

# On Convergence Analysis of Policy Iteration Algorithms for Entropy-Regularized Stochastic Control Problems

Jin Ma\*, Gaozhan Wang<sup>†</sup> and Jianfeng Zhang<sup>‡</sup>

July 29, 2024

## Abstract

In this paper we investigate the issues regarding the convergence of the *Policy Iteration Algorithm* (PIA) for a class of general continuous-time entropy-regularized stochastic control problems. In particular, instead of employing sophisticated PDE estimates for the iterative PDEs involved in the PIA (see, e.g., Huang-Wang-Zhou [6]), we shall provide a simple proof from scratch for the convergence of the PIA. Our approach builds on probabilistic representation formulae for solutions of PDEs and their derivatives. Moreover, in the infinite horizon model with large discount factor and in the finite horizon model, the similar arguments lead to the exponential rate of convergence of PIA without tear. Finally, with some extra efforts we show that our approach can also be extended to the case when diffusion contains control, in the one dimensional setting but without much extra constraints on the coefficients. We believe that these results are new in the literature.

**Keywords.** Reinforcement learning, model uncertainty, entropy-regularization, policy iteration algorithm, Feynman-Kac formula, Bismut-Elworthy-Li formula, rate of convergence.

*2020 AMS Mathematics subject classification:* 93E35, 60H30, 35Q93.

---

\*Department of Mathematics, University of Southern California, Los Angeles, 90089; email: jinma@usc.edu. This author is supported in part by NSF grant #DMS-2205972.

<sup>†</sup>Department of Mathematics, University of Southern California, Los Angeles, 90089; email: gaozhanw@usc.edu.

<sup>‡</sup>Department of Mathematics, University of Southern California, Los Angeles, 90089; email: jianfenz@usc.edu. This author is supported in part by NSF grant #DMS-2205972.

# 1 Introduction

The *Policy Iteration Algorithm* (PIA), also known as the *Policy Improvement Algorithm*, is a well-known approach in numerical optimal control theory, see, e.g., Jacka-Mijatović [7], Kerimkulov-Siska-Szpruch [8, 9], Puterman [12]. Its main idea is to construct an iteration scheme for the control actions that traces the maximizers/minimizers of the Hamiltonian, so that the corresponding returns are naturally improving. Mathematically this amounts to a type of Picard iteration for the associated HJB equations. Motivated by the above scheme but with model uncertainty, the *Reinforcement Learning* (RL) algorithms for the entropy-regularized stochastic control problems have received very strong attention in recent years. By using the idea of relaxed control, the control problem is regularized (or “penalized”) by Shannon’s entropy in order to capture the trade-off between exploitation (to optimize) and exploration (to learn the model). Such entropy-regularized problem in a continuous time model was introduced by Wang-Zariphopoulou-Zhou [16], see also Guo-Xu-Zariphopoulou [5], Reisinger-Zhang [13], and Tang-Zhang-Zhou [14] in this direction, especially on the relation between the entropy-regularized problem and the original control problem.

The convergence of the PIA for the entropy-regularized problem, in terms of both the value and the optimal strategy, is clearly a central issue in the theory. In a linear quadratic model, Wang-Zhou [17] solved the problem explicitly and the convergence is immediate. Our paper is mainly motivated by the work Huang-Wang-Zhou [6], which established the desired convergence in a general infinite horizon diffusion model with drift controls. Note that the values of the iterative sequence in PIA are by nature increasing (and bounded), so the main issue is to identify its limit with the true value function of the entropy-regularized control problem. By using some sophisticated Sobolev estimates, [6] established uniform regularity for the iterative value functions, especially uniform bounds for their derivatives, and then derived the convergence by compactness arguments. We also refer to Bai-Gamage-Ma-Xie [1] and Dong [3] for some related works.

Our original purpose of this paper is to provide a simple proof from scratch for the convergence results that first appeared in [6], but without using any heavy PDE machinery. Our proof builds on the Bismut-Elworthy-Li representation formulae [2, 4] for derivatives of functions, see also Ma-Zhang [10] and Zhang [18]. These formulae enable us to establish the uniform bounds of the derivatives of the iterative value functions rather easily (see §3.2 Step 1 below), and then the desired convergence under the  $C^2$ -norm follows immediately.

It turns out that, in the case that the discount factor is sufficiently large, our argument

can lead to an exponential rate of convergence of PIA, which is new in the literature<sup>1</sup>. In fact, instead of applying the compactness arguments, in this case the representation formulae can yield the rate of convergence directly. In particular, while the involved derivatives still have uniform bounds, we do not require them for the proof of the convergence here. We would also like to note that, when the discount factor is small, in general it may not be reasonable to expect a good rate of convergence, see Remark 3.5 and Example 3.6 below.

A natural question is then whether the approach works also for the finite horizon case, for which the associated PDEs become parabolic. The answer is affirmative. In particular, we are able to obtain the exponential rate of convergence under the  $C^{1,2}$ -norm, by first considering small time duration and then extending to arbitrary time duration. In this case, we do not need a constraint corresponding to the large discount factor in the infinite horizon case. To the best of our knowledge, this result is also new in the literature.

The widely recognized and much more challenging question is the convergence analysis for the case when the diffusion term contains control. In this paper we shall argue that, in an infinite horizon setting when the state process is scalar, with some extra effort our approach can still lead to the  $C^2$ -convergence of the value function, which implies the convergence of the optimal strategy. In fact, in a further special case we are able to obtain again the exponential rate of convergence. However, we must note that the general higher dimensional case is much more subtle, and our current arguments may be ineffective as they rely heavily on the scalar assumption.

When finalizing the present paper, we learned the very interesting recent paper Tran-Wang-Zhang [15]. In the base case of infinite horizon model with drift control and sufficiently large discount factor, [15] obtained the same exponential rate of convergence. The main ideas are similar, however, they used Schauder estimates from PDE literature while we proved the required estimates from scratch by using the probabilistic representation formulae. It is remarkable that [15] established the convergence for multi-dimensional diffusion controlled models when the control is small in a certain sense and the discount factor is sufficiently large. The key is again the crucial uniform estimates for the associated iterative fully nonlinear PDEs, in the spirit of Evans-Krylov theorem. It will be very interesting to combine our approaches and to explore more general models, especially when there is diffusion control, which we shall leave for future research.

The rest of the paper is organized as follows. In §2 we formulate the problem and prove the main results for the finite horizon case. In §3 we prove the main results in the infinite

---

<sup>1</sup>The same rate is obtained independently by [15], as we will comment in details soon.

horizon case, first in the case when the discounting factor is large, and then for the general case when the discounting factor is small. In §4 we investigate the problem with diffusion control in a scalar setting.

**Notations.** To end this section and to facilitate the reader, we list the following notations that will be used frequently throughout the paper. Let  $E$  be a generic Euclidean space, whose usual inner product is denoted by  $x \cdot y$ , for  $x, y \in E$ , where all  $x \in E$  are column vectors. In particular, for  $A, B \in \mathbb{R}^{d \times d}$ , we denote  $A : B := \text{trace}(AB^\top)$ , where  $B^\top$  is the transpose of  $B$ . Moreover,  $I_d$  denotes the  $d \times d$  identity matrix. For two Euclidean spaces  $E_1, E_2$ , and  $m, k \geq 0$ , we denote  $C_b^{m,k}([0, T] \times E_1; E_2)$  to be the set of functions  $\phi : [0, T] \times E_1 \mapsto E_2$  which is  $m$ -th order continuous differentiable in  $t \in [0, T]$  and  $k$ -th order continuously differentiable in  $x \in E_1$ , such that  $\phi$  as well as all its derivatives involved above are bounded. Moreover,  $C_b^k(E_1; E_2)$  denotes the subspace where  $\phi : E_1 \rightarrow E_2$  is independent of the temporal variable  $t$ . Furthermore, for  $\phi \in C^2(E_1; E_2)$  and  $\psi \in C^{1,2}([0, T] \times E_1; E_2)$ , we denote

$$\begin{aligned} \|\phi\|_0 &:= \sup_{x \in E_1} |\phi(x)|, & \|\phi\|_2 &:= \|\phi\|_0 + \|\phi_x\|_0 + \|\phi_{xx}\|_0; \\ \|\psi\|_0 &:= \sup_{t \in [0, T]} \|\phi(t, \cdot)\|_0, & \|\psi\|_{1,2} &:= \sup_{t \in [0, T]} \left[ \|\psi(t, \cdot)\|_2 + \|\psi_t(t, \cdot)\|_0 \right]. \end{aligned} \quad (1.1)$$

We use both notation  $\partial_x \phi = \phi_x$  for derivatives, whichever is more convenient.

Finally, let  $A \subseteq E$  be a domain. We denote  $\mathcal{P}_0(A)$  to be the set of all probability densities  $\pi$  on  $A$ , namely  $\pi : A \rightarrow \mathbb{R}_+$  such that  $\int_A \pi(a) da = 1$ . For each  $\pi \in \mathcal{P}_0(A)$ , we denote its corresponding Shannon's entropy by

$$\mathcal{H}(\pi) := - \int_A \pi(a) \ln \pi(a) da. \quad (1.2)$$

Moreover, for  $\phi \in \mathbb{L}^1(E_1 \times A; E_2)$  with generic Euclidean spaces  $E_1, E_2$ , we denote:

$$\tilde{\phi}(x, \pi) := \int_A \phi(x, a) \pi(a) da, \quad x \in E_1. \quad (1.3)$$

## 2 The Finite Horizon Case

We begin our discussion by fixing a finite time horizon  $[0, T]$ , as in this case we have complete results with simple arguments and we shall consider the infinite horizon case ( $T = \infty$ ) in the next two sections. Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space on which is defined a standard  $d$ -dimensional Brownian motion  $W$ ,  $\mathbb{F} := \mathbb{F}^W$ , and the control set  $A$  be a bounded domain with smooth boundary in some Euclidean space, which in particular has finite volume:  $0 < |A| < \infty$ . We note that in this section we consider only drift controls.

Given  $(t, x) \in [0, T] \times \mathbb{R}^d$ , our underlying control problem is as follows:

$$\begin{aligned} X_s^{t,x,\alpha} &= x + \int_t^s b(X_l^{t,x,\alpha}, \alpha_l) dl + \int_t^s \sigma(X_l^{t,x,\alpha}) dW_l, \quad s \in [t, T]; \\ u_0(t, x) &:= \sup_{\alpha} \mathbb{E} \left[ g(X_T^{t,x,\alpha}) + \int_t^T r(X_s^{t,x,\alpha}, \alpha_s) ds \right], \end{aligned} \quad (2.1)$$

where  $b : \mathbb{R}^d \times A \rightarrow \mathbb{R}^d$ ,  $\sigma : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ ,  $r : \mathbb{R}^d \times A \rightarrow \mathbb{R}$ ,  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  are measurable functions, and  $\alpha$  is an appropriate  $A$ -valued admissible control. Here for notational simplicity we assume  $b, \sigma, r$  are time homogeneous. All the results in this section will remain true when they depend on  $t$ . It is well known that, under certain technical conditions,  $u_0$  satisfies an HJB equation, and there is a vast literature on numerical methods for  $u_0$ , provided that the coefficients  $b, \sigma, r, g$  are known.

Strongly motivated by numerical methods for the above control problem but with model uncertainty, namely when the coefficients  $b, \sigma, r, g$  are unknown, we consider instead the *entropy-regularized exploratory optimal control problem*. That is, we consider an associated *relaxed control problem* regularized by the *Shannon's entropy* for the purpose of exploration. More precisely, let  $\mathcal{A}_T$  denote the set of functions  $\pi : [0, T] \times \mathbb{R}^d \rightarrow \mathcal{P}_0(A)$ , and recall (1.2), (1.3). Our entropy-regularized exploratory optimal control problem associated to (2.1) takes the form:

$$\begin{aligned} X_s^{t,x,\pi} &= x + \int_t^s \tilde{b}(X_l^{t,x,\pi}, \pi(l, X_l^{t,x,\pi})) dl + \int_t^s \sigma(X_l^{t,x,\pi}) dW_l, \quad s \in [t, T]; \\ J(t, x; \pi) &:= \mathbb{E} \left[ g(X_T^{t,x,\pi}) + \int_t^T [\tilde{r}(X_s^{t,x,\pi}, \pi(s, X_s^{t,x,\pi})) + \lambda \mathcal{H}(\pi(s, X_s^\pi))] ds \right]; \\ u(t, x) &:= \sup_{\pi \in \mathcal{A}_T} J(t, x; \pi). \end{aligned} \quad (2.2)$$

Here  $\lambda > 0$  is the exogenous “temperature” parameter capturing the trade-off between exploitation and exploration. We remark that  $u \rightarrow u_0$  when  $\lambda \downarrow 0$ , see [14].

