a promising direction is to leverage the formalism of viscosity solutions in conjunction with neural networks.

In this paper, we revisit CTRL approaches and analyze how the latest advances in deep learning can be applied and adjusted to solve the HJB equation in the viscosity sense. Training a CTRL agent remains technically very challenging and it is not yet possible to cope with tasks that are as large as those dealt by DTRL. Yet, it is our goal to make CTRL agents suitable for solving larger tasks and we provide some recipes to make it eventually as easy and accessible as training a DQN or PPO for DTRL. For that purpose, we need to understand the theory on which CTRL is rooted and, in particular, the HJB equation (similarly, DTRL development is due to understanding of the Bellman equation) and finally put theory into practice.

Further, we focus on the case of deterministic environments with known dynamics. The extension to stochastic environments and unknown dynamics is possible (see [32, 35]), but left for the future work. In particular, we consider how NN-based methods can be combined to find viscosity solutions, offering a novel training scheme. We start with the analysis of the existing literature in Section 2. Then, we carefully introduce the definitions and notations related to CTRL in Section 3.1, describe the HJB equation in Section 3.2, define viscosity solutions in Section 3.3 and discuss boundary conditions in Section 3.4. Then, different ways of integrating neural networks in the process of solving the HJB equation are discussed in Section 4. Finally, we demonstrate performances of some approaches from Section 4 on the inverted pendulum, a classical use case from RL, in Section 5.

2 Related Works

Among the first papers to study CTRL are [8, 25, 5]. Those works introduced the HJB equation as a key equation for finding the optimal policy. In [8], different discretization schemes and algorithms are analysed for computing the value function, including continuous TD(λ) and continuous Actor-Critic. In his Ph.D. [5], Coulom studied applicability of Doya’s methods [8] to a large class of control problems. Munos [25] tackled CTRL by the study of viscosity solutions and their properties. He demonstrated the challenges of solving the HJB equation such as the non-uniqueness of solutions and inequality boundary conditions. In addition, convergent numerical schemes based on dynamic programming were derived to approximate the value function. One of the problems of numerical schemes is the curse of dimensionality. To mitigate this problem, sparse grids [14] can be used instead

of naive uniform grid. In this article, we take the formalism proposed in [25], but we consider neural network based approaches to find viscosity solutions, while [25] considers tabular algorithms.

NNs can be applied for solving the HJB equation [24, 18, 4, 30, 20, 12, 1]. It was first demonstrated in [24]. In [30], the training is regularized to avoid finding “bad” solutions. Moreover, the same work raised the problem of falling in a local minimum of the squared HJB residual, and some solutions to avoid it were proposed. Compared to previous papers, [18, 1] have emphasized the robustness of the resulting controller and the stability of the proposed algorithm. Adapting the former methods for more practical cases such as the inverted pendulum and the cartpole was done in [20]. We consider similar approaches as [30, 20], but unlike them we use the formalism proposed in [25] to develop an approach that can converge to viscosity solutions.

Several extensions to the classical HJB equation were also introduced such as HJB for an explorative reward function [32], the soft HJB equation with maximum entropy regularization [17, 11] and the distributional HJB equation [33]. Those works extend the existing theory to other definitions of the value function, however experiments are conducted on a limited set of simple problems. Furthermore, it is also possible to extend the HJB equation to continuous-time partially observable Markov decision processes (CTPOMDPs). In [2], they proposed a formalism to describe CTPOMDPs, including the CTPOMDP HJB equation. [16] tackles the problem of CTRL by adapting the well-known DQN algorithm to this framework. A definition for the Q-function in the continuous case is given and the “HJB equation for the Q-function” is derived, which results in a DQN-like algorithm for the semi-discrete time setting. However, this approach is limited to Lipschitz continuous control. The HJB equation was used to improve DTRL algorithms like PPO in [23] and it resulted in a significant improvement on MuJoCo, proving that HJB loss is better adapted for learning value functions of dynamical systems.

There were some attempts to propose alternative ways for solving the HJB equation. An other NN approach to approximate the value function based on Pontryagin’s maximum principle has been proposed [26]. Conversely, [7] considers some special neural network architectures that can work with min-plus algebra, though this approach is suitable only for some optimal control problems.

Another line of research exploits the continuous time formulation to do more accurate Monte Carlo estimations of value function [22, 21, 35] rather than using the HJB equation. [22] adapted the fitted value iteration algorithm to the CTRL problem. In [35], a model-based algorithm was introduced, which aims at solving the problem in an actor-critic manner. Bayesian neural ODEs are used in order to learn the dynamics of the system.

3 Hamiltonian-Jacobi-Bellman Equation

3.1 Formalism for Reinforcement Learning in the Continuous Case

In this work, we consider infinite-horizon deterministic problems of finding the optimal control in the dynamical systems, the dynamics and rewards of which are known. The extension to unknown dynamics and rewards is left for future works.

In what follows, we denote with $O \subseteq \mathbb{R}^d$, the open set of *controlable* states of our system, and then the *admissible* control $u \in U$ is the one that keeps the state trajectory inside the domain O . In this case, $U \subset \mathbb{R}^m$ corresponds to the action/control space. We assume that U is bounded. Further, we use $g \in C(O)$ to denote that g is a continuous function on O , while $g \in C^1(O)$ that g is a continuously differentiable function on O .

The main difference between the continuous time and the discrete time cases is that transitions depend on time t , which is a continuous variable, generating trajectories of states continuously in time. To not confuse with Semi-Markov Decision Processes [29], that is an MDP with variable-duration actions. More formally, in the continuous case the states do not form a sequence $(X_t)_t$ but a trajectory $x : \mathbb{R}_+ \rightarrow \bar{O} \subset \mathbb{R}^d$. If there is a time when the state trajectory reaches and crosses the boundary ∂O , then it means that the system has exited the control domain, which corresponds to reaching a terminal state in the discrete time case. The exit time is denoted with τ , $0 < \tau \leq \infty$ (the case $\tau = \infty$ is also possible in the situation where exiting $\bar{O} = O \cup \partial O$ is not optimal, see Section 3.4). Similarly, the actions are defined for any $t < \tau$: $u : \mathbb{R}_+ \rightarrow U$. The dynamics of the environment (the transition

function) is defined through the following ordinary differential equation (ODE):

$$\frac{dx(t)}{dt} = f(x(t), u(t)) \quad \text{or} \quad x(t') = x(t) + \int_t^{t'} f(x(s), u(s)) ds, \quad t' > t \quad (1)$$

The function $f : \bar{O} \times U \rightarrow \mathbb{R}^d$ is called the state dynamics. To compare, $x(t)$ is only defined at times $\{0, dt, 2dt, \dots\}$ in discrete time case, where dt is a time step. While the state x is inside the control domain O (i.e. $t < \tau$), the reward $r : \bar{O} \times U \rightarrow \mathbb{R}$ is received, conversely once the system exits the control domain, then the exit reward $R : \partial O \rightarrow \mathbb{R}$ is obtained. Given the initial state $x(0) = x_0$, we define the cumulative discounted reward for a given control function $u(t)$ as follows

$$J(x_0; u(t)) = \int_0^\tau \gamma^t r(x(t), u(t)) dt + \gamma^\tau R(x(\tau)), \quad (2)$$

where $\gamma \in [0, 1)$ is a discount factor. J is called the reinforcement functional. In the continuous time case, the value function is defined as:

$$V(x) = \sup_{u(t)} J(x; u(t)). \quad (3)$$

Here, the value function is defined as the maximum value that J can take when we start at x at time $t = 0$. Our goal is to find an optimal policy $\pi : \bar{O} \rightarrow U$, such that $u^*(t) = \pi(x^*(t))$ for any t , where $u^*(t)$ and $x^*(t)$ are control and state of the optimal trajectory respectively, such that $V(x) = J(x; u^*(t))$.