In the rest of this section we shall assume:

**Assumption 2.1.** (i)  $b, \sigma, r$  are measurable in  $a$  and twice continuously differential in  $x$ ,<sup>2</sup> and both the functions and their derivatives are bounded by a constant  $C_0 > 0$ ; and  $g \in C_b^2(\mathbb{R}^d; \mathbb{R})$  with  $\|g\|_2 \leq C_0$ .

(ii)  $\sigma$  is uniform non-degenerate:  $[\sigma \sigma^\top](x) \geq \frac{1}{C_0} I_d$ ,  $x \in \mathbb{R}^d$ .

Throughout this paper, we shall denote  $C > 0$  to be a generic constant depending only on  $d, \lambda, |A|$ , and  $C_0$ , but not on  $T$ , and it is allowed to vary from line to line. In particular, when the constant does depend on  $T$ , we shall denote it as  $C_T$ .

---

<sup>2</sup>When  $b, \sigma, r$  depend on  $t$ , we require only their continuity in  $t$ .

Clearly, under Assumption 2.1, the SDE for  $X^{t,x,\pi}$  has a unique weak solution and  $u(t, x)$  is well defined and satisfies the following exploratory HJB equation:

$$u_t + \frac{1}{2}\sigma\sigma^\top : u_{xx} + H(x, u_x) = 0, \quad u(T, x) = g(x), \quad (2.3)$$

where  $H(x, z) := \sup_{\pi \in \mathcal{P}_0(A)} [\tilde{b}(x, \pi) \cdot z + \tilde{r}(x, \pi) + \lambda \mathcal{H}(\pi)]$ .

Moreover, a straightforward calculation along the lines of calculus of variation for the Hamiltonian  $H$  shows that the optimal relaxed control  $\pi^*$  takes form  $\pi^*(t, x, a) := \Gamma(x, u_x(t, x), a)$ ,  $(t, x, a) \in [0, T] \times \mathbb{R}^d \times A$ , with  $\Gamma$  being the Gibbs function:

$$\Gamma(x, z, a) := \frac{\gamma(x, z, a)}{\int_A \gamma(x, z, a') da'}, \quad \gamma(x, z, a) := \exp\left(\frac{1}{\lambda}[b(x, a) \cdot z + r(x, a)]\right). \quad (2.4)$$

Further, it is easily seen that the Hamiltonian  $H$  can be written in the following form:

$$H(x, z) = \lambda \ln\left(\int_A \gamma(x, z, a) da\right). \quad (2.5)$$

We then have the following simple result.

**Lemma 2.2.** *Let Assumption 2.1 hold. Then,*

(i)  *$H$  is twice continuously differentiable in  $(x, z)$ , and there exists a constant  $C > 0$ , independent of  $T$ , such that*

$$|H_z| \leq C, \quad 0 \leq H_{zz} \leq CI_d; \quad [ |H| + |H_x| + |H_{xz}| ](x, z) \leq C[1 + |z|]. \quad (2.6)$$

(ii) *The PDE (2.3) has a unique classical solution  $u$  with  $\|u\|_{1,2} \leq Ce^{CT}$ .*

**Proof** (ii) is standard in the PDE literature, given the uniform non-degeneracy of  $\sigma$ , provided (2.6) holds. It remains to check (i). Since  $b, r$  are twice continuously differentiable in  $x$ , by (2.4) it is clear that  $H$  is also twice continuously differentiable in  $(x, z)$ . Moreover, it is easy to check that (suppressing variables)

$$H_z = \frac{\int_A b \gamma da}{\int_A \gamma da}, \quad H_{zz} = \frac{\int_A b b^\top \gamma da \int_A \gamma da - \int_A b \gamma da \int_A b^\top \gamma da}{(\int_A \gamma da)^2};$$

$$H_x = \frac{\int_A [b_x z + r_x] \gamma da}{\int_A \gamma da}, \quad H_{xz} = \frac{\int_A b_x \gamma da}{\int_A \gamma da} - \frac{\int_A [b_x z + r_x] \gamma da \int_A b^\top \gamma da}{(\int_A \gamma da)^2}.$$

Then it is straightforward to verify (2.6). ■

We now introduce the *Policy Iteration Algorithm* (PIA) for solving PDE (2.3):

*Step 0.* Set  $u^0(t, x) := -C_0 - [C_0 - \lambda(\ln |A|)^+](T - t)$ ;<sup>3</sup>

*Step n.* Define  $\pi^n(t, x, a) := \Gamma(x, u_x^{n-1}(t, x), a)$  and  $u^n(t, x) := J(t, x, \pi^n)$ ,  $n \geq 1$ .

Then, using (2.4) and (2.5), one shows that

$$H_z(x, u_x^{n-1}) = \tilde{b}(x, \pi^n(t, x, \cdot)) = \int_A b(x, a) \Gamma(x, u_x^{n-1}, a) da,$$

and that  $u^n$  satisfies the following recursive linear PDE (suppressing variables):

$$\begin{cases} u_t^n + \frac{1}{2} \sigma \sigma^\top : u_{xx}^n + H_z(x, u_x^{n-1}) \cdot (u_x^n - u_x^{n-1}) + H(x, u_x^{n-1}) = 0, \\ u^n(T, x) = g(x). \end{cases} \quad (2.7)$$

The following result is more or less standard (see e.g. [6]), and we omit the proof.

**Proposition 2.3.** *Let Assumption 2.1 hold. Then*

- (i)  $u^n$  is increasing in  $n$  and  $u^n(t, x) \leq C_0 + [C_0 + \lambda(\ln |A|)^+](T - t)$ ;
- (ii) For  $n \geq 1$ ,  $u^n \in C_b^{1,2}([0, T] \times \mathbb{R}^d; \mathbb{R})$  is the unique classical solution of (2.7).

This clearly indicates that  $u^n \uparrow u^*$  for some function  $u^*$ . Our purpose is to argue that  $u^* = u$  and to obtain the rate of convergence. Our main result is as follows.

**Theorem 2.4.** *Under Assumption 2.1,<sup>4</sup> we have*

$$\|\Delta u^n\|_{1,2} \leq \frac{C_T}{2^n}, \quad \text{where } \Delta u^n := u^n - u. \quad (2.8)$$

Consequently, the iterative strategy  $\pi^n(t, x, a) = \Gamma(x, u_x^{n-1}(t, x), a)$  converges to the optimal strategy  $\pi^*(t, x, a) = \Gamma(x, u_x(t, x), a)$  for the entropy-regularized exploratory optimal control problem (2.2).

**Proof** We proceed in five steps.

**Step 1.** In this step we provide probabilistic representation formulae for  $u, u^n$  and their derivatives, which will be crucial for our estimates. Fix  $(t, x) \in [0, T] \times \mathbb{R}^d$  and denote

$$X_s^{t,x} = x + \int_t^s \sigma(X_l^{t,x}) dW_l, \quad s \in [t, T]. \quad (2.9)$$

<sup>3</sup>Note that  $\sup_{\pi \in \mathcal{P}_0(A)} \mathcal{H}(\pi) = (\ln |A|)^+$ . We set  $u^0$  in this way so that  $u^0(t, x) \leq J(t, x, \pi)$  for all  $\pi$ . In particular, while it is not crucial for the remaining analysis, this will imply that  $u^0 \leq u^1$ .

<sup>4</sup>We assume the twice differentiability in Assumption 2.1 in order to get the  $C^{1,2}$ -convergence in the theorem. If we content ourselves with the  $C^{0,1}$ -convergence, from our proofs one can easily see that the second order differentiability is not required. We also note that the  $C^{0,1}$ -convergence is sufficient for the convergence of the optimal strategies.

Let  $u^n$  be the solution to the recursive PDE (2.7), then by standard Feynman-Kac formula we have

$$\begin{aligned} u^n(t, x) &= \mathbb{E} \left[ g(X_T^{t,x}) + \int_t^T f^n(s, X_s^{t,x}) ds \right], \\ u(t, x) &= \mathbb{E} \left[ g(X_T^{t,x}) + \int_t^T f(s, X_s^{t,x}) ds \right], \end{aligned} \quad (2.10)$$

where

$$\begin{aligned} f^n(t, x) &:= [H_z(x, u_x^{n-1}) \cdot (u_x^n - u_x^{n-1})](t, x) + H(x, u_x^{n-1}(t, x)), \\ f(t, x) &:= H(x, u_x(t, x)). \end{aligned} \quad (2.11)$$

Next, for any  $\phi \in C_b^2(\mathbb{R}^d; \mathbb{R})$ , applying the Bismut-Elworthy-Li formula [2, 4] or the representation formula in [10], and following the arguments in [18] we have <sup>5</sup>

$$\begin{aligned} \partial_x \mathbb{E}[\phi(X_s^{t,x})] &= \mathbb{E} \left[ \phi_x(X_s^{t,x}) \nabla X_s^{t,x} \right] = \mathbb{E} \left[ \phi(X_s^{t,x}) N_s^{t,x} \right], \\ \partial_{xx} \mathbb{E}[\phi(X_s^{t,x})] &= \mathbb{E} \left[ \phi_{xx}(X_s^{t,x}) (\nabla X_s^{t,x})^2 + \phi_x(X_s^{t,x}) \nabla^2 X_s^{t,x} \right] \\ &= \mathbb{E} \left[ (\phi_x(X_s^{t,x}))^\top \nabla X_s^{t,x} N_s^{t,x} + \phi(X_s^{t,x}) \nabla N_s^{t,x} \right]. \end{aligned} \quad (2.12)$$

Here in the above, by using the Einstein summation for repeated indices:

$$\begin{aligned} \nabla X_s^{t,x} &= I_d + \int_t^s \sigma_x^i(X_l^{t,x}) \nabla X_l^{t,x} dW_l^i, \\ \nabla_j^2 X_s^{t,x} &= \int_t^s [\sigma_{xxk}^i(X_l^{t,x}) (\nabla X_l^{t,x})^{kj} \nabla X_s^{t,x} + \sigma_x^i(X_l^{t,x}) \nabla_j^2 X_s^{t,x}] dW_l^i, \end{aligned} \quad (2.13)$$

where  $\sigma^i$  is the  $i$ -th column of  $\sigma$ ,  $\nabla_j^2 X = \partial_{x_j} \nabla X^x$ , and the  $i$ -th column of  $\nabla X^{t,x}$  stands for  $\partial_{x_i} X^{t,x}$ . Similarly, denoting  $\check{\sigma} := \sigma^{-1}$  to be the inverse matrix, we have

$$\begin{aligned} N_s^{t,x} &:= \frac{1}{s-t} \int_t^s (\check{\sigma}(X_l^{t,x}) \nabla X_l^{t,x})^\top dW_l; \\ \nabla_i N_s^{t,x} &:= \frac{1}{s-t} \int_t^s ((\nabla X_l^{t,x})^{ij} \check{\sigma}_{x_j}(X_l^{t,x}) \nabla X_l^{t,x} + \check{\sigma}(X_l^{t,x}) \nabla_i^2 X_l^{t,x})^\top dW_l. \end{aligned} \quad (2.14)$$

Furthermore, one can easily check that

$$\begin{aligned} \mathbb{E} [|\nabla X_s^{t,x}|^2 + |\nabla^2 X_s^{t,x}|^2] &\leq C e^{C(s-t)}, \\ \mathbb{E} [|N_s^{t,x}|^2 + |\nabla N_s^{t,x}|^2] &\leq \frac{C}{s-t} e^{C(s-t)}. \end{aligned} \quad (2.15)$$

---

<sup>5</sup>[10] provides only the first one in (2.12). The second one follows the same arguments as in [10, 18] directly by differentiating the first one with respect to the initial value  $x$ . We also note that, the  $N$  in [10, 18] is a row vector, corresponding to the transpose of the  $N$  here.