Furthermore, we introduce additional assumptions on functions $f(x, u)$ and $r(x, u)$ and domain O .

Assumption 3.1 (i) U is bounded,

(ii) f, r are bounded, f is continuous on $\mathbb{R}^d \times U$ and r is uniform continuous on $\mathbb{R}^d \times U$,

(iii) there exists L_f , such that $\|f(x, u) - f(y, u)\| \leq L_f \|x - y\|$ for any $x, y \in \mathbb{R}^d$.

Assumption 3.2 For any $x \in \partial O$ and its normal vector $\eta(x)$

(i) $\exists u(x) \in U$ with $f(x, u(x))^T \eta(x) < 0$;

(ii) $\exists u(x) \in U$ with $f(x, u(x))^T \eta(x) > 0$.

Assumption 3.3 (Regularity condition) There exist $\epsilon_0, r > 0$ and $\hat{\eta}(x)$, a bounded and uniform continuous map of \bar{O} , satisfying

$$B(x + \epsilon \hat{\eta}(x), r\epsilon) \subset O, \quad \forall x \in O, \epsilon \in (0, \epsilon_0] \quad (4)$$

with $B(x, r) = \{y \in \mathbb{R}^d : \|x - y\| < r\}$.

Assumptions 3.1-3.3 are necessary for proving the known theoretical results stated in Sections 3.3-3.4 and Appendix A.2. Even though some assumptions can be relaxed in practice, the case of non-continuous dynamics (i.e. $f \notin C(O)$) is much more complicated and requires other approaches than the ones considered in the paper. As the latter case is more general and interesting in real life applications, it should be studied in the future work.

3.2 Hamilton-Jacobi-Bellman Equation

We can prove that the value function defined in Eq. (3) satisfies the following result (see [9]):

Theorem 3.1 (Hamilton-Jacobi-Bellman). *If the value function V is differentiable at x , then the Hamilton-Jacobi-Bellman (HJB) equation holds at any $x \in O$:*

$$V(x) \ln(\gamma) + \sup_{u \in U} \{ \nabla_x V(x)^T f(x, u) + r(x, u) \} = 0. \quad (5)$$

Computing the value function allows us to find the optimal control. Indeed, from the value function we can define a feed-back control policy $\pi : \bar{O} \rightarrow U$ such that $\pi(x(t)) = u^*(t)$ by setting:

$$\pi(x) \in \operatorname{argsup}_{u \in U} \{ \nabla_x V(x)^T f(x, u) + r(x, u) \} \quad (6)$$

Therefore, it is crucial to find an efficient way to compute V to get the optimal control u^* . However, solving the HJB equation involves several challenges (see [25]). First, the value function V is often non-smooth, and only continuous on O . Second, the HJB equation may have multiple generalized continuous solutions. Third, it is common that the HJB equation (5) has to be solved under inequality boundary conditions. To address those points, we introduce the viscosity property in Sections 3.3-3.4.

3.3 Viscosity

Here, we present a crucial and yet complex notion of viscosity (refer to Appendix A.1 for an intuition behind). Let $H(x, W, \nabla_x W) = -W(x) \ln \gamma - \sup_{u \in \mathcal{U}} [\nabla_x W(x)^T f(x, u) + r(x, u)]$, then the HJB equation can be expressed as

$$H(x, W, \nabla_x W) = 0. \quad (7)$$

Using these notations, let us define what a viscosity solution is.

Definition 3.1 (Viscosity solution)

- $W \in C(O)$ is a viscosity subsolution of the HJB equation in O if $\forall \psi \in C^1(O)$ and $\forall x \in O$ local maximum of $W - \psi$ such that $W(x) = \psi(x)$, we have:

$$H(x, \psi(x), \nabla_x \psi(x)) \leq 0$$

- $W \in C(O)$ is a viscosity supersolution of the HJB equation in O if $\forall \psi \in C^1(O)$ and $\forall x \in O$ local minimum of $W - \psi$ such that $W(x) = \psi(x)$, we have:

$$H(x, \psi(x), \nabla_x \psi(x)) \geq 0$$

- If $W \in C(O)$ is a viscosity subsolution and a supersolution then it is a viscosity solution.

Viscosity allows us to separate “good” solutions of the HJB equation from “bad” ones. Viscosity solutions were first introduced in [6] and they are proven to be unique for multiple types of PDEs including the HJB equation. Next, we state a few useful results proven in [9].

Theorem 3.2 *Given Assumption 3.1, the value function V is uniformly continuous and bounded in \mathbb{R} and then it is a unique viscosity solution of the HJB equation (7).*

This theorem is given for the case when $O = \mathbb{R}^d$ and therefore there is no boundary condition. The other cases are discussed in the next sections.

The proof of Theorem 3.2 consists of several parts. First, one can prove that under Assumption 3.1, the value function is indeed uniform continuous and bounded. Then, one can show that it is a viscosity solution due to the dynamic programming principle that holds in continuous time case as well. The uniqueness comes from the comparison principle. It states that under Assumptions 3.1, if W and V are viscosity subsolution and supersolution respectively and are bounded and uniformly continuous functions, then $W \leq V$. The comparison principle implies that if such W and V are viscosity solutions (both subsolutions and supersolutions), then $W \leq V$ and $W \geq V$, therefore $W \equiv V$. Thus, V is a unique viscosity solution of the HJB equation in \mathbb{R}^d .

3.4 Discussion of Boundary Conditions

Now let us consider the situation of $O \subset \mathbb{R}^d$. Then it becomes necessary to take into account boundary conditions. To guarantee the existence of a continuous value function that satisfies the boundary conditions, it is important to impose Assumption 3.2 that verifies that the boundary is reachable from the interior and the interior is reachable from the boundary. The formulation of the boundary conditions can vary depending on the problem type. In this section, we concentrate on the case when the system should remain within the domain O , *i.e.* $\tau = \infty$. This is a common requirement in the dynamical systems. For example, we may want to limit the maximal angular speed in the inverted pendulum to mitigate the risk of wearing off or breaking the mechanism. Other boundary conditions are considered in Appendix A.2.

To provide the mathematical formulation, let us denote the external normal vector at point $x \in \partial O$ as $\eta(x)$, then this boundary can be expressed as $f(x, u^*(x))^T \eta(x) \leq 0$ for any $x \in \partial O$. The optimal

policy is not known a priori and thus it is hard to verify this constraint. In [9, 28], it was shown that it can be reformulated as:

$$-H(x, W, \nabla_x W + \alpha \eta(x)) \leq 0 \quad \forall \alpha \leq 0, x \in \partial O. \quad (8)$$

This allows to extend Definition 3.1 to the points at the boundary $x \in \partial O$.

Definition 3.2 (Constrained viscosity solution) $W \in C(\bar{O})$ is called a constrained viscosity solution of the HJB equation (5) if it is a viscosity subsolution in O and a viscosity supersolution in \bar{O} , i.e. if $\forall \psi \in C^1(\bar{O})$ and $\forall x \in \bar{O} \cup \arg \min\{(W - \psi)(x) : x \in \bar{O}\}$ with $W(x) = \psi(x)$, we have:

$$H(x, \psi(x), \nabla_x \psi(x)) \geq 0.$$

This definition is obtained from Eq. (8) (see [9]). Similarly to Theorem 3.2, we can prove that there exists a continuous value function provided that Assumption 3.2(i) holds and the set of admissible actions is not empty for any state of the system. This value function is a constrained viscosity solution. Under the additional Assumption 3.3, there is only one constrained viscosity solution (see [28]).

Some of the assumptions can be relaxed and it is possible to obtain more general uniqueness results (see [9, 3, 13]). However, the assumptions mentioned earlier are verified for the large class of control problems that appear in practice, like classical control or MuJoCo [31].

4 Possible Approaches

In this section we discuss several ways of solving the HJB equation (5). In what follows, we assume that the dynamics of the state determined by $f(x, u)$ is known and the control space is discrete, the case of unknown $f(x, u)$ and the continuous control space is left for future work. With regards to boundary conditions, we illustrate further methods for solving HJB under boundary conditions from Section 3.4 (see Eq. (8)), the extension to other cases (see Appendix A.2) is straightforward.

4.1 Dynamic Programming

Several solvers exist such as Finite Difference method (FD) or Finite Element Method (FEM). These methods require the discretization of the domain (a grid for FD or a triangulation for FEM). Moreover, the dynamics of the system has to be available to compute the solution of the PDE. The work of [25] establishes the connection between solving the HJB equation and the classical reinforcement learning framework by deriving an MDP from the discretization of the HJB equation, using either FD or FEM schemes (see Appendix A.3.1 for a short summary of the method). We show the performance of the FEM based dynamic programming in Section 5. However, despite many efforts, we were not able to make the algorithm based on FD work in our experiments, therefore, its results are missing.

Despite convergence guarantees provided in [25], the main problem of these methods is that they are mesh-dependent. Thus, they suffer from the curse of dimensionality, i.e. the computational cost increases exponentially with the dimension of the problem. For example, if the state space of the environment is a subdomain of \mathbb{R}^4 then a naive approach that divides all axes uniformly in N parts results in N^4 states to handle. If $N = 32$, which may not be sufficient to solve the problem, it leads to 2^{20} states for the algorithm to process. In our experiments performed on a single CPU, we are not able to go beyond $N = 20$ for $O \subseteq \mathbb{R}^4$. There exist smarter ways of defining the grid such as a sparse grid [14], which may help to use the available compute and memory budget more efficiently, but we leave sparse grids usage out of scope of this paper.

4.2 Approach Based on PINNs

The idea of PINNs was proposed in [27] and applied to some simple PDEs like one-dimensional nonlinear Schrödinger equation, but not for optimal control. In PINNs framework, the neural network acts as a solution of the PDE that needs to be solved. Being randomly initialized, the neural network is gradually fit to satisfy the PDE and its corresponding boundary conditions with the help of optimization and automatic differentiation that allows to compute precise derivatives. For example, the solution of the equation $W'_x - W = 0$ for $W \in C^1(O = [0, 1])$ with the boundary condition $W(0) = 1$ can be found by minimizing the loss $\min_W \{\|W'_x - W\|_2^2 + \lambda \|W(0) - 1\|_2^2\}$. The first term of this loss is called a PDE loss and the second term a boundary loss, where λ is a hyperparameter that

weighs a boundary loss against a PDE loss. PINNs can be trained in a self-supervised manner as a dataset can be generated by simply drawing random samples from the domain O . Still, if the solution is known at some points of the domain, then the training can be augmented with a data-driven loss. Refer to Appendix A.3.2 for a more detailed introduction to PINNs.

Let $W(\cdot, \theta)$ be a neural network that computes the solution of Eq. (5) with θ being its parameters. In this section, we describe how to train $W(\cdot, \theta)$ in PINNs-like manner. To apply PINNs to the HJB equation, it is necessary to search among viscosity solutions to avoid non-uniqueness issues (see Section 3). Checking the conditions of Definition 3.1 is not feasible in practice. Instead, we use the next property of viscosity solutions.

Lemma 4.1 (Stability) *Let W^ϵ be a viscosity subsolution (resp. a super solution) of*

$$W^\epsilon(x) + F^\epsilon(x, W^\epsilon(x), \nabla_x W^\epsilon(x), \nabla_x^2 W^\epsilon(x)) = 0$$

in O . Suppose that F^ϵ converges to F uniformly on every compact subset of its domain, and W^ϵ converges to W uniformly on compact subsets of \bar{O} . Then W is a viscosity subsolution (resp. a supersolution) of the limiting equation.

This lemma is proven in [9]. In our case, we are interested in equation:

$$H(x, W^\epsilon(x), \nabla_x W^\epsilon(x)) = \epsilon \nabla_x^2 W^\epsilon(x), \quad (9)$$

where the left hand side is the same as in Eq. (5), while the right hand side depends linearly on $\epsilon > 0$. Therefore, $F^\epsilon(x, W^\epsilon(x), \nabla_x W^\epsilon(x), \nabla_x^2 W^\epsilon(x)) = -\frac{1}{\ln \gamma} H(x, W^\epsilon(x), \nabla_x W^\epsilon(x)) - W^\epsilon(x) + \frac{\epsilon}{\ln \gamma} \nabla_x^2 W^\epsilon(x)$ and $F(x, W(x), \nabla_x W(x), \nabla_x^2 W(x)) = -\frac{1}{\ln \gamma} H(x, W(x), \nabla_x W(x)) - W(x)$ in Lemma 4.1. In [9], it is shown that Eq. (9) has a unique smooth solution $W^\epsilon(x)$, *i.e.* it admits a classical solution, which is a viscosity solution at the same time. Therefore, if $W^\epsilon(x)$ converges uniformly to $W(x)$ then $W(x)$ is a viscosity solution of the original HJB equation (5).

To ensure the uniform convergence of Eq. (9) to Eq. (5), we need to define a sequence of $\{\epsilon_n\}_{n=0}^\infty$, such that $\epsilon_n \rightarrow 0$ and $\|\epsilon_n \nabla_x^2 W^\epsilon(x)\| \rightarrow 0$ uniformly for all $x \in \bar{O}$ when $n \rightarrow \infty$. Let us denote $\delta(\epsilon, \theta) = \frac{1}{N_S} \sum_i \|\epsilon \nabla_x^2 W^\epsilon(x_i, \theta)\|^2$. Further, we use the following scheme to generate a sequence $\{\epsilon\}_{n=0}^\infty$. Given the first element of the sequence ϵ_0 , we get all the consecutive elements with the next update rule

$$\epsilon_{n+1} = \begin{cases} \frac{k_\epsilon \delta(\epsilon_n, \theta_{n-1})}{\delta(\epsilon_n, \theta_n)} \epsilon_n & \text{if } k_\epsilon \delta(\epsilon_n, \theta_{n-1}) \leq \delta(\epsilon_n, \theta_n), \\ & \text{and } \mathcal{L}(\theta_i) \geq \mathcal{L}(\theta_{i-1}) \forall i : n - n_\epsilon + 1 \leq i \leq n \\ \epsilon_n & \text{otherwise,} \end{cases} \quad (10)$$

where n corresponds to the current iteration number and $k_\epsilon \in (0, 1)$ is a coefficient defining the speed with which $\|\epsilon_n \nabla_x^2 W^\epsilon(x)\|^2$ should decrease. Moreover, n_ϵ serves as a number of iterations since θ_i does not improve the loss, *i.e.* it also specifies the minimum number of iterations with fixed ϵ to obtain $W^\epsilon(x, \theta) \approx W^\epsilon(x)$ that solves Eq. (9). We also tried more simple update rules like a linear update rule $\epsilon_{n+1} = k_\epsilon \epsilon_n$, but this resulted in much worse results, while our scheme results in $\|\epsilon_n \nabla_x^2 W^\epsilon(x)\|$ steadily going towards 0.