Then we have the representation formulae for the first order derivatives

$$\begin{aligned} u_x^n(t, x) &= \mathbb{E} \left[ g_x(X_T^{t,x}) \nabla X_T^{t,x} + \int_t^T f^n(s, X_s^{t,x}) N_s^{t,x} ds \right]; \\ u_x(t, x) &= \mathbb{E} \left[ g_x(X_T^{t,x}) \nabla X_T^{t,x} + \int_t^T f(s, X_s^{t,x}) N_s^{t,x} ds \right]; \end{aligned} \quad (2.16)$$

and that for the second order derivatives

$$\begin{aligned} u_{xx}^n(t, x) &= \mathbb{E} \left[ g_{xx}(X_T^{t,x}) (\nabla X_T^{t,x})^2 + g_x(X_T^{t,x}) \nabla^2 X_T^{t,x} \right. \\ &\quad \left. + \int_t^T [N_s^{t,x} (\nabla X_s^{t,x} f_x^n(s, X_s^{t,x}))^\top + f^n(s, X_s^{t,x}) \nabla N_s^{t,x}] ds \right], \\ u_{xx}(t, x) &= \mathbb{E} \left[ g_{xx}(X_T^{t,x}) (\nabla X_T^{t,x})^2 + g_x(X_T^{t,x}) \nabla^2 X_T^{t,x} \right. \\ &\quad \left. + \int_t^T [N_s^{t,x} (\nabla X_s^{t,x} f_x(s, X_s^{t,x}))^\top + f(s, X_s^{t,x}) \nabla N_s^{t,x}] ds \right]. \end{aligned} \quad (2.17)$$

**Remark 2.5.** If  $d = 1$ ,  $\sigma \equiv 1$ , and  $X_t^x := x + W_t$ , then  $\nabla X_t^x \equiv 1$ , and  $N_t^x = \frac{W_t}{t}$ . The first formula in (2.12) is a direct consequence of the integration by parts formula:

$$\partial_x \mathbb{E}[\phi(X_t^x)] = \int_{\mathbb{R}} \phi'(x+y) \frac{1}{\sqrt{2\pi t}} e^{-\frac{y^2}{2t}} dy = \int_{\mathbb{R}} \phi(x+y) \frac{1}{\sqrt{2\pi t}} e^{-\frac{y^2}{2t}} \frac{y}{t} dy = \mathbb{E}[\phi(X_t^x) N_t^x].$$

The general formula follows from the integration by parts formula for Malliavin derivatives (cf. [2, 4, 10]).

**Step 2.** In this step we first assume  $T \leq \delta$ , for some  $\delta > 0$  which will be specified later. We shall estimate  $\varepsilon_1^n := \|\Delta u_x^n\|_0$ , here the subscript  $_1$  stands for the first order derivative  $\partial_x$  (the meaning of  $\varepsilon_2^n$  below is therefore clear). Note that, by (2.11),

$$|\Delta f^n(t, x)| \leq |H_z(x, u_x^{n-1})|(\varepsilon_1^n + \varepsilon_1^{n-1}) + |H(x, u_x) - H(x, u_x^{n-1})|,$$

where  $\Delta f^n := f^n - f$ . Recall (2.6) that  $|H_z| \leq C$ , we have

$$|\Delta f^n(t, x)| \leq C(\varepsilon_1^n + \varepsilon_1^{n-1}). \quad (2.18)$$

Now for any  $n \geq 1$  and  $(t, x) \in [0, T] \times \mathbb{R}^d$ , by (2.16) we have

$$\begin{aligned} |\Delta u_x^n(t, x)| &\leq \mathbb{E} \left[ |\Delta u_x^n(T, X_T^{t,x})| |\nabla X_T^{t,x}| + \int_t^T |\Delta f^n(s, X_s^{t,x})| |N_s^{t,x}| ds \right] \\ &\leq C \|\Delta u_x^n(T, \cdot)\|_0 \mathbb{E}[|\nabla X_T^{t,x}|] + C(\varepsilon_1^n + \varepsilon_1^{n-1}) \int_t^T \mathbb{E}[|N_s^{t,x}|] ds. \end{aligned}$$

We remark that here obviously  $\Delta u_x^n(T, \cdot) \equiv 0$ , however, for the sake of argument later let us keep this term. Then by (2.15), for some constant  $C_1 > 0$  independent of  $T$ ,

$$|\Delta u_x^n(t, x)| \leq C_1 e^{C_1(T-t)} \left[ \|\Delta u_x^n(T, \cdot)\|_0 + (\varepsilon_1^n + \varepsilon_1^{n-1}) \sqrt{T-t} \right].$$

Since  $x$  is arbitrary, we obtain that

$$\varepsilon_1^n \leq C_1 e^{C_1(T-t)} \left[ \|\Delta u_x^n(T, \cdot)\|_0 + (\varepsilon_1^n + \varepsilon_1^{n-1})\sqrt{T-t} \right]. \quad (2.19)$$

We now set  $\delta > 0$  small such that

$$C_1 e^{C_1 \delta} \sqrt{\delta} \leq \frac{1}{3}. \quad (2.20)$$

Then, for  $T \leq \delta$ , (2.19) reads  $\varepsilon_1^n \leq \frac{1}{3}(\varepsilon_1^n + \varepsilon_1^{n-1}) + C\|\Delta u_x^n(T, \cdot)\|_0$ , and thus

$$\varepsilon_1^n \leq \frac{1}{2}\varepsilon_1^{n-1} + C\|\Delta u_x^n(T, \cdot)\|_0. \quad (2.21)$$

Note that  $u_x^0 \equiv 0$ , then it follows from Lemma 2.2 (ii) that  $\varepsilon_1^0 = \|u_x\|_0 \leq C e^{C\delta} \leq C$ . Note further that  $\Delta u_x^n(T, \cdot) = 0$ . Then (2.21) becomes

$$\varepsilon_1^n \leq \frac{\varepsilon_1^0}{2^n} + C\|\Delta u_x^n(T, \cdot)\|_0 \leq \frac{C}{2^n} + C\|\Delta u_x^n(T, \cdot)\|_0 = \frac{C}{2^n}. \quad (2.22)$$

**Step 3.** We next estimate  $\varepsilon_1^n$  for general  $T$ . First, let  $0 = t_0 < \dots < t_m = T$  be such that  $t_i - t_{i-1} \leq \delta$ , where  $\delta$  satisfies (2.20) and is independent of  $T$ . For each  $i$ , apply the arguments in Step 2 on  $[t_{i-1}, t_i]$ , then the first inequality in (2.22) leads to

$$\sup_{t \in [t_{i-1}, t_i]} \|\Delta u_x^n(t, \cdot)\|_0 \leq \frac{C_T}{2^n} + C\|\Delta u_x^n(t_i, \cdot)\|_0.$$

Note that  $\|\Delta u_x^n(t_m, \cdot)\|_0 = 0$ . Then by a backward induction on  $i = m, \dots, 1$ , we obtain immediately that

$$\|\Delta u_x^n\|_0 \leq \frac{C_T}{2^n}. \quad (2.23)$$

Moreover, by (2.10), (2.18), and (2.22), the above leads to

$$|\Delta u^n(t, x)| \leq \mathbb{E} \left[ \int_t^T |\Delta f^n(s, X_s^{t,x})| ds \right] \leq (T-t)\|\Delta f^n\|_0 \leq \frac{C_T}{2^n}.$$

This implies

$$\|\Delta u^n\|_0 \leq \frac{C_T}{2^n}. \quad (2.24)$$

**Step 4.** We now estimate the difference of the second order derivatives. Denote  $\varepsilon_2^n := \|\Delta u_{xx}^n\|_0$ , and we again first assume  $T \leq \delta'$  for some  $\delta' > 0$  to be specified later. By (2.17) we have

$$\begin{aligned} |\Delta u_{xx}^n(t, x)| &\leq \mathbb{E} \left[ \|\Delta u_{xx}^n(T, \cdot)\|_0 |\nabla X_T^{t,x}|^2 + \|\Delta u_x^n(T, \cdot)\|_0 |\nabla^2 X_T^{t,x}| \right. \\ &\quad \left. + \int_t^T [N_s^{t,x} (\nabla X_s^{t,x} \Delta f_x^n(s, X_s^{t,x}))^\top + \Delta f^n(s, X_s^{t,x}) \nabla N_s^{t,x}] ds \right]. \end{aligned} \quad (2.25)$$

Here again we keep the terminal difference term for our argument. Note that

$$\begin{aligned}\Delta f_x^n(t, x) &= [H_{xz}(x, u_x^{n-1}) + H_{zz}(x, u_x^{n-1})u_{xx}^{n-1}][u_x^n - u_x^{n-1}] \\ &\quad + H_z(x, u_x^{n-1})[u_{xx}^n - u_{xx}^{n-1}] + [H_x(x, u_x^{n-1}) - H_x(x, u_x)] \\ &\quad + [H_z(x, u_x^{n-1})u_{xx}^{n-1} - H_z(x, u_x)u_{xx}].\end{aligned}$$

Then it follows from Lemma 2.2 that

$$\begin{aligned}|\Delta f_x^n(t, x)| &\leq C \left\{ [1 + \varepsilon_2^{n-1}][\varepsilon_1^n + \varepsilon_1^{n-1}] + [\varepsilon_2^n + \varepsilon_2^{n-1}] + \varepsilon_1^{n-1} + [\varepsilon_1^{n-1} + \varepsilon_2^{n-1}] \right\} \\ &\leq C [1 + \varepsilon_1^n + \varepsilon_1^{n-1}][\varepsilon_2^n + \varepsilon_2^{n-1}] + C[\varepsilon_1^n + \varepsilon_1^{n-1}].\end{aligned}\tag{2.26}$$

Recall (2.18) and (2.22), we see that (2.25) leads to that

$$\begin{aligned}|\Delta u_{xx}^n(t, x)| &\leq C e^{C\delta'} \left[ \|\Delta u_{xx}^n(T, \cdot)\|_0 + \|\Delta u_x^n(T, \cdot)\|_0 \right] \\ &\quad + C_1 e^{C_1\delta'} \sqrt{\delta'} [1 + \varepsilon_1^n + \varepsilon_1^{n-1}][\varepsilon_2^n + \varepsilon_2^{n-1}] + C[\varepsilon_1^n + \varepsilon_1^{n-1}].\end{aligned}$$

Since  $x$  is arbitrary, we obtain

$$\begin{aligned}\varepsilon_2^n &\leq C e^{C\delta'} \left[ \|\Delta u_{xx}^n(T, \cdot)\|_0 + \|\Delta u_x^n(T, \cdot)\|_0 \right] \\ &\quad + C_1 e^{C_1\delta'} \sqrt{\delta'} [1 + \varepsilon_1^n + \varepsilon_1^{n-1}][\varepsilon_2^n + \varepsilon_2^{n-1}] + C[\varepsilon_1^n + \varepsilon_1^{n-1}].\end{aligned}\tag{2.27}$$

**Step 5.** Finally we estimate  $\varepsilon_2^n$  for arbitrary  $T$ . Let  $0 = t'_0 < \dots < t'_{m'} = T$  be another partition with  $t'_i - t'_{i-1} \leq \delta'$ , where, for the  $C_1$  in (2.27) and  $C_T$  in (2.23),

$$C_1 e^{C_1\delta'} \sqrt{\delta'} \left[ 1 + \frac{C_T}{2^n} + \frac{C_T}{2^{n-1}} \right] \leq \frac{1}{3}.$$

We remark that here we allow  $\delta'$  to depend on  $T$ . For each  $i$ , apply the arguments in Step 4 on  $[t'_{i-1}, t'_i]$ , then by (2.27) and (2.23) we have

$$\begin{aligned}\sup_{t \in [t'_{i-1}, t'_i]} \|\Delta u_{xx}^n(t, \cdot)\|_0 &\leq C e^{C\delta'} \|\Delta u_{xx}^n(t'_i, \cdot)\|_0 \\ &\quad + \frac{1}{3} \left[ \sup_{t \in [t'_{i-1}, t'_i]} \|\Delta u_{xx}^n(t, \cdot)\|_0 + \sup_{t \in [t'_{i-1}, t'_i]} \|\Delta u_{xx}^{n-1}(t, \cdot)\|_0 \right] + \frac{C_T}{2^n}.\end{aligned}$$

By standard arguments, this leads to

$$\sup_{t \in [t'_{i-1}, t'_i]} \|\Delta u_{xx}^n(t, \cdot)\|_0 \leq C e^{C\delta'} \|\Delta u_{xx}^n(t'_i, \cdot)\|_0 + \frac{C_T}{2^n}.$$

Note that  $\|\Delta u_{xx}^n(t'_{m'}, \cdot)\|_0 = 0$ . Then by backward induction on  $i = m', \dots, 1$ , we obtain immediately that

$$\|\Delta u_{xx}^n\|_0 \leq \frac{C_T}{2^n}.\tag{2.28}$$

Finally, by the PDEs (2.3) and (2.7), we obtain from (2.24), (2.23), and (2.28) the desired estimate for  $\|\Delta u_t^n\|_0$ , and thus prove (2.8).  $\blacksquare$

### 3 The Infinite Horizon Case

In this section we consider the infinite horizon case:  $T = \infty$ . We note that in this case the PDE's involved will become purely elliptic, but defined on the whole space. Our argument will be slightly different, albeit along the same line. We shall first still consider only the drift control case, but will extend our result to some special diffusion control cases in the next section.