Finally, we define losses corresponding to Eq. (9) and Eq. (8). Let us define $\mathcal{S}_O \sim \mathcal{U}(O, N_F)$, a sample of points drawn uniformly from O of size N_F , and $\mathcal{S}_{\partial O} \sim \mathcal{U}(\partial O, N_B)$. Further, if ϵ and θ are indicated without an iteration number, then they correspond to the current iteration n . We have:

$$\mathcal{L}_O(\theta, \mathcal{S}_O) = \frac{1}{N_F} \sum_{i=1}^{N_F} (H(x_i, W^\epsilon(x_i, \theta), \nabla W^\epsilon(x_i, \theta)) - \epsilon \nabla_x^2 W^\epsilon(x_i, \theta))^2 \quad x_i \in \mathcal{S}_O, \quad (11)$$

$$\mathcal{L}_{\partial O}(\theta, \mathcal{S}_{\partial O}) = \frac{1}{N_B} \sum_{i=1}^{N_B} \left([-H(x_i, W^\epsilon(x_i, \theta), \nabla W^\epsilon(x_i, \theta)) + \alpha \eta(x_i)]^+ \right)^2 \quad x_i \in \mathcal{S}_{\partial O}, \quad (12)$$

where $[f(x)]^+ = \max\{f(x), 0\}$. In addition to PDE-related and boundary-related losses, we introduce an MSE regularization loss that should encourage uniform convergence:

$$\mathcal{L}_R(\theta, \mathcal{S}_O) = \frac{1}{N_F} \sum_{i=1}^{N_F} (W^\epsilon(x_i, \theta) - W^{\epsilon_{n-1}}(x_i, \theta_{\epsilon_{n-1}}))^2 \quad x_i \in \mathcal{S}_O \quad (13)$$

Algorithm 1 ϵ -HJBPINNs

```
Set  $\epsilon = \epsilon_0, \theta = \theta_0$ , initialize  $W^\epsilon(\cdot, \theta)$ 
for iteration  $n$  in  $\{1, \dots, \text{NB\_ITER}\}$  do
  Generate datasets  $\mathcal{D}_O(N_D)$  and  $\mathcal{D}_{\partial O}(N_D)$  of  $N_D$  states uniformly sampled from  $O$  and  $\partial O$  respectively
  for batches  $\mathcal{S}_O \in \mathcal{U}(\mathcal{D}_O, N_S)$  and  $\mathcal{S}_{\partial O} \in \mathcal{U}(\mathcal{D}_{\partial O}, N_S)$  do
    Update  $\theta_n := \theta_n - \nu \nabla_{\theta} \mathcal{L}(\theta, \mathcal{S}_O, \mathcal{S}_{\partial O})$ , where  $\mathcal{L}(\theta, \mathcal{S}_O, \mathcal{S})$  is computed with Eq. (14)
  end for
  Compute  $\epsilon_n$  using Eq. (10)
  Set  $\epsilon = \epsilon_n$  and  $\theta = \theta_n$ 
end for
```

where $W^{\epsilon_{n-1}}(x, \theta_{\epsilon_{n-1}})$ is the best function obtained for ϵ_{n-1} . The final loss is:

$$\mathcal{L}(\theta, \mathcal{S}_O, \mathcal{S}_{\partial O}) = \mathcal{L}_O(\theta, \mathcal{S}_O) + \lambda \mathcal{L}_{\partial O}(\theta, \mathcal{S}_{\partial O}) + \lambda_R \mathcal{L}_R(\theta, \mathcal{S}_O). \quad (14)$$

Bringing everything together results in Algorithm 1.

5 Experimental Results

In this section we evaluate the performance of methods from Section 4 on the inverted pendulum, which is a classical RL control task. We use a continuous-time adaptation of the inverted pendulum¹. The state space consists of angle ϕ and angular speed $\dot{\phi}$. We consider the domain $O = [-\pi, \pi] \times [-10, 10]$. The set of actions is $\{-2, 0, 2\}$. In this paper, we consider only the boundary condition from Section 3, described with Eq. (8).

5.1 PINNs Results

In this section, we give the results of the training with Algorithm 1. We tried several types of architectures, such as Multi-Layer Perceptron (MLP) and Fourier-Feature Network (FFN) [34] with different number of layers and neurons. We have obtained the best performances with FFN, consisting of 3 layers, where the first layer is of size $d \times 40$ (this size of the first layer was recommended in the original article, where d is the dimensionality of state space) and other layers are 100 neurons each. The best performing activation function is `torch.tanh`. It is important to use smooth activation functions for solutions of Eq. (9). Using non-smooth activations like `torch.relu` causes the training to fail. The hyperparameters for Algorithm 1 are given in Appendix A.4.2 together with the discussion on how they are chosen.

We have executed a training for 1000 iterations (`NB_ITER = 1000`). In Figure 1, we present the value function, control and loss maps that we have obtained after the training. We observe that PINNs are able to find the general structure of the value function, though it is less precise to approximate the non-smooth zones because of the smooth activation functions used in the neural networks. Thus, further research on how to improve this approach is necessary.

5.2 Comparison with DTRL and Dynamic Programming Algorithms

We evaluate the above methods on discrete time inverted pendulum with small $dt = 0.001^2$ to compare their performance with PPO and A2C, which are state-of-the-art DTRL algorithms. A complete discussion of our experiments with dynamic programming can be found in A.4.1. To test the algorithms, we perform 100 rollouts, each rollout being made of 5000 time steps. In the case of dynamic programming, we also vary the grid size N to test its performance (then the discretization step is $\delta = 1/N$). We report in Figure 2 (right) the mean and the standard deviation. We also present the value function and the policy obtained by the dynamic programming algorithm with a resolution of 200 by 200 ($N = 200$) in Figure 2 (left).

¹We used a slightly modified environment taken from <https://github.com/cagatayyildiz/oder1>

²Note that in the gym inverted pendulum environment dt is set to 1/20

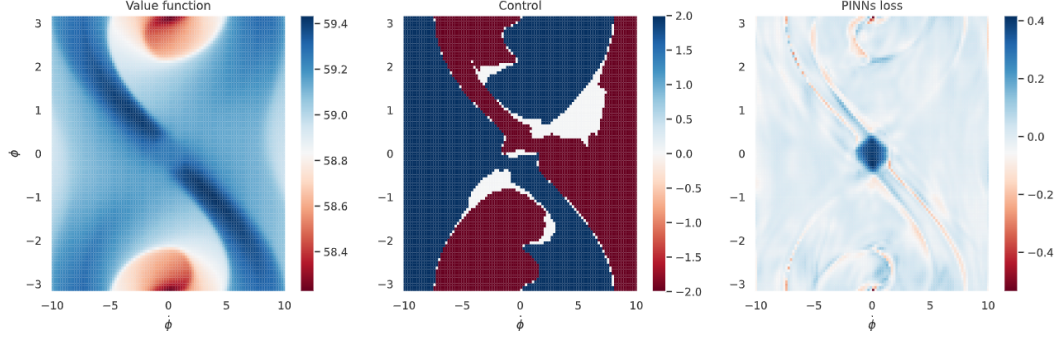


Figure 1: Value function, control and loss maps obtained with PINNs after 1000 iterations.

As expected for the dynamic programming algorithm, the cumulative reward grows when N is increased. However, for $N > 30$, the gain in performance is marginal with respect to the increase of computational cost. Conversely, PINNs is performing worse than VI as it struggles to reach the optimal control in the area where the value function is non-smooth, but it still performs much better than PPO or A2C that are unable to learn the optimal policy function for such a small dt .

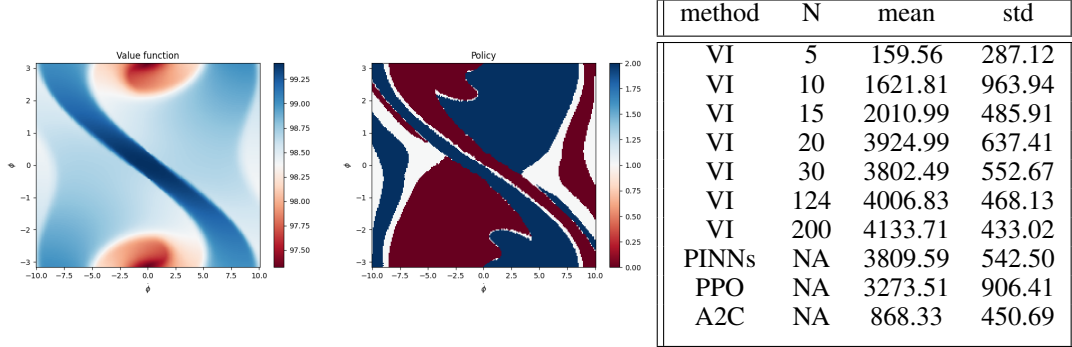


Figure 2: The value function and the optimal policy for the inverted pendulum (left). Mean and standard deviation of the cumulative reward for different methods and several grid sizes for VI (right).

6 Conclusion

In this article, we consider the problem of finding the value function of the optimal policy in the context of continuous time for deterministic dynamical system relying on the HJB equation. There are multiple ways for solving this equation, but one should make sure that the final solution is the viscosity solution to get the value function. We propose methods to compute viscosity solutions, making a particular focus on neural-based approaches. We compare some of those methods on the inverted pendulum. On the one hand, our results have shown that the dynamic programming algorithm is efficient for getting a precise solution once the discretization time step is small enough. However, scaling this algorithm to more difficult problems is limited as the complexity grows exponentially with respect to the dimensionality of the state space. On the other hand, PINNs can be applied to high dimensional problems, but performance depends on the optimization methods that do not always converge to the global minimum. The PINNs algorithm proposed in this paper is designed to find the viscosity solution by approximating it with smooth solutions of other PDE equations (9). Even though it is able to approximate well the solution where it is smooth, it still struggles to match precisely the solution at the points of non-smoothness, which causes a bad control at those points. Possible improvements include considering more sophisticated neural network architectures such as Mixture of Experts or Physics Informed Neural Operators, augmenting the PINNs loss with the data driven loss (e.g. by integrating the discrete time solutions into the observational loss) or considering more advanced sampling techniques.

Acknowledgments Ph. Preux acknowledges the support of the Métropole Européenne de Lille (MEL), ANR, Inria, Université de Lille, through the AI chair Apprenf number R-PILOTE-19-004-APPRENF". All authors would like to thank the Scool research group for providing an outstanding research environment".

References

- [1] Dipak M Adhyaru, IN Kar, and Madan Gopal. Bounded robust control of nonlinear systems using neural network–based hjb solution. *Neural Computing and Applications*, 20:91–103, 2011.
- [2] Bastian Alt, Matthias Schultheis, and Heinz Koepl. Pomdps in continuous time and discrete spaces. *Advances in Neural Information Processing Systems*, 33:13151–13162, 2020.
- [3] Piermarco Cannarsa, Fausto Gozzi, and Halil Mete Soner. A boundary-value problem for hamilton-jacobi equations in hilbert spaces. *Applied Mathematics and Optimization*, 24(1):197–220, 1991.
- [4] Tao Cheng, Frank L Lewis, and Murad Abu-Khalaf. Fixed-final-time-constrained optimal control of nonlinear systems using neural network hjb approach. *IEEE Transactions on Neural Networks*, 18(6):1725–1737, 2007.
- [5] Rémi Coulom. *Reinforcement Learning Using Neural Networks, with Applications to Motor Control*. Theses, Institut National Polytechnique de Grenoble - INPG, June 2002.
- [6] Michael G Crandall and Pierre-Louis Lions. Viscosity solutions of hamilton-jacobi equations. *Transactions of the American mathematical society*, 277(1):1–42, 1983.
- [7] Jérôme Darbon, Peter M Dower, and Tingwei Meng. Neural network architectures using min-plus algebra for solving certain high-dimensional optimal control problems and hamilton–jacobi pdes. *Mathematics of Control, Signals, and Systems*, 35(1):1–44, 2023.
- [8] Kenji Doya. Reinforcement learning in continuous time and space. *Neural computation*, 12(1):219–245, 2000.
- [9] Wendell H Fleming and Halil Mete Soner. *Controlled Markov processes and viscosity solutions*, volume 25. Springer Science & Business Media, 2006.
- [10] Christian Grossmann, Hans-Görg Roos, and Martin Stynes. *Numerical treatment of partial differential equations*, volume 154. Springer, 2007.
- [11] Igor Halperin. Distributional offline continuous-time reinforcement learning with neural physics-informed pdes (sciphy rl for doctrl). *arXiv preprint arXiv:2104.01040*, 2021.
- [12] Jiequn Han, Arnulf Jentzen, and Weinan E. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018.
- [13] Hitoshi Ishii. Uniqueness of unbounded viscosity solution of hamilton-jacobi equations. *Indiana University Mathematics Journal*, 33(5):721–748, 1984.
- [14] Wei Kang and Lucas C. Wilcox. Mitigating the curse of dimensionality: Sparse grid characteristics method for optimal feedback control and hjb equations, 2016.
- [15] George Em Karniadakis, Ioannis G Kevrekidis, Lu Lu, Paris Perdikaris, Sifan Wang, and Liu Yang. Physics-informed machine learning. *Nature Reviews Physics*, 3(6):422–440, 2021.
- [16] Jeongho Kim, Jaek Shin, and Insoon Yang. Hamilton-jacobi deep q-learning for deterministic continuous-time systems with lipschitz continuous controls. *The Journal of Machine Learning Research*, 22(1):9363–9396, 2021.
- [17] Jeongho Kim and Insoon Yang. Hamilton-jacobi-bellman equations for maximum entropy optimal control. *arXiv preprint arXiv:2009.13097*, 2020.

- [18] Derong Liu, Ding Wang, Fei-Yue Wang, Hongliang Li, and Xiong Yang. Neural-network-based online hjb solution for optimal robust guaranteed cost control of continuous-time uncertain nonlinear systems. *IEEE transactions on cybernetics*, 44(12):2834–2847, 2014.
- [19] Lu Lu, Xuhui Meng, Zhiping Mao, and George Em Karniadakis. Deepxde: A deep learning library for solving differential equations. *SIAM Review*, 63(1):208–228, 2021.
- [20] Michael Lutter, Boris Belousov, Kim Listmann, Debora Clever, and Jan Peters. Hjb optimal feedback control with deep differential value functions and action constraints. In *Conference on Robot Learning*, pages 640–650. PMLR, 2020.
- [21] Michael Lutter, Shie Mannor, Jan Peters, Dieter Fox, and Animesh Garg. Robust value iteration for continuous control tasks. *arXiv preprint arXiv:2105.12189*, 2021.
- [22] Michael Lutter, Shie Mannor, Jan Peters, Dieter Fox, and Animesh Garg. Value iteration in continuous actions, states and time. *arXiv preprint arXiv:2105.04682*, 2021.
- [23] Amartya Mukherjee and Jun Liu. Bridging physics-informed neural networks with reinforcement learning: Hamilton-jacobi-bellman proximal policy optimization (hjbppo). *arXiv preprint arXiv:2302.00237*, 2023.
- [24] R. Munos, L.C. Baird, and A.W. Moore. Gradient descent approaches to neural-net-based solutions of the hamilton-jacobi-bellman equation. In *IJCNN'99. International Joint Conference on Neural Networks. Proceedings (Cat. No.99CH36339)*, volume 3, pages 2152–2157 vol.3, 1999.
- [25] Remi Munos. A Study of Reinforcement Learning in the Continuous Case by the Means of Viscosity Solutions. *Machine Learning*, 40:265–299, 2000.
- [26] Tenavi Nakamura-Zimmerer, Qi Gong, and Wei Kang. Adaptive deep learning for high-dimensional hamilton–jacobi–bellman equations. *SIAM Journal on Scientific Computing*, 43(2):A1221–A1247, 2021.
- [27] M. Raissi, P. Perdikaris, and G.E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [28] Halil Mete Soner. Optimal control with state-space constraint i. *SIAM Journal on Control and Optimization*, 24(3):552–561, 1986.
- [29] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- [30] Yuval Tassa and Tom Erez. Least squares solutions of the hjb equation with neural network value-function approximators. *IEEE transactions on neural networks*, 18(4):1031–1041, 2007.
- [31] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE, 2012.
- [32] Haoran Wang, Thaleia Zariphopoulou, and Xun Yu Zhou. Reinforcement learning in continuous time and space: A stochastic control approach. *The Journal of Machine Learning Research*, 21(1):8145–8178, 2020.
- [33] Harley E Wiltzer, David Meger, and Marc G Bellemare. Distributional hamilton-jacobi-bellman equations for continuous-time reinforcement learning. In *International Conference on Machine Learning*, pages 23832–23856. PMLR, 2022.
- [34] Ge Yang, Anurag Ajay, and Pulkit Agrawal. Overcoming the spectral bias of neural value approximation. In *International Conference on Learning Representations*, 2022.
- [35] Çağatay Yıldız, Markus Heinonen, and Harri Lähdesmäki. Continuous-time model-based reinforcement learning, 2021.