Consider the following entropy-regularized exploratory optimal control problem:

$$\begin{aligned} X_t^\pi &= x + \int_0^t \tilde{b}(X_s^\pi, \pi(s, X_s^\pi)) ds + \int_0^t \sigma(X_s^\pi) dW_s; \\ J(x, \pi) &:= \mathbb{E} \left[ \int_0^\infty e^{-\rho t} [\tilde{r}(X_t^\pi, \pi(s, X_s^\pi)) + \lambda \mathcal{H}(\pi(t, X_t^\pi))] dt \right], \\ v(x) &:= \sup_{\pi \in \mathcal{A}} J(x, \pi), \end{aligned} \quad (3.1)$$

where  $\rho > 0$  is the discount factor, and  $\mathcal{A}$  denote the set of  $\pi : [0, \infty) \times \mathbb{R}^d \rightarrow \mathcal{P}_0(A)$ . We shall assume that the coefficients  $b, \sigma, r$  satisfy Assumption 2.1 with the constant  $C_0$ . Furthermore, throughout this section, the generic constant  $C > 0$  does not depend on  $\rho$ . In particular, when the constant does depend on  $\rho$ , we shall denote it as  $C_\rho$ .

Clearly, for the same Gibbs form  $\Gamma$  and function  $H$  as in (2.4), (2.5),  $v$  is well defined and satisfies the following HJB equation:

$$\rho v(x) = \frac{1}{2} [\sigma \sigma^\top](x) : v_{xx}(x) + H(x, v_x), \quad x \in \mathbb{R}^d. \quad (3.2)$$

Similar to Lemma 2.2, it is standard to show that under Assumption 2.1, the PDE (3.2) has a unique classical solution  $v \in C_b^2(\mathbb{R}^d; \mathbb{R})$ , with  $\|v\|_2 \leq C_\rho$ .

Consider now the *Policy Iteration Algorithm* (PIA) for solving (3.2) recursively:

*Step 0.* Set  $v^0 := -\frac{1}{\rho} [C_0 - \lambda(\ln |A|)^+]$ ;

*Step n.* For  $n \geq 1$ , define  $\pi^n(x, a) := \Gamma(x, v_x^{n-1}(x), a)$  and  $v^n(x) := J(x, \pi^n)$ .

Then, similar to the analysis in the last section, we can easily check that each  $v^n$  satisfies the following recursive linear PDE:

$$\rho v^n = \frac{1}{2} \sigma \sigma^\top : v_{xx}^n + H_z(x, v_x^{n-1}) \cdot (v_x^n - v_x^{n-1}) + H(x, v_x^{n-1}). \quad (3.3)$$

Then we have the following analogue of Proposition 2.3.

**Proposition 3.1.** *Let Assumption 2.1 hold. Then*

- (i)  $v^n$  is increasing in  $n$  and  $v^n \leq \frac{1}{\rho} [C_0 + \lambda(\ln |A|)^+]$ ;
- (ii) For each  $n \geq 1$ ,  $v^n \in C_b^2(\mathbb{R}^d; \mathbb{R})$  is a classical solution of (3.3).

Our main result of this section is the following analogue of Theorem 2.4,

**Theorem 3.2.** *Let Assumption 2.1 hold and denote  $\Delta v^n := v^n - v$ .*

(i) *There exists a constant  $\rho_0$ , depending only on  $d, \lambda, |A|$ , and  $C_0$ , such that,*

$$\|\Delta v^n\|_2 \leq \frac{C}{2^n}, \quad \text{whenever } \rho \geq \rho_0. \quad (3.4)$$

(ii) *In the general case,  $v^n \rightarrow v$  in  $C^2$ , uniformly on compacts. That is, for any compact set  $K \Subset \mathbb{R}^d$ , it holds that*

$$\lim_{n \rightarrow \infty} \|\Delta v^n\|_{2,K} = 0, \quad (3.5)$$

where  $\|\varphi\|_{0,K} := \sup_{x \in K} |\varphi(x)|$ , and  $\|\varphi\|_{2,K} := \|\varphi\|_{0,K} + \|\varphi_x\|_{0,K} + \|\varphi_{xx}\|_{0,K}$ .

Since the arguments for the proof of Theorem 3.2 will depend crucially on the “size” of the discounting factor  $\rho$ , we shall carry it out separately in the two subsections below, for “large” and “small”  $\rho$ , respectively.

**Remark 3.3.** *In the setting of Theorem 3.2, it is clear that the iterative optimal strategies  $\pi^n(t, x, a) = \Gamma(x, v_x^{n-1}(x), a)$  also converge to the optimal strategy  $\pi^*(t, x, a) = \Gamma(x, v_x(x), a)$  for the entropy-regularized problem (3.1).*

### 3.1 Proof of Theorem 3.2 (i)

In this subsection we prove (3.4), assuming that  $\rho$  is sufficiently large. We emphasize again that in this subsection the generic constant  $C$  does not depend on  $\rho$ . We proceed in three steps.

**Step 1.** We begin by recalling the probabilistic representation formulae for  $v, v^n$  and their derivatives, which are crucial for our arguments. Denote

$$(X^x, \nabla X^x, N^x, \nabla_j^2 X^x, \nabla_i N^x) := (X^{0,x}, \nabla X^{0,x}, N^{0,x}, \nabla_j^2 X^{0,x}, \nabla_i N^{0,x}).$$

Then by standard Feynman-Kac formula we derive from (3.2) and (3.3) that

$$v(x) = \mathbb{E} \left[ \int_0^\infty e^{-\rho t} f(X_t^x) dt \right], \quad v^n(x) = \mathbb{E} \left[ \int_0^\infty e^{-\rho t} f^n(X_t^x) dt \right], \quad (3.6)$$

where  $f(x) := H(x, v_x(x))$ , and

$$f^n(x) := H_z(x, v_x^{n-1}) \cdot (v_x^n - v_x^{n-1}) + H(x, v_x^{n-1}).$$

Next, applying (2.12) on above, we obtain

$$\begin{aligned}
v_x(x) &= \mathbb{E} \left[ \int_0^\infty e^{-\rho t} f(X_t^x) N_t^x dt \right], \quad v_x^n(x) = \mathbb{E} \left[ \int_0^\infty e^{-\rho t} f^n(X_t^x) N_t^x dt \right]; \\
v_{xx}(x) &= \mathbb{E} \left[ \int_0^\infty e^{-\rho t} [N_t^x (\nabla X_t^x f_x(X_t^x))^\top + f(X_t^x) \nabla N_t^x] dt \right]; \\
v_{xx}^n(x) &= \mathbb{E} \left[ \int_0^\infty e^{-\rho t} [N_t^x (\nabla X_t^x f_x^n(X_t^x))^\top + f^n(X_t^x) \nabla N_t^x] dt \right].
\end{aligned} \tag{3.7}$$

**Step 2.** We first estimate  $\varepsilon_1^n := \|\Delta v_x^n\|_0$ . Similarly to (2.18) we have

$$|\Delta f^n(x)| \leq C(\varepsilon_1^n + \varepsilon_1^{n-1}). \tag{3.8}$$

Now for any  $n \geq 1$  and  $x \in \mathbb{R}^d$ , by (3.7) and (2.15) we have

$$\begin{aligned}
|\Delta v_x^n(x)| &\leq \mathbb{E} \left[ \int_0^\infty e^{-\rho t} |\Delta f^n(X_t^x)| |N_t^x| dt \right] \leq C(\varepsilon_1^n + \varepsilon_1^{n-1}) \int_0^\infty e^{-\rho t} \mathbb{E}[|N_t^x|] dt \\
&\leq C(\varepsilon_1^n + \varepsilon_1^{n-1}) \left[ \int_0^\infty \frac{1}{\sqrt{t}} e^{-\rho t + C_1 t} dt \right] = \frac{C_2}{\sqrt{\rho - C_1}} (\varepsilon_1^n + \varepsilon_1^{n-1}),
\end{aligned}$$

where  $C_1, C_2 > 0$  are generic constants independent of  $\rho$  and we assumed  $\rho > C_1$ . Set

$$\rho_0 := C_1 + 9|C_2|^2 \quad \text{so that} \quad \frac{C_2}{\sqrt{\rho_0 - C_1}} = \frac{1}{3}. \tag{3.9}$$

Since  $x$  is arbitrary, then for  $\rho \geq \rho_0$  we obtain

$$\varepsilon_1^n \leq \frac{C_2}{\sqrt{\rho - C_1}} (\varepsilon_1^n + \varepsilon_1^{n-1}) \leq \frac{1}{3} (\varepsilon_1^n + \varepsilon_1^{n-1}), \quad \text{and thus} \quad \varepsilon_1^n \leq \frac{1}{2} \varepsilon_1^{n-1}. \tag{3.10}$$

Moreover, by (3.6) and (2.6) we have

$$|v_x(x)| \leq C(1 + \|v_x\|_0) \int_0^\infty e^{-\rho t} \mathbb{E}[|N_t^x|] dt \leq \frac{C_2}{\sqrt{\rho - C_1}} (1 + \|v_x\|_0) \leq \frac{1}{3} (1 + \|v_x\|_0).$$

This implies  $\|v_x\|_0 \leq \frac{1}{2}$ . Note further that  $v_x^0 \equiv 0$ . Then, by (3.10) we get

$$\varepsilon_1^n \leq \frac{\varepsilon_1^0}{2^n} = \frac{\|v_x\|_0}{2^n} \leq \frac{C}{2^n}, \quad \text{and thus} \quad \|v_x^n\|_0 \leq \|v_x\|_0 + \varepsilon_1^n \leq C. \tag{3.11}$$

Furthermore, it follows from (3.6) that

$$|\Delta v^n(x)| \leq \mathbb{E} \left[ \int_0^\infty e^{-\rho t} |\Delta f^n(X_t^x)| dt \right].$$

Then, by (3.8) we have

$$|\Delta v^n(x)| \leq C(\varepsilon_1^n + \varepsilon_1^{n-1}) \int_0^\infty e^{-\rho t} dt = \frac{C}{\rho} (\varepsilon_1^n + \varepsilon_1^{n-1}).$$

Plug (3.11) into it and note again that  $\rho \geq \rho_0$ , we obtain

$$\|\Delta v^n\|_0 \leq \frac{C}{2^n}. \quad (3.12)$$

**Step 3.** We now estimate  $\varepsilon_2^n := \|\Delta v_{xx}^n\|_0$ . By (3.7) we have, for any  $x \in \mathbb{R}^d$ ,

$$|\Delta v_{xx}^n(x)| \leq \mathbb{E} \left[ \int_0^\infty e^{-\rho t} [|\Delta f_x^n(X_t^x)| |\nabla X_t^x| |N_t^x| + |\Delta f^n(X_t^x)| |\nabla N_t^x|] dt \right].$$

Similarly to (2.26) we have

$$\left| \Delta f_x^n(x) \right| \leq C [1 + \varepsilon_1^n + \varepsilon_1^{n-1}] [\varepsilon_2^n + \varepsilon_2^{n-1}] + C [\varepsilon_1^n + \varepsilon_1^{n-1}].$$

Thus, for possibly larger  $C_1$  and  $C_2$ , by (3.8), (3.11), and (2.15) we have

$$\begin{aligned} |\Delta v_{xx}^n(x)| &\leq C [\varepsilon_2^n + \varepsilon_2^{n-1} + \frac{1}{2^n}] \int_0^\infty e^{-\rho t} \mathbb{E} [|\nabla X_t^x| |N_t^x| + |\nabla N_t^x|] dt \\ &\leq C [\varepsilon_2^n + \varepsilon_2^{n-1} + \frac{1}{2^n}] \int_0^\infty \frac{1}{\sqrt{t}} e^{-\rho t + C_1 t} dt \leq \frac{C_2}{\sqrt{\rho - C_1}} [\varepsilon_2^n + \varepsilon_2^{n-1} + \frac{1}{2^n}]. \end{aligned}$$