A Appendix

A.1 Intuition for Viscosity Solutions

In this section, we aim at providing the intuition behind the viscosity solutions. For that, we draw some parallels between DTRL and CTRL.³

Let us consider the DTRL formulation of the problem. We know that the optimal value function V in DTRL should satisfy the Bellman equation

$$V(x) = \max_u \{r(x, u) + \gamma \sum_{x'} p(x'|x, u) V(x')\}. \quad (15)$$

From that, we can introduce the Bellman operator as

$$T(\psi)(x) = \max_u \{r(x, u) + \gamma \sum_{x'} p(x'|x, u) \psi(x')\} \quad (16)$$

where ψ is an arbitrary function defined on the state space. This operator is known to be monotonic, *i.e.* for any functions ψ, ψ' we have

$$\psi \geq \psi' \Rightarrow T(\psi) \geq T(\psi'). \quad (17)$$

Moreover, from Eq. (15) follows that V should satisfy $V = T(V)$. Therefore, from Eq. (15)-(17) we get the alternative definition for the solution of the Bellman equation 15.

Definition A.1 *Let $V \in C(O)$, then V is the optimal value function if and only if*

- $\forall \psi \in C^1(O)$ such that $\psi \geq V$

$$\max_u \{r(x, u) + \gamma \sum_{x'} p(x'|x, u) \psi(x')\} \geq V(x), \forall x \in O,$$

- $\forall \psi \in C^1(O)$ such that $\psi \leq V$

$$\max_u \{r(x, u) + \gamma \sum_{x'} p(x'|x, u) \psi(x')\} \leq V(x), \forall x \in O.$$

This definition can be seen as the discrete-time version of the viscosity solution definition. Therefore, in the discrete-time case, satisfying the fixed point equation is equivalent to satisfying the "discrete-time" viscosity solution definition.

As mentioned in the paper, $V \in C(O)$ can be non differentiable at some points of O , thus it is impossible to verify whether HJB equation is satisfied everywhere. Therefore, the main idea behind viscosity solutions is to replace V by some smooth functions where V is non differentiable.

In the following, first, we suppose that V is differentiable everywhere and we show a connection between Hamilton-Jacobi-Bellman equation and Bellman equation. Then, for the case when V is non smooth, we replace V by a smooth function and we show that it is possible to derive the notions of viscosity super/subsolutions.

Let us discretize our continuous-time problem with a time-step dt . For simplicity, we consider that for any $x \in O$ there exists an optimal control $u^* = \pi(x) \in U$ so that $V(x) = J(x; u^*)$. Therefore, we replace sup with max in the definition of the value function, though it is possible to show that the next results also hold in case of sup. From the definition of the value function, we get

$$\begin{aligned} V(x(t)) &= \max_u \left\{ \int_t^{t+dt} \gamma^{(s-t)} r(x(s), u(s)) ds + \gamma^{dt} V(x(t+dt)) \right\} \\ &= \max_u \left\{ \int_t^{t+dt} \gamma^{(s-t)} r(x(s), u(s)) ds + e^{dt \ln(\gamma)} V(x(t+dt)) \right\} \\ &\approx \max_u \left\{ r(x(t), u(t)) dt + e^{dt \ln(\gamma)} V(x(t+dt)) \right\} \\ &\approx \max_u \left\{ r(x(t), u(t)) dt + (1 + \ln(\gamma) dt) V(x(t+dt)) \right\} \end{aligned}$$

³This section is based on https://benjaminmoll.com/wp-content/uploads/2020/02/viscosity_for_dummies.pdf.

So we derive this discrete-time dynamic programming problem:

$$V(x_t) = \max_u \{r(x_t, u)dt + (1 + \ln(\gamma)dt)V(x_{t+dt})\}, \quad (18)$$

where $x_{t+dt} = f(x_t, u)dt + x_t$.

Let us suppose that V is differentiable for all $x \in O$ and that $dt \in \left(0, -\frac{1}{\ln(\gamma)}\right)$. By subtracting $(1 + \ln(\gamma)dt)V(x_t)$ from both sides of Eq. (18) and then dividing by dt , we obtain

$$-\ln(\gamma)V(x_t) = \max_u \left\{ r(x_t, u) + \left(\frac{1}{dt} + \ln(\gamma)\right) (V(x_{t+dt}) - V(x_t)) \right\}.$$

If dt goes toward 0, we have

$$\ln(\gamma)V(x_t) = -\max_u \{r(x_t, u) + \nabla_x V(x_t)^T f(x_t, u)\}.$$

This is exactly the Hamilton-Jacobi-Bellman equation, the continuous time equivalent of the Bellman equation.

Now, let us assume that V is non differentiable. As mentioned before, V should be replaced by a smooth function at the points where $\nabla_x V$ does not exist. Let ψ be an arbitrary smooth function on O such that $V - \psi$ has a local maximum at x_t and $V(x_t) = \psi(x_t)$. Therefore, $V \leq \psi$ in a neighborhood of x_t . If $1 + \ln(\gamma)dt > 0$, then

$$\begin{aligned} V(x_t) &= \max_u \{r(x_t, u)dt + (1 + \ln(\gamma)dt)V(x_{t+dt})\} \\ &\leq \max_u \{r(x_t, u)dt + (1 + \ln(\gamma)dt)\psi(x_{t+dt})\} \end{aligned}$$

Let us subtract $(1 + \ln(\gamma)dt)\psi(x_t)$ from both sides and use $\psi(x_t) = V(x_t)$, as a result we have

$$-\ln(\gamma)V(x_t)dt \leq \max_u \{r(x_t, u)dt + (1 + \ln(\gamma)dt)(\psi(x_{t+dt}) - \psi(x_t))\}$$

Then, let us divide by dt and let dt goes toward 0, we have

$$\begin{aligned} -\ln(\gamma)V(x_t) &\leq \max_u \{r(x_t, u) + \nabla_x \psi(x_t)^T f(x_t, u)\} \\ \Leftrightarrow \ln(\gamma)V(x_t) - \max_u \{r(x_t, u) + \nabla_x \psi(x_t)^T f(x_t, u)\} &\leq 0 \\ \Leftrightarrow H(x_t, \psi(x_t), \nabla_x \psi(x_t)) &\leq 0. \end{aligned}$$

This gives us the definition of a viscosity subsolution.

It is possible to obtain the definition of a viscosity supersolution in a like manner, by performing the same derivations for an arbitrary $\psi \in C^1(O)$ such that $V - \psi$ has a local minimum in x_t and $V(x_t) = \psi(x_t)$. In both cases, we use the monotonicity of $\max_u \{r(x_t, u)dt + (1 + \ln(\gamma)dt)\psi(x_{t+dt})\}$ in the function ψ , which is a counterpart of the Bellman operator in Definition A.1.

Thus, we recover the definition of a viscosity solution. The intuition is whenever a solution V of the HJB equation is non differentiable at some point $x \in O$, it should also satisfy other conditions imposed by viscosity for it to be a proper value function. In this way, the viscosity property serves as a regularizer to help to eliminate "bad" solutions of the HJB equation.

A.2 Boundary Conditions. Supplementary

In this section, we cover three different cases of boundary conditions that appear in control problems when $O \subset \mathbb{R}^d$, which are illustrated in Figure 3.

Case 1 If the system exits at any boundary point $x \in \partial O$ once the boundary is reached, *e.g.* it is the case when the exit reward $R(x)$ is sufficiently high to prefer to leave the area, *e.g.* when $R(x) \equiv R \geq r$ for any $x \in \partial O$ and $r \geq r(\bar{x}, u(\bar{x}))$ for any $\bar{x}, u(\bar{x})$. Then, the boundary condition is described with the following equation:

$$V(x) - R(x) = 0 \quad \forall x \in \partial O. \quad (19)$$

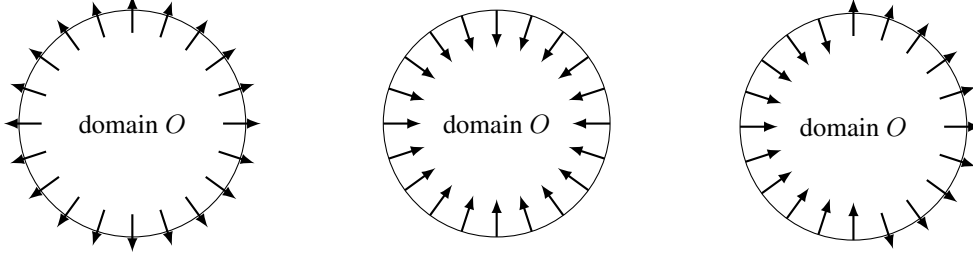


Figure 3: Boundary conditions. Case 1(left), Case 2 (middle) and Case 3(right).

The uniqueness result holds due to the comparison principle that states that under Assumptions 3.1 and provided that W and V are bounded and uniform continuous functions, if W and V are viscosity subsolution and supersolution respectively, then $\sup_{x \in \bar{O}} (W(x) - V(x)) \leq \sup_{x \in \partial O} (W(x) - V(x))$. The existence of such value function is assured with Assumption 3.2-(ii).

Case 2 If the system never exits the control domain. Let us denote the external normal vector at point $x \in \partial O$ as $\eta(x)$, then this boundary can be expressed as $f(x, u^*(x))^T \eta(x) \leq 0$ for any $x \in \partial O$. The optimal policy is not known a priori and thus it is hard to verify this constraint. In [9, 28], it was shown that it can be reformulated as:

$$-H(x, W, \nabla_x W + \alpha \eta(x)) \leq 0 \quad \forall \alpha \leq 0, x \in \partial O. \quad (20)$$

This allows to extend Definition 3.1.

Definition A.2 (Constrained viscosity solution) $W \in C(\bar{O})$ is called a constrained viscosity solution of the HJB equation (5) if it is a viscosity subsolution in O and a viscosity supersolution in \bar{O} , i.e. if $\forall \psi \in C^1(\bar{O})$ and $\forall x \in \bar{O} \cup \arg \min\{(W - \psi)(x) : x \in \bar{O}\}$ with $W(x) = \psi(x)$, we have:

$$H(x, \psi(x), \nabla_x \psi(x)) \geq 0.$$

It is also possible to prove that there exists a continuous value function provided that Assumption 3.2(i) holds and the set of admissible actions is not empty for any state of the system. Under the additional Assumption 3.3, there exists a unique constrained viscosity solution (see [28]).

Case 3 If there exists a subset of points of the boundary at which the system exits the control domain. This is the same boundary condition considered in [25]. This boundary is formulated as follows:

$$R(x) - V(x) \leq 0 \quad \forall x \in \partial O. \quad (21)$$

However, this boundary is not sufficient to have uniqueness, therefore we redefine viscosity for this inequality constraint (21).

Definition A.3 (Viscosity solution with the boundary condition (21)) • $W \in C(\bar{O})$ is a viscosity subsolution of the HJB equation in O with the boundary condition (21) if it is a viscosity subsolution in O and $\forall \psi \in C^1(\bar{O})$ and $\forall x \in \partial O$ local maximum of $W - \psi$ such that $W(x) = \psi(x)$, we have:

$$\min\{H(x, \psi(x), \nabla_x \psi(x)), R(x) - W(x)\} \leq 0$$

• $W \in C(\bar{O})$ is a viscosity supersolution of the HJB equation in O with the boundary condition (21) if it is a viscosity supersolution and $\forall \psi \in C^1(O)$ and $\forall x \in \partial O$ local minimum of $W - \psi$ such that $W(x) = \psi(x)$, we have:

$$\max\{H(x, \psi(x), \nabla_x \psi(x)), R(x) - W(x)\} \geq 0$$

• If $W \in C(\bar{O})$ is a viscosity subsolution and a supersolution with the boundary condition (21) then it is a viscosity solution with the boundary condition (21).

It is easy to check that when Eq. (21) is verified then a viscosity subsolution $W(x)$ in O is a viscosity subsolution with the boundary condition (21). However, when $W(x) > R(x)$ for some point $x \in \partial O$ then definition A.3 imposes an additional constraint that $W(x)$ should be a viscosity supersolution at such boundary points. Then similarly to Case 2, boundary condition (8) should be also satisfied, which can be interpreted as the system not being able to exit at those points. Similarly to Case 2, there is a uniqueness result:

Theorem A.1 *Let us assume that Assumptions 3.1-3.3 hold, then the value function V is in $C(\bar{O})$ and it is the unique viscosity solution of the HJB equation in O with the boundary condition (21).*

The proof of this theorem can be found in [9, 3].

A.3 Possible Approaches. Supplementary

A.3.1 Dynamic Programming

Further, we consider only FEM based dynamic programming proposed in [25]. In the FEM case, we use a triangulation Σ^δ to cover the state space. It is also possible to discretize the control space, denoted by U^δ . The vertices of the triangulation Σ^δ are denoted $\{\xi_1, \xi_2, \dots, \xi_{N_\delta}\}$ with $N_\delta \in \mathbb{N}$. In this setting, V is approximated by a piecewise linear function V^δ . Thus, for $x \in \text{Simplex}(\xi_0, \dots, \xi_d)$, we have

$$V^\delta(x) = \sum_{i=0}^d \lambda_{\xi_i}(x) V^\delta(\xi_i)$$

where $\lambda_{\xi_i}(x)$ is the barycentric coordinates inside the simplex (ξ_0, \dots, ξ_d) .

By using a FEM approximation scheme, the HJB equation is transformed into:

$$V^\delta(\xi) = \sup_{u \in U^\delta} [\gamma^{\tau(\xi, u)} V^\delta(\eta(\xi, u)) + \tau(\xi, u) r(\xi, u)]$$

where $\eta(\xi, u) = \xi + \tau(\xi, u) f(\xi, u)$ and $\tau(\xi, u)$ is a time discretization function that should satisfy:

$$\exists k_1, k_2 > 0, \forall \xi \in \Sigma^\delta, \forall u \in U^\delta, k_1 \delta \leq \tau(\xi, u) \leq k_2 \delta$$

If F^δ is defined as $F^\delta[\phi](\xi) = \sup_{u \in U^\delta} [\gamma^{\tau(\xi, u)} \sum_{i=0}^d \lambda_{\xi_i}(\eta(\xi, u)) \phi(\xi_i) + \tau(\xi, u) r(\xi, u)]$, it is possible to show that F^δ satisfies a contraction property, and since $V^\delta(\xi) = F^\delta[V^\delta](\xi)$ holds, dynamic programming techniques can be applied to compute V^δ . Moreover, it can be proved that $V^\delta \xrightarrow[\delta \rightarrow 0]{} V$ uniformly on any compact of the state space. With this method, one can derive algorithms that converge towards V , without even knowing the dynamics of the system. Thus, this is one of the approaches that allows us to find a viscosity solution of the HJB equation (see [25] for more details).

A.3.2 Introduction to PINNs

Here, we provide a short introduction to PINNs for those readers who are not familiar with this method. The adaptation of PINNs to solving the HJB equation in the viscosity sense is covered in Section 4.2.