Since  $x$  is arbitrary, then for  $\rho \geq \rho_0$  we have

$$\varepsilon_2^n \leq \frac{C_2}{\sqrt{\rho - C_1}} [\varepsilon_2^n + \varepsilon_2^{n-1} + \frac{1}{2^n}] \leq \frac{1}{3} [\varepsilon_2^n + \varepsilon_2^{n-1} + \frac{1}{2^n}].$$

This implies

$$\varepsilon_2^n \leq \frac{1}{2} \varepsilon_2^{n-1} + \frac{C}{2^n}, \quad \text{and thus} \quad \varepsilon_2^n \leq \frac{\varepsilon_2^0}{2^n} + \frac{C}{2^n}. \quad (3.13)$$

Moreover, by (2.6) and noting from Step 2 that  $\|v_x\|_0 \leq C$  for  $\rho \geq \rho_0$ , we have

$$\begin{aligned} |f(x)| &= |H(x, v_x)| \leq C[1 + |v_x|] \leq C, \\ |f_x(x)| &\leq |H_x(x, v_x)| + |H_z(x, v_x)| |v_{xx}| \leq C[1 + \|v_{xx}\|_0]. \end{aligned}$$

Thus by (3.7) and (2.15) we have, again for  $\rho \geq \rho_0$ ,

$$\begin{aligned} |v_{xx}(x)| &\leq C[1 + \|v_{xx}\|_0] \int_0^\infty e^{-\rho t} \mathbb{E} [|\nabla X_t^x| |N_t^x|] + C \int_0^\infty e^{-\rho t} \mathbb{E} [|\nabla N_t^x|] dt \\ &\leq \frac{C_2}{\sqrt{\rho - C_1}} [1 + \|v_{xx}\|_0] \leq \frac{1}{3} [1 + \|v_{xx}\|_0]. \end{aligned}$$

By the arbitrariness of  $x$ , we have  $\|v_{xx}\|_0 \leq \frac{1}{3} [1 + \|v_{xx}\|_0]$  and thus  $\|v_{xx}\|_0 \leq \frac{1}{2}$ . Note further that  $v_{xx}^0 \equiv 0$ . Then  $\varepsilon_2^0 = \|v_{xx}\|_0 \leq \frac{1}{2}$ , and thus it follows from (3.13) that  $\varepsilon_2^n \leq \frac{C}{2^n}$ . This, together with (3.11) and (3.12), proves (3.4).  $\blacksquare$

### 3.2 Proof of Theorem 3.2 (ii)

We now prove (3.5) for arbitrary  $\rho$ . Let  $\rho_1 > \rho$  be a large constant which will be specified later. We remark that here we allow  $\rho_1$  to depend on  $\rho$ . We proceed in two steps.

**Step 1.** In this step we estimate  $\|v^n\|_2$ . Note that we may rewrite (3.3) as

$$\rho_1 v^n = \frac{1}{2} \sigma \sigma^\top : v_{xx}^n + H_z(x, v_x^{n-1}) \cdot (v_x^n - v_x^{n-1}) + H(x, v_x^{n-1}) + (\rho_1 - \rho) v^n.$$

Denote

$$L_1^n := \|v_x^n\|_0, \quad L_2^n := \|v_{xx}^n\|_0, \quad \tilde{f}^n := f^n + (\rho_1 - \rho) v^n.$$

First, similarly to (3.7) we have

$$v_x^n(x) = \mathbb{E} \left[ \int_0^\infty e^{-\rho_1 t} \tilde{f}^n(X_t^x) N_t^x dt \right].$$

By Proposition 3.1 (ii) we have  $\|v^n\|_0 \leq \frac{C}{\rho}$ . Then, by (3.6) and Lemma 2.2,

$$|\tilde{f}^n(x)| \leq C(L_1^n + L_1^{n-1}) + C_\rho(\rho_1 - \rho). \quad (3.14)$$

Thus, by (2.15),

$$\begin{aligned} |v_x^n(x)| &\leq C[L_1^n + L_1^{n-1} + C_\rho(\rho_1 - \rho)] \int_0^\infty e^{-\rho_1 t} \mathbb{E}[|N_t^x|] dt \\ &\leq C[L_1^n + L_1^{n-1} + C_\rho(\rho_1 - \rho)] \int_0^\infty \frac{1}{\sqrt{t}} e^{-\rho_1 t + Ct} dt \\ &\leq \frac{C_2}{\sqrt{\rho_1 - C_1}} (L_1^n + L_1^{n-1}) + \frac{C_\rho(\rho_1 - \rho)}{\sqrt{\rho_1 - C_1}}. \end{aligned}$$

Here  $C_1, C_2$  are generic constants independent of  $\rho$ . Now set  $\rho_1 > \rho$  large enough as in (3.9) such that  $\frac{C_2}{\sqrt{\rho_1 - C_1}} \leq \frac{1}{3}$ . Then, by the arbitrariness of  $x$ , we obtain

$$L_1^n \leq \frac{1}{3}(L_1^n + L_1^{n-1}) + C_\rho.$$

Note further that  $v_x^0 = 0$  and thus  $L_1^0 = 0$ . Then by standard arguments we have

$$L_1^n \leq \frac{L_1^0}{2^n} + C_\rho = C_\rho. \quad (3.15)$$

Next, similarly to (3.7) we have

$$v_{xx}^n(x) = \mathbb{E} \left[ \int_0^\infty e^{-\rho_1 t} [N_t^x (\nabla X_t^x \tilde{f}_x^n(X_t^x))^\top + \tilde{f}^n(X_t^x) \nabla N_t^x] dt \right]. \quad (3.16)$$

By (3.14) and (3.15) it is clear that  $|\tilde{f}^n(x)| \leq C_\rho(\rho_1 - \rho + 1)$ . Moreover, following similar arguments as in (2.26) and by using (3.15) again, we have

$$|\tilde{f}_x^n(x)| \leq |f_x^n(x)| + (\rho_1 - \rho)|v_x^n| \leq C_\rho[L_2^n + L_2^{n-1} + \rho_1 - \rho + 1].$$

Then, by (2.15) and by the arbitrariness of  $x$ ,

$$\begin{aligned} L_2^n &\leq C_\rho[L_2^n + L_2^{n-1} + \rho_1 - \rho + 1] \mathbb{E}\left[\int_0^\infty \frac{1}{\sqrt{t}} e^{-\rho_1 t + Ct} dt\right] \\ &\leq \frac{C_\rho}{\sqrt{\rho_1 - C_1}} [L_2^n + L_2^{n-1} + \rho_1 - \rho + 1] \\ &\leq \frac{1}{3} [L_2^n + L_2^{n-1} + \rho_1 - \rho + 1], \end{aligned}$$

where in the last inequality we set  $\rho_1 := C_1 + 9|C_\rho|^2$  so that  $\frac{C_\rho}{\sqrt{\rho_1 - C_1}} = \frac{1}{3}$ . Note that  $L_2^0 = \|v_{xx}^0\|_0 = 0$ . By standard arguments this implies that

$$L_2^n \leq \frac{1}{2} [L_2^{n-1} + \rho_1 - \rho + 1], \text{ and thus, } L_2^n \leq C(\rho_1 - \rho + 1) \leq C_\rho. \quad (3.17)$$

**Step 2.** We now prove the desired convergence. First, by the monotonicity and boundedness of  $v^n$ , there exists bounded  $v^*$  such that  $v^n \uparrow v^*$ . By (3.15)  $\{v^n\}_{n \geq 1}$  are equicontinuous, then the above convergence is uniform on compacts. Next, by (3.15) and (3.17) we see that  $\{v_x^n\}_{n \geq 1}$  are bounded and equicontinuous. Then by applying Arzella-Ascoli Theorem there exist a subsequence  $\{n_k\}_{k \geq 1}$  such that  $v_x^{n_k}$  converge uniformly on compacts. Note that differentiation is a closed operator, and since  $v^n \rightarrow v^*$ , we must have  $v_x^{n_k} \rightarrow v_x^*$ . This implies that the limit of the subsequence  $\{v_x^{n_k}\}_{k \geq 1}$  is unique, then we must have the convergence of the whole sequence  $v_x^n$ , namely  $v_x^n \rightarrow v_x^*$  uniformly on compacts. In particular, this implies that

$$f^n \rightarrow f^* \text{ uniformly on compacts, where } f^*(x) := H(x, v_x^*).$$

Moreover, by (3.17) it is clear that  $v_x^*$ , whence  $f^*$ , is uniformly Lipschitz continuous.

Note that  $v^n$  is the classical solution of the following PDE:

$$\rho v^n = \frac{1}{2} \sigma \sigma^\top : v_{xx}^n + f^n.$$

Let  $\tilde{v}$  denote the unique viscosity solution of the PDE:

$$\rho \tilde{v} = \frac{1}{2} \sigma \sigma^\top : \partial_{xx} \tilde{v} + f^*. \quad (3.18)$$

By the stability of the viscosity solution, we see that  $v^* = \lim_{n \rightarrow \infty} v^n$  is a viscosity solution of (3.18). Moreover, since  $v_x^*$  is (Lipschitz) continuous, for any smooth test function  $\phi$

of  $v^*$  at  $x$  in the definition of viscosity solution, we must have  $\phi_x(x) = v_x^*(x)$  and thus  $H(x, \phi_x(x)) = f^*(x)$ , then  $v^*$  is also a viscosity solution of the PDE:

$$\rho \tilde{v} = \frac{1}{2} \sigma \sigma^\top : \partial_{xx} \tilde{v} + H(x, \partial_x \tilde{v}).$$

This PDE identifies with (3.2), then by the uniqueness of its viscosity solution, we obtain  $v^* = v$ . That is,  $(v^n, v_x^n) \rightarrow (v, v_x)$  uniformly on compacts.

It remains to prove the desired convergence of  $v_{xx}^n$ .<sup>6</sup> To this end, we shall first introduce another representation formula for  $v_{xx}^n$ . Let us recall (2.13) and denote

$$R_t^x := N_t^x (N_t^x)^\top - \frac{1}{t} \int_0^t D_s N_t^x \check{\sigma}(X_s^x) \nabla X_s ds + \nabla N_t^x, \quad (3.19)$$

where  $D_s N_t^x$  is the Malliavin derivative, see [11], and  $\check{\sigma} := \sigma^{-1}$ . Note that, denoting by  $D_s^i N^x$  (resp.  $\nabla_i X^x$ ) the  $i$ -th column of  $D_s N^x$  (resp.  $\nabla X^x$ ),

$$D_s^i N_t^x = \frac{1}{t} \check{\sigma}(X_s^x) \nabla_i X_s^x + \frac{1}{t} \int_s^t \left( \check{\sigma}_{x_j}(X_l^x) D_s^i X_l^{x,j} \nabla X_s + \check{\sigma}(X_s^x) D_s^i \nabla X_l \right)^\top dW_l,$$

$$D_s^i X_t^x = \sigma^i(X_s^x) + \int_s^t \sigma_{x_j}(X_l^x) D_s^i X_l^{x,j} dW_l,$$

$$D_s^i \nabla X_t^x = \sigma_{x^i}^i(X_s^x) \nabla X_s^x + \int_s^t \left[ \sigma_{x_k}^j(X_l^x) D_s^i X_l^{x,k} \nabla X_l^x + \sigma_x^j(X_l^x) D_s^i \nabla X_l^x \right] dW_l^i.$$

Here we used the Einstein summation again. Fix  $s$ , and consider  $D_s^i X^x, D_s^i \nabla X^x$  as the solution to the above linear SDE systems for  $t \in [s, \infty)$ . One can easily check that

$$\mathbb{E}[|D_s N_t^x|^4] \leq \frac{C}{t^4} e^{Ct}, \quad \text{and thus} \quad \mathbb{E}[|R_t^x|^2] \leq \frac{C}{t^2} e^{Ct}. \quad (3.20)$$

Then, for any  $\phi \in C_b^0(\mathbb{R}^d; \mathbb{R})$ , by [18, Chapter 2] we have<sup>7</sup>

$$\partial_{xx} \mathbb{E}[\phi(X_t^x)] = \mathbb{E}[\phi(X_t^x) R_t^x], \quad (3.21)$$

Thus, for any  $\delta > 0$  small, we may rewrite (3.16) as

$$v_{xx}^n(x) = \mathbb{E} \left[ \int_0^\delta e^{-\rho_1 t} [N_t^x (\nabla X_t^x \tilde{f}_x^n(X_t^x))^\top + \tilde{f}_x^n(X_t^x) \nabla N_t^x] dt + \int_\delta^\infty e^{-\rho_1 t} \tilde{f}_x^n(X_t^x) R_t^x dt \right].$$

<sup>6</sup>This can be done by PDE arguments. In particular, once we have a uniform Hölder continuity of  $v_{xx}^n$ , then it follows from the same compactness argument at above to derive the convergence of  $v_{xx}^n$ . Nevertheless, we provide a probabilistic proof for the convergence directly here.

<sup>7</sup>As in Remark 2.5, when  $d = 1$  and  $\sigma \equiv 1$ , we have  $\nabla X_t^x = 1$ ,  $N_t^x = \frac{W_t}{t}$ ,  $D_s N_t^x = \frac{1}{t}$ ,  $\nabla N_t^x = 0$ , then  $R_t^x = \frac{W_t^2 - t}{t^2}$ , and thus

$$\phi_{xx}(X_t^x) = \partial_{xx} \int_{\mathbb{R}} \phi(y) \frac{1}{\sqrt{2\pi t}} e^{-\frac{(y-x)^2}{2t}} dy = \int_{\mathbb{R}} \phi(y) \frac{1}{\sqrt{2\pi t}} e^{-\frac{(y-x)^2}{2t}} \frac{(y-x)^2 - t}{t^2} dy = \mathbb{E}[\phi(X_t^x) R_t^x].$$

We remark that  $\mathbb{E}[|\tilde{f}^n(X_t^x)R_t^x|] \leq \frac{C}{t}e^{Ct}$  which is not integrable around  $t = 0$ , so at above we use different representations for small  $t$  and large  $t$ . Similarly we have

$$v_{xx}(x) = \mathbb{E}\left[\int_0^\delta e^{-\rho_1 t} [N_t^x (\nabla X_t^x \tilde{f}_x(X_t^x))^\top + \tilde{f}(X_t^x) \nabla N_t^x] dt + \int_\delta^\infty e^{-\rho_1 t} \tilde{f}(X_t^x) R_t^x dt\right];$$

where  $\tilde{f}(x) = H(x, v_x) + (\rho_1 - \rho)v$ .

Now fix a compact set  $K \Subset \mathbb{R}^d$ , and assume that  $K \subseteq B_{M_0}(0)$  for some  $M_0 > 0$ . Denote, for all  $M > M_0$ ,

$$\varepsilon_M^n := \sup_{|x| \leq M} [|\Delta v^n(x)| + |\Delta v_x^n(x)|] \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

Then one can easily see that, for any  $x \in K$ ,

$$\begin{aligned} & |\Delta v_{xx}^n(x)| \\ & \leq \mathbb{E}\left[\int_0^\delta e^{-\rho_1 t} [(|\tilde{f}^n| + |\tilde{f}_x|)(X_t^x) |\nabla X_t^x| |N_t^x| + (|\tilde{f}^n| + |\tilde{f}|)(X_t^x) |\nabla N_t^x|] dt \right. \\ & \quad \left. + \int_\delta^\infty e^{-\rho_1 t} |\Delta \tilde{f}^n(X_t^x)| |R_t^x| dt\right] \\ & \leq C_\rho \int_0^\delta e^{-\rho_1 t} \mathbb{E}[|\nabla X_t^x| |N_t^x| + |\nabla N_t^x|] dt + \int_\delta^\infty e^{-\rho_1 t} \mathbb{E}[(\varepsilon_M^n + C_\rho \mathbf{1}_{\{|X_t^x| \geq M\}}) |R_t^x|] dt \\ & \leq C_\rho \sqrt{\delta} + C \varepsilon_M^n \int_\delta^\infty \frac{1}{t} e^{-\rho_1 t + Ct} dt + \frac{C_\rho}{M} \mathbb{E}\left[\int_\delta^\infty e^{-\rho_1 t} |X_t^x| |R_t^x| dt\right] \\ & \leq C_\rho \sqrt{\delta} + \varepsilon_M^n C_\rho \ln \frac{1}{\delta} + \frac{C_\rho}{M} (1 + |x|) \mathbb{E}\left[\int_\delta^\infty \frac{1}{t} e^{-\rho_1 t + Ct} dt\right] \\ & \leq C_\rho \sqrt{\delta} + \varepsilon_M^n C_\rho \ln \frac{1}{\delta} + \frac{C_\rho}{M} (1 + |M_0|) \ln \frac{1}{\delta}. \end{aligned}$$

Thus

$$\sup_{x \in K} |\Delta v_{xx}^n(x)| \leq C_\rho \sqrt{\delta} + \varepsilon_M^n C_\rho \ln \frac{1}{\delta} + \frac{C_\rho}{M} (1 + |M_0|) \ln \frac{1}{\delta}.$$

Fix  $M, \delta$  and send  $n \rightarrow \infty$ , we obtain

$$\overline{\lim}_{n \rightarrow \infty} \sup_{x \in K} |\Delta v_{xx}^n(x)| \leq C_\rho \sqrt{\delta} + \frac{C_\rho}{M} (1 + |M_0|) \ln \frac{1}{\delta}.$$

By first sending  $M \rightarrow \infty$  and then  $\delta \rightarrow 0$ , we obtain the desired estimate:

$$\overline{\lim}_{n \rightarrow \infty} \sup_{x \in K} |\Delta v_{xx}^n(x)| = 0. \quad \blacksquare$$

**Remark 3.4.** When  $d = 1$ , the uniform estimate for  $\|v_{xx}^n\|_0$  in Step 1 and the convergence of  $v_{xx}^n$  in Step 2 become trivial. Indeed, in this case we have

$$v_{xx}^n = \frac{2}{\sigma^2(x)} \left[ \rho v^n - H_z(x, v_x^{n-1})(v_x^n - v_x^{n-1}) - H(x, v_x^{n-1}) \right].$$

Then the desired boundedness and convergence of  $v_{xx}^n$  follow directly from those of  $v^n, v_x^n$ . We shall use this feature to study a diffusion control model in the next section.

**Remark 3.5.** The convergence in this case relies heavily on the fact that  $v^n$  is monotone and hence converging. When  $\rho$  is small, in general Picard iteration may not converge, not to mention rate of convergence, as we see in the following example.

**Example 3.6.** Let  $d = 1$ . Consider the following (linear) PDE with unique bounded classical solution  $v \equiv 0$ :

$$\rho v = \frac{1}{2}v_{xx} + v_x.$$

Set  $v^0(x) := -\cos x$ , and define  $v^n$  recursively by Picard iteration:

$$\rho v^n = \frac{1}{2}v_{xx}^n + v_x^{n-1}.$$

Then

$$v_x^n(0) = \begin{cases} 0, & n \text{ is even;} \\ \frac{(-1)^{\frac{n-1}{2}}}{(\rho + \frac{1}{2})^n}, & n \text{ is odd.} \end{cases} \quad (3.22)$$

In particular,  $|v_x^{2m+1}(0)| \rightarrow \infty$  as  $m \rightarrow \infty$ , whenever  $\rho < \frac{1}{2}$ .

**Proof** By (3.7), we have

$$\begin{aligned} v_x^n(x) &= \mathbb{E} \left[ \int_0^\infty e^{-\rho t_1} v_x^{n-1}(x + W_{t_1}) \frac{W_{t_1}}{t_1} dt_1 \right], \\ v_x^{n-1}(x + W_{t_1}) &= \mathbb{E} \left[ \int_0^\infty e^{-\rho t_2} v_x^{n-2}(x + W_{t_1+t_2}) \frac{W_{t_1+t_2} - W_{t_1}}{t_2} dt_2 \middle| \mathcal{F}_{W_{t_1}} \right]. \end{aligned}$$

Plug the second formula into the first one, we get

$$v_x^n(x) = \mathbb{E} \left[ \int_{\mathbb{R}_+^2} e^{-\rho(t_1+t_2)} v_x^{n-2}(x + W_{t_1+t_2}) \frac{W_{t_1}}{t_1} \frac{W_{t_1+t_2} - W_{t_1}}{t_2} dt_2 dt_1 \right].$$

Repeat the arguments and note that  $v_x^0(x) = \sin x$ , we obtain

$$v_x^n(x) = \mathbb{E} \left[ \int_{\mathbb{R}_+^n} e^{-\rho T_n} \sin(x + W_{T_n}) \frac{W_{t_1}}{t_1} \cdots \frac{W_{T_n} - W_{T_{n-1}}}{t_n} dt_n \cdots dt_1 \right],$$

where  $T_i := t_1 + \cdots + t_i$ . Note that  $E[\cos W_t \frac{W_t}{t}] = 0$ , and

$$\mathbb{E}[\sin W_t \frac{W_t}{t}] = \sum_{k=0}^{\infty} (-1)^k \mathbb{E} \left[ \frac{W_t^{2k+1}}{(2k+1)!} \frac{W_t}{t} \right] = \sum_{k=0}^{\infty} (-t)^k \frac{(2k+1)!!}{(2k+1)!} = \sum_{k=0}^{\infty} \frac{(-t)^k}{2^k k!} = e^{-\frac{t}{2}}.$$

Let  $\text{Im}$  denote the imaginary part of a complex number. Then

$$\begin{aligned}
v_x^n(0) &= \text{Im} \mathbb{E} \left[ \int_{\mathbb{R}_+^n} e^{-\rho T_n} e^{\sqrt{-1}W_{T_n}} \frac{W_{t_1}}{t_1} \dots \frac{W_{T_n} - W_{T_{n-1}}}{t_n} dt_n \dots dt_1 \right] \\
&= \text{Im} \int_{\mathbb{R}_+^n} \prod_{i=1}^n e^{-\rho t_i} \mathbb{E} \left[ e^{\sqrt{-1}(W_{T_i} - W_{T_{i-1}})} \frac{W_{T_i} - W_{T_{i-1}}}{t_i} \right] dt_n \dots dt_1 \\
&= \text{Im} \int_{\mathbb{R}_+^n} \prod_{i=1}^n e^{-\rho t_i} \mathbb{E} \left[ e^{\sqrt{-1}W_{t_i}} \frac{W_{t_i}}{t_i} \right] dt_n \dots dt_1 = \text{Im} \left( \int_0^\infty e^{-\rho t} \mathbb{E} \left[ e^{\sqrt{-1}W_t} \frac{W_t}{t} \right] dt \right)^n \\
&= \text{Im} \left( \int_0^\infty e^{-\rho t} \sqrt{-1} e^{-\frac{t}{2}} dt \right)^n = \text{Im} \left( \frac{\sqrt{-1}}{\rho + \frac{1}{2}} \right)^n.
\end{aligned}$$

This implies (3.22) immediately. ■

## 4 The Scalar Case with Diffusion Control

In this section we consider the diffusion control case, where the corresponding HJB equation becomes fully nonlinear. It has been widely recognized that the general case is much more challenging than the drift control case, which we shall leave to future research. In this section we consider only the *one-dimensional case*, i.e,  $d = 1$ . Recall Remark 3.4.

Consider the setting in §3, our entropy-regularized problem is:

$$\begin{aligned}
X_t^\pi &= x + \int_0^t \tilde{b}(X_s^\pi, \pi(s, X_s^\pi)) ds + \int_0^t \sqrt{\widetilde{\sigma}^2(X_s^\pi, \pi(s, X_s^\pi))} dW_s; \\
J(x, \pi) &:= \mathbb{E} \left[ \int_0^\infty e^{-\rho t} [\tilde{r}(X_t^\pi, \pi(s, X_s^\pi)) + \lambda \mathcal{H}(\pi(t, X_t^\pi))] dt \right], \\
v(x) &:= \sup_{\pi \in \mathcal{A}_{Lip}} J(x, \pi).
\end{aligned} \tag{4.1}$$

Here in the above, to simplify the arguments we restrict the admissible relaxed controls to  $\mathcal{A}_{Lip}$ , the set of  $\pi \in \mathcal{A}$  that is Lipschitz continuous in  $x$ , so that the  $X^\pi$  above has a unique strong solution. Furthermore, In this section we shall assume:

**Assumption 4.1.**  $d = 1$ ;  $b, \sigma, r$  satisfy Assumption 2.1, with  $\sigma$  depending on  $a$ ; and  $b, \sigma$  are uniformly continuous in  $a$ , uniformly in  $x$ .