To solve a differential equation

$$F(x, W(x), \nabla_x W(x), \nabla_x^2 W(x)) = 0, \quad W : \bar{O} \rightarrow \mathbb{R}, x \in O, \quad (22)$$

with K_1 equality boundary conditions

$$B_i(x, W(x), \nabla_x W(x), \nabla_x^2 W(x)) = 0, \quad x \in \partial O, i \leq K_1, \quad (23)$$

and K_2 inequality boundary conditions

$$G_i(x, W(x), \nabla_x W(x), \nabla_x^2 W(x)) \leq 0, \quad x \in \partial O, i \leq K_2. \quad (24)$$

one can assume that $W(x)$ lies in the class of functions $\mathcal{F}_\theta = \{f_\theta(x) = NN(x, \theta) : \theta \in \Theta\}$ represented by neural networks of a fixed architecture and parametrized with weights $\theta \in \Theta$. If it is the case then there exists θ such that $W(x, \theta)$ should satisfy Eq. (22) and thus minimize the loss

$$\mathcal{L}_{PDE}(\theta) = \frac{1}{N_F} \sum_{i=1}^{N_F} (F(x_i, W(x_i, \theta), \nabla_x W(x_i, \theta), \nabla_x^2 W(x_i, \theta)))^2 \quad \forall x_i \in \mathcal{S}_u(O, N_F), \quad (25)$$

with $\mathcal{S}_u(O, N_F)$ denoting a sample of N_F points drawn uniformly from O . If the solution $W(x, \theta)$ should satisfy some additional boundary constraints then it should also minimize the boundary losses for all $k \leq K_1$ and $k' \leq K_2$

$$\mathcal{L}_{B_k}(\theta) = \frac{1}{N_B^k} \sum_{i=1}^{N_B^k} (B_k(x_i, W(x_i, \theta), \nabla_x W(x_i, \theta), \nabla_x^2 W(x_i, \theta)))^2 \quad \forall x_i \in \mathcal{S}_u(\partial O, N_B^k) \quad (26)$$

$$\mathcal{L}_{G_{k'}}(\theta) = \frac{1}{N_G^{k'}} \sum_{i=1}^{N_G^{k'}} \left([G_{k'}(x_i, W(x_i, \theta), \nabla_x W(x_i, \theta), \nabla_x^2 W(x_i, \theta))]^+ \right)^2 \quad \forall x_i \in \mathcal{S}_u(\partial O, N_G^{k'}), \quad (27)$$

where $[f(x)]^+ = \max\{f(x), 0\}$.

To put everything together, when solving a PDE in a PINNs-like manner, one should train a neural network $W(x, \theta)$ that minimizes:

$$\mathcal{L}(\theta) = \mathcal{L}_{PDE}(\theta) + \sum_{k=1}^{K_1} \lambda_k \mathcal{L}_{B_k}(\theta) + \sum_{k=1}^{K_2} \lambda'_k \mathcal{L}_{G_k}(\theta). \quad (28)$$

where $\lambda_k, \lambda'_k > 0$ are some mixing coefficients for different boundary conditions.

A.4 Experimental Results. Supplementary

A.4.1 Dynamic Programming Experimental Results

In this section, we present the results obtained with one of the algorithms proposed in [25]. First, a grid is built by dividing each axis by N points. Then, we use the Delaunay's triangulation over the grid and apply the Value Iteration algorithm (VI) to the FEM-MDP derived in [25].

We set $\delta = \frac{1}{N}$, δ being the discretization step. The stopping criterion used at the step n is $\|V_n - V_{n-1}\|_\infty \leq \epsilon$ where ϵ is a chosen tolerance. In our experiments we work with $\epsilon = 10^{-5}$.

When δ goes towards 0, our approximated value function, V^δ , converges towards the true value function. In our case, $\delta \rightarrow 0$ is equivalent to $N \rightarrow +\infty$. Empirically, we can see in Figure 4 that this property is satisfied. Indeed, as we increase N , we obtain a more accurate V^δ , and as a result, a better control that leads to a higher cumulative reward.

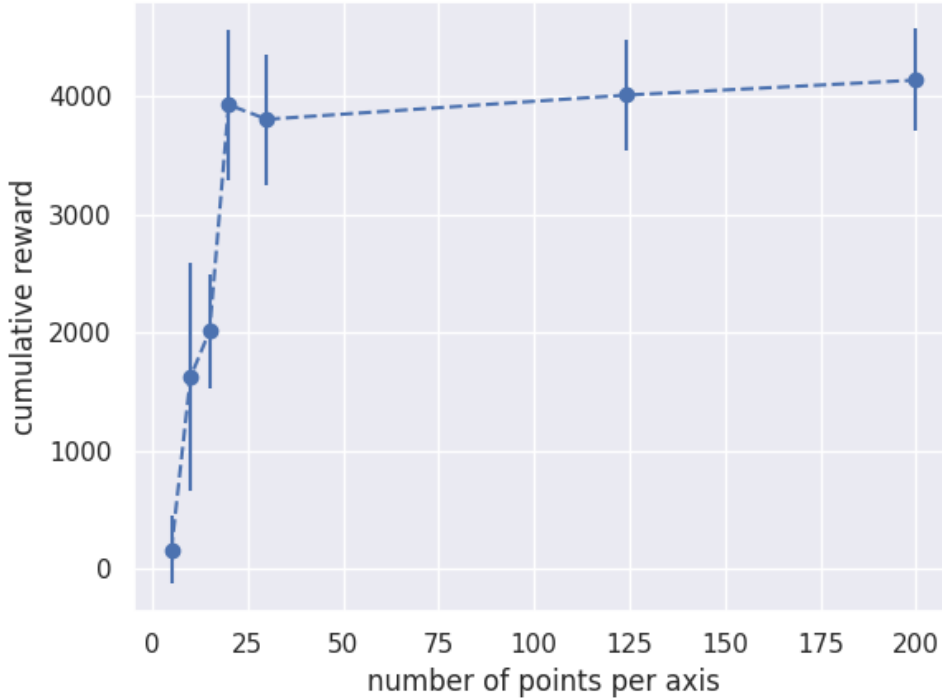


Figure 4: The cumulative reward obtained on the inverted pendulum for different grid sizes N .

Hyperparameters	values
N_D	200000
N_S	100
ν	0.00085
n_ϵ	7
λ	10^{-2}
λ_R	10^{-6}
ϵ_0	10^{-2}
k_ϵ	0.99

Table 1: Hyperparameters for Algorithm 1.



Figure 5: PINNs loss, boundary loss and $\delta(\epsilon, \theta)$ with respect to different λ_R .

A.4.2 Ablation Study on Regularization

We have tested Algorithm 1 on the inverted pendulum environment with multiple hyperparameters. The best performing hyperparameters are gathered in Table 1.

For the inverted pendulum training, we tried several $\epsilon_0 \in (10^{-7}, 10^{-1})$ and $k_\epsilon \in \{0.9, 0.99, 0.999, 0.9999\}$, and we did not observe any significant change in performance for different set of parameters. In what follows, we keep $\epsilon_0 = 10^{-2}$ and $k_\epsilon = 0.99$.

In addition, we study how regularization term λ_R affects the training. We demonstrate the performance of Algorithm 1 with 3 metrics: PDE loss \mathcal{L}_O , boundary loss $\mathcal{L}_{\partial O}$ and ϵ -added term $\delta(\epsilon, \theta)$. All measurements are performed using 5 seeds. From Figure 5, one can conclude that $\lambda_R = 10^{-6}$ achieves the best performance for all metrics, though the values close to it result in the similar behaviour. From the same figure, we see that high λ_R can be too restrictive, resulting into slow convergence. However, putting λ_R too low can lead to more instabilities. Thus, $\lambda_R = 10^{-6}$ is a good trade-off, though it does not remove all the risk of instabilities.