It is well-known that, in this case  $v$  satisfies the fully-nonlinear HJB equation:

$$\begin{aligned}
\rho v &= H(x, v_x, v_{xx}), \quad x \in \mathbb{R}, \quad \text{where} \\
H(x, z, q) &:= \sup_{\pi \in \mathcal{P}_0(A)} \left[ \frac{1}{2} \widetilde{\sigma}^2(x, \pi) q + \tilde{b}(x, \pi) z + \tilde{r}(x, \pi) + \lambda \mathcal{H}(\pi) \right].
\end{aligned} \tag{4.2}$$

Moreover, the maximizer of the Hamiltonian  $H$  has the Gibbs form  $\Gamma$ :

$$\Gamma(x, z, q, a) := \frac{\gamma(x, z, q, a)}{\int_A \gamma(x, z, q, a') da'}, \quad (4.3)$$

where  $\gamma(x, z, q, a) := \exp\left(\frac{1}{\lambda}\left[\frac{1}{2}\sigma^2(x, a)q + b(x, a)z + r(x, a)\right]\right)$ ,

and consequently, we have

$$H(x, z, q) = \lambda \ln \left( \int_A \gamma(x, z, q, a) da \right). \quad (4.4)$$

We have the following simple extension of Lemma 2.2:

**Lemma 4.2.** *Let Assumption 4.1 hold. Then  $H$  is twice continuously differentiable in  $(x, z, q)$ ; jointly convex in  $(z, q)$ ; and, for some constant  $C > 0$ ,*

$$|H_z|, |H_q|, |H_{zz}|, |H_{zq}|, |H_{qq}| \leq C, \quad H_q \geq \frac{1}{C}, \quad |H(x, z, q)| \leq C[1 + |z| + |q|]. \quad (4.5)$$

From Lemma 4.2 we see that  $H$  is convex, strictly increasing, and also has linear growth in  $q$ . The following observation about the asymptotic behavior of  $H$  in  $q$  is crucial for our convergence analysis on the recursive PDEs.

**Proposition 4.3.** *Assume that Assumption 4.1 is in force. Then, for any  $\varepsilon > 0$ , there exists  $C_\varepsilon > 0$  such that*

$$|h(x, z, q)| \leq \varepsilon|q| + C|z| + C_\varepsilon, \quad \text{where } h := H - H_z z - H_q q. \quad (4.6)$$

*In particular, if  $b, \sigma$  are uniformly Hölder continuous in  $a$ , uniformly in  $x$ , then*

$$|h(x, z, q)| \leq C[1 + |z| + \ln(1 + |q|)]. \quad (4.7)$$

**Proof** We prove the result only for  $q > 0$ . The case  $q < 0$  can be proved similarly<sup>8</sup>.

First, since  $H$  is jointly convex in  $(z, q)$ , we have

$$h(x, z, q) \leq H(x, 0, 0) \leq C. \quad (4.8)$$

To see the opposite inequality, denote  $\theta(x, z, q, a) := \frac{1}{2}\sigma^2(x, a) + b(x, a)\frac{z}{q}$ . When there is no confusion, we omit the variables  $(x, z, q)$  in  $\theta, \gamma$ . Then, by (4.3), (4.4) we have

$$-h(x, z, q) = \frac{q \int_A \theta(a) \gamma(a) da}{\int_A \gamma(a) da} - \lambda \ln \left( \int_A \gamma(a) da \right). \quad (4.9)$$

---

<sup>8</sup>Actually this case is not needed, because later on we can easily show that  $v_{xx}^n \geq -C$  for all  $n$ .

Denote  $C_1 := \sup_{a \in A, x \in \mathbb{R}} \frac{2|b(x,a)|}{|\sigma^2(x,a)|} < \infty$ . When  $q \leq 2C_1|z|$ , the result is obviously true. We now assume  $q \geq 2C_1|z|$ , which implies  $\theta(a) \geq \frac{1}{C_2} > 0$ , for some  $C_2 > 0$ , for all  $(x, a)$ . Fix  $(x, q, z)$  and denote

$$\bar{\theta} := \sup_{a \in A} \theta(a) > \frac{1}{C_2}, \quad A_\varepsilon := \{a \in A : \theta(a) \geq \bar{\theta} - \varepsilon\}.$$

Let  $a^\varepsilon \in A$  be such that  $\theta(a^\varepsilon) \geq \bar{\theta} - \frac{\varepsilon}{2}$ . Since  $b, \sigma$  are uniformly continuous in  $a$ , uniformly in  $x$ , and  $\sigma$  is bounded and  $q \geq 2C_1|z|$ , we see that  $\theta$  is also uniformly continuous in  $a$ . Thus there exists  $\delta_\varepsilon > 0$ , independent of  $(x, z, q)$ , such that  $\theta(a) \geq \bar{\theta} - \varepsilon$  whenever  $|a - a^\varepsilon| \leq \delta_\varepsilon$ . That is,  $A_\varepsilon \supset A \cap B_{\delta_\varepsilon}(a^\varepsilon)$ . Then, since  $A$  has smooth boundary, there exists  $\mu_\varepsilon > 0$ , depending only on the model parameters, such that  $|A_\varepsilon| \geq \mu_\varepsilon$  for all  $(x, z, q)$  with  $q \geq 2C_1|z|$ . Note that

$$\begin{aligned} \int_A \gamma(a) da &\geq \int_{A_\varepsilon} \gamma(a) da \geq e^{\frac{(\bar{\theta}-\varepsilon)q}{\lambda}-C} \mu_\varepsilon; \\ \frac{\int_{A \setminus A_\varepsilon} \gamma(a) da}{\int_{A_\varepsilon} \gamma(a) da} &\leq \frac{\int_{A \setminus A_\varepsilon} \gamma(a) da}{\int_{A_\varepsilon} \gamma(a) da} \leq \frac{e^{\frac{(\bar{\theta}-\varepsilon)q}{\lambda}+C} |A \setminus A_\varepsilon|}{e^{\frac{(\bar{\theta}-\varepsilon)q}{\lambda}-C} |A_\varepsilon|} \leq e^{-\frac{\varepsilon}{2\lambda}q+C} \frac{|A|}{\mu_\varepsilon}; \\ \int_A \theta(a) \gamma(a) da &\leq \bar{\theta} \int_{A_\varepsilon} \gamma(a) da + (\bar{\theta} - \varepsilon) \int_{A \setminus A_\varepsilon} \gamma(a) da \\ &\leq \left[ \bar{\theta} + (\bar{\theta} - \varepsilon) e^{-\frac{\varepsilon}{2\lambda}q+C} \frac{|A|}{\mu_\varepsilon} \right] \int_{A_\varepsilon} \gamma(a) da. \end{aligned}$$

Then, by (4.9) and assuming without loss of generality that  $\varepsilon < \frac{1}{C_2}$  so that  $\bar{\theta} - \varepsilon > 0$ ,

$$\begin{aligned} -h(x, z, q) &\leq q \left[ \bar{\theta} + (\bar{\theta} - \varepsilon) e^{-\frac{\varepsilon}{2\lambda}q+C} \frac{|A|}{\mu_\varepsilon} \right] - \lambda \ln \left( e^{\frac{(\bar{\theta}-\varepsilon)q}{\lambda}-C} \mu_\varepsilon \right) \\ &= q \left[ \bar{\theta} + (\bar{\theta} - \varepsilon) e^{-\frac{\varepsilon}{2\lambda}q+C} \frac{|A|}{\mu_\varepsilon} \right] - [(\bar{\theta} - \varepsilon)q - \lambda C + \lambda \ln \mu_\varepsilon] \leq \varepsilon q + C_\varepsilon. \end{aligned} \tag{4.10}$$

This, together with (4.8), proves (4.6).

Finally, if  $b, \sigma$  are Hölder- $\beta$  continuous in  $a$ , then we can easily see that  $\mu_\varepsilon \geq \frac{1}{C} \varepsilon^{\frac{1}{\beta}}$ . Thus from the second line of (4.10) we have

$$\begin{aligned} -h(x, z, q) &\leq q \left[ \bar{\theta} + (\bar{\theta} - \varepsilon) e^{-\frac{\varepsilon}{2\lambda}q+C} C \varepsilon^{-\frac{1}{\beta}} \right] - [(\bar{\theta} - \varepsilon)q - C + \lambda \ln \varepsilon^{\frac{1}{\beta}}] \\ &\leq \varepsilon q + C q e^{-\frac{\varepsilon}{2\lambda}q} \varepsilon^{-\frac{1}{\beta}} - C \ln \varepsilon + C. \end{aligned}$$

Set  $\varepsilon := 2\lambda(1 + \frac{1}{\beta}) \frac{\ln q}{q}$ . Then, assuming  $q \geq e$  without loss of generality,

$$\begin{aligned} -h(x, z, q) &\leq C \ln q + C q e^{-(1+\frac{1}{\beta}) \ln q} \left( \frac{q}{\ln q} \right)^{\frac{1}{\beta}} + C \ln q - C \ln \ln q + C \\ &= C \ln q + C (\ln q)^{-\frac{1}{\beta}} + C \ln q - C \ln \ln q + C \leq C \ln q + C. \end{aligned}$$

This, together with the arguments in (4.10) and (4.8), leads to (4.7).  $\blacksquare$

For the PIA corresponding to (4.2), we set  $v^0$  the same as in §3, and for  $n \geq 1$ , define  $\pi^n(x, a) := \Gamma(x, v_{xx}^{n-1}(x), a)$  and  $v^n(x) := J(x, \pi^n)$ . Then, using (4.3) and (4.4), it is easy to check that  $v^n$  satisfies the following recursive linear PDE:

$$\begin{aligned} \rho v^n &= H_q(x, v_x^{n-1}, v_{xx}^{n-1})(v_{xx}^n - v_{xx}^{n-1}) \\ &\quad + H_z(x, v_x^{n-1}, v_{xx}^{n-1})(v_x^n - v_x^{n-1}) + H(x, v_x^{n-1}, v_{xx}^{n-1}) \\ &= H_q(x, v_x^{n-1}, v_{xx}^{n-1})v_{xx}^n + H_z(x, v_x^{n-1}, v_{xx}^{n-1})v_x^n + h(x, v_x^{n-1}, v_{xx}^{n-1}). \end{aligned} \quad (4.11)$$

The following result is similar to Proposition 3.1.

**Proposition 4.4.** *Let Assumption 4.1 (i) hold. Then*

- (i) *For each  $n \geq 1$ ,  $v^n \in C_b^2(\mathbb{R}^d; \mathbb{R})$  is a classical solution of (4.11);*
- (ii)  *$v^n$  is increasing in  $n$  and  $v^n \leq \frac{1}{\rho} [C_0 + \lambda(\ln |A|)^+]$ .*

Our main result is as follows.

**Theorem 4.5.** *Let Assumption 4.1 hold. Then  $v^n \rightarrow v$  in  $C^2$ , uniformly on compacts in the sense of (3.5). Consequently,  $\pi^n \rightarrow \Gamma$  as well.*

**Proof** Again, we proceed in several steps. Denote

$$L_1^n := \|v_x^n\|_0, \quad L_2^n := \|v_{xx}^n\|_0, \quad \bar{L}_1^n := \sum_{k=1}^n \frac{L_1^k}{3^{n-k+1}}.$$

**Step 1.** First, since  $d = 1$ , by (4.11) we may write down  $v_{xx}^n$  explicitly:

$$v_{xx}^n = \frac{1}{H_q(x, v_x^{n-1}, v_{xx}^{n-1})} \left[ \rho v^n - H_z(x, v_x^{n-1}, v_{xx}^{n-1})v_x^n - h(x, v_x^{n-1}, v_{xx}^{n-1}) \right]. \quad (4.12)$$

Then, for  $\varepsilon > 0$ , by (4.5) and by the arbitrariness of  $x$  we have

$$L_2^n \leq \frac{\varepsilon}{C_1} L_2^{n-1} + C(L_1^n + L_1^{n-1}) + C_{\rho, \varepsilon} \leq \frac{1}{3} L_2^{n-1} + C(L_1^n + L_1^{n-1}) + C_\rho,$$

where we set  $\varepsilon := \frac{C_1}{3}$  in the second inequality. Then by standard arguments we have

$$L_2^n \leq C\bar{L}_1^n + C_\rho. \quad (4.13)$$

**Step 2.** Let  $\rho_1 > 0$  be a large constant. Rewrite (4.11) as:

$$\rho_1 v^n = \frac{1}{2} v_{xx}^n + f_n(x) + (\rho_1 - \rho) v^n, \quad \text{where} \quad (4.14)$$

$$f_n(x) := H_q(x, v_x^{n-1}, v_{xx}^{n-1})v_{xx}^n + H_z(x, v_x^{n-1}, v_{xx}^{n-1})v_x^n + h(x, v_x^{n-1}, v_{xx}^{n-1}) - \frac{1}{2} v_{xx}^n.$$

Then, denoting  $X_t^x := x + W_t$ , by Remark 2.5 we have

$$v_x^n(x) := \mathbb{E} \left[ \int_0^\infty e^{-\rho_1 t} [f_n(X_t^x) + (\rho_1 - \rho)v^n(X_t^x)] \frac{W_t}{t} dt \right]. \quad (4.15)$$

By (4.13) we get

$$\begin{aligned} |v_x^n(x)| &\leq C \left[ L_2^n + L_2^{n-1} + L_1^n + L_1^{n-1} + 1 + C_\rho |\rho_1 - \rho| \right] \int_0^\infty e^{-\rho_1 t} \frac{1}{\sqrt{t}} dt \\ &\leq \frac{C}{\sqrt{\rho_1}} \left[ \bar{L}_1^n + \bar{L}_1^{n-1} + L_1^n + L_1^{n-1} + C_\rho + C_\rho |\rho_1 - \rho| \right] \\ &\leq \frac{C_1}{\sqrt{\rho_1}} \left[ L_1^n + \bar{L}_1^{n-1} \right] + \frac{C_\rho}{\sqrt{\rho_1}} [|\rho_1 - \rho| + 1]. \end{aligned}$$

Setting  $\rho_1 = 16C_1^2$  and by the arbitrariness of  $x$ , we get

$$L_1^n \leq \frac{1}{4}(L_1^n + \bar{L}_1^{n-1}) + C_\rho, \quad \text{and thus} \quad L_1^n \leq \frac{1}{3}\bar{L}_1^{n-1} + C_\rho$$

Note that

$$\bar{L}_1^n = \frac{1}{3}L_1^n + \frac{1}{3}\bar{L}_1^{n-1} \leq \frac{1}{9}\bar{L}_1^{n-1} + \frac{1}{3}\bar{L}_1^{n-1} + C_\rho \leq \frac{1}{2}\bar{L}_1^{n-1} + C_\rho.$$

This, together with (4.13), implies immediately that

$$\bar{L}_1^n \leq C_\rho, \quad \text{and thus} \quad L_1^n \leq C_\rho, \quad L_2^n \leq C_\rho. \quad (4.16)$$

**Step 3.** Follow the arguments in the beginning of Step 2 in §3.2, we have  $(v^n, v_x^n) \rightarrow (v^*, v_x^*)$  uniformly on compacts for some function  $v^* \in C_b^1(\mathbb{R}; \mathbb{R})$  such that  $v_x^*$  is Lipschitz continuous. Fix an arbitrary  $x_0$ , and denote  $q_* := \underline{\lim}_{n \rightarrow \infty} v_{xx}^n(x_0) = \lim_{k \rightarrow \infty} v_{xx}^{n_k}(x_0)$  for some subsequence  $\{n_k\}_{k \geq 1}$ , which may depend on  $x_0$ . By (4.11) we have, at  $x_0$ ,

$$\begin{aligned} v_{xx}^n &= v_{xx}^{n-1} + \frac{1}{H_q(x_0, v_x^{n-1}, v_{xx}^{n-1})} \times \\ &\quad \left[ \rho v^n - H_z(x_0, v_x^{n-1}, v_{xx}^{n-1})(v_x^n - v_x^{n-1}) - H(x_0, v_x^{n-1}, v_{xx}^{n-1}) \right] \\ &\leq v_{xx}^{n-1} + \frac{\rho v^* - H(x_0, v_x^*, v_{xx}^{n-1})}{H_q(x_0, v_x^*, v_{xx}^{n-1})} + C\varepsilon_n, \end{aligned} \quad (4.17)$$

where

$$\varepsilon_n' := |v^n(x_0) - v^*(x_0)| + |v_x^n(x_0) - v_x^*(x_0)|, \quad \varepsilon_n := \varepsilon_n' + \varepsilon_{n-1}' \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

Since  $H$  is convex in  $(z, q)$ , by (4.11) again we have

$$\begin{aligned} \rho v^n &= H(x_0, v_x^n, v_{xx}^n) - \left[ H(x_0, v_x^n, v_{xx}^n) - H(x_0, v_x^{n-1}, v_{xx}^{n-1}) \right. \\ &\quad \left. - H_q(x_0, v_x^{n-1}, v_{xx}^{n-1})(v_{xx}^n - v_{xx}^{n-1}) \right. \\ &\quad \left. - H_z(x_0, v_x^{n-1}, v_{xx}^{n-1})(v_x^n - v_x^{n-1}) \right] \leq H(x_0, v_x^n, v_{xx}^n). \end{aligned} \quad (4.18)$$

Set  $n = n_k + 1$  and send  $k \rightarrow \infty$ , we have  $\rho v^* \leq H(x_0, v_x^*, q_*)$ . Then, by (4.17),

$$v_{xx}^n \leq v_{xx}^{n-1} + \frac{H(x_0, v_x^*, q_*) - H(x_0, v_x^*, v_{xx}^{n-1})}{H_q(x_0, v_x^*, v_{xx}^{n-1})} + C\varepsilon_n.$$

Denote  $\tilde{\varepsilon}_n := q_* - \inf_{m \geq n} v_{xx}^{m-1}(x_0) \geq 0$ . Then  $\lim_{n \rightarrow \infty} \tilde{\varepsilon}_n = 0$  and  $q_* \leq v_{xx}^{n-1}(x_0) + \tilde{\varepsilon}_n$ , thus

$$v_{xx}^n \leq v_{xx}^{n-1} - \frac{1}{C_1}(v_{xx}^{n-1} - q_*) + C(\varepsilon_n + \tilde{\varepsilon}_n).$$

This implies that

$$v_{xx}^n - q_* \leq (1 - \frac{1}{C_1})(v_{xx}^{n-1} - q_*) + C(\varepsilon_n + \tilde{\varepsilon}_n).$$

Then by standard arguments we have  $\overline{\lim}_{n \rightarrow \infty} (v_{xx}^n(x_0) - q_*) \leq 0$ . This, together with the definition of  $q_*$ , implies the limit  $\lim_{n \rightarrow \infty} v_{xx}^n(x_0) = q_*$  exists. Then it follows from the closeness of the differentiation operator that  $\lim_{n \rightarrow \infty} v_{xx}^n(x) = v_{xx}^*(x)$ . Now send  $n \rightarrow \infty$  in (4.11), we see that  $v^*$  satisfies (4.2). Note further that  $q \mapsto H(x, z, q)$  has an inverse function, then from (4.2) we conclude that  $v_{xx}^*$  is uniformly Lipschitz continuous. Finally, it follows from the uniqueness of classical solutions to (4.2) that  $v^* = v$ .  $\blacksquare$

#### 4.1 Rate of convergence in a further special case

In this subsection we obtain the rate of convergence under the following extra assumption.

**Assumption 4.6.** For  $\phi = b, \sigma, r$ , there exist  $\underline{\phi}, \overline{\phi} \in C_b^0(\mathbb{R}; \mathbb{R})$  such that

$$\lim_{x \rightarrow -\infty} \sup_{a \in A} |\phi(x, a) - \underline{\phi}(x)| = \lim_{x \rightarrow \infty} \sup_{a \in A} |\phi(x, a) - \overline{\phi}(x)| = 0.$$

**Theorem 4.7.** Let Assumption 4.1 and 4.6 hold true. Then

$$\lim_{n \rightarrow \infty} \|\Delta v^n\|_2 = 0. \tag{4.19}$$

Moreover, there exists  $\rho_0 > 0$  and  $C > 0$  such that, whenever  $\rho \geq \rho_0$ ,

$$\|\Delta v^n\|_2 \leq \frac{C}{2^n}. \tag{4.20}$$

**Proof** We proceed in three steps. Denote

$$\varepsilon_0^n := \|\Delta v^n\|_0, \quad \varepsilon_1^n := \|\Delta v_x^n\|_0, \quad \varepsilon_2^n := \|\Delta v_{xx}^n\|_0.$$

**Step 1.** Denote  $\bar{v}(x) := \frac{1}{\rho}[\bar{r}(x) + \lambda(\ln |A|)^+]$  and, for  $R > 0$ ,

$$\delta_R := \sup_{x \geq R, a \in A} [|b(x, a) - \bar{b}(x)| + |\sigma(x, a) - \bar{\sigma}(x)| + |r(x, a) - \bar{r}(x)|] \rightarrow 0,$$

as  $R \rightarrow \infty$ . For any  $x \geq 2R$  and any  $\pi \in \mathcal{A}$ , by (4.1) we have

$$\begin{aligned}
|J(x, \pi) - \bar{v}(x)| &\leq \sup_{\pi \in \mathcal{A}} \int_0^\infty e^{-\rho t} \mathbb{E}[|\tilde{r}(X_t^\pi, \pi(s, X_s^\pi)) - \bar{r}(x)|] dt \\
&\leq \frac{\delta_R}{\rho} + C \sup_{\pi \in \mathcal{A}} \int_0^\infty e^{-\rho t} \mathbb{P}(X_t^\pi \leq R) dt \leq \frac{\delta_R}{\rho} + C \sup_{\pi \in \mathcal{A}} \int_0^\infty e^{-\rho t} \mathbb{P}(|X_t^\pi - x| \geq R) dt \\
&\leq \frac{\delta_R}{\rho} + \frac{C}{R^2} \sup_{\pi \in \mathcal{A}} \int_0^\infty e^{-\rho t} \mathbb{E}[|X_t^\pi - x|^2] dt \leq \frac{\delta_R}{\rho} + \frac{C}{R^2} \int_0^\infty e^{-\rho t + C_1 t} dt \\
&\leq \frac{\delta_R}{\rho} + \frac{C}{R^2(\rho - C_1)},
\end{aligned}$$

for some  $C_1 > 0$ , and here we assume  $\rho \geq \rho_0 > C_1$ . This clearly implies that

$$\lim_{x \rightarrow \infty} |v(x) - \bar{v}(x)| = 0. \quad (4.21)$$

Next, by (4.3), (4.4), and (4.9) one can easily show that

$$\begin{aligned}
\sup_{x \geq R} \left[ |H_z(x, z, q) - \bar{b}(x)| + |H_q(x, z, q) - \frac{1}{2}\bar{\sigma}^2(x)| \right. \\
\left. + |h(x, z, q) + \bar{v}(x)| \right] \leq C_{z,q} \delta_R,
\end{aligned} \quad (4.22)$$

where  $C_{z,q}$  depends on the bound of  $z, q$ . By (4.11) we have

$$v^n(x) = \mathbb{E} \left[ \int_0^\infty e^{-\rho t} h(X_t^n, v_x^{n-1}(X_t^n), v_{xx}^{n-1}(X_t^n)) dt \right].$$

From the Step 1 in the proof of Theorem 4.5 we can easily see that, for  $\rho$  large,  $v_x^{n-1}$  and  $v_{xx}^{n-1}$  are uniformly bounded, uniformly in  $n$  and  $\rho$ . Then, similar to (4.21) we can show that  $\lim_{x \rightarrow \infty} \sup_n |v^n(x) - \bar{v}| = 0$ , which in turn shows that  $\lim_{x \rightarrow \infty} \sup_n |\Delta v^n(x)| = 0$ . Similarly we have  $\lim_{x \rightarrow -\infty} \sup_n |\Delta v^n(x)| = 0$ . These, together with Theorem 4.5, lead easily to that

$$\lim_{n \rightarrow \infty} \varepsilon_0^n = 0. \quad (4.23)$$

**Step 2.** Let  $\rho_1 > \rho$  be a large number. Recall (4.15), similarly we have

$$v_x(x) := \mathbb{E} \left[ \int_0^\infty e^{-\rho_1 t} [f + (\rho_1 - \rho)v](X_t^x) \frac{W_t}{t} dt \right], \quad f(x) := H(x, v_x, v_{xx}) - \frac{1}{2}v_{xx}.$$

Then,

$$|\Delta v_x^n(x)| \leq \int_0^\infty e^{-\rho_1 t} [C + (\rho_1 - \rho)\varepsilon_0^n] \frac{\mathbb{E}[|W_t|]}{t} dt \leq \frac{C}{\sqrt{\rho_1}} + \frac{\rho_1 - \rho}{\sqrt{\rho_1}} \varepsilon_0^n.$$

By the arbitrariness of  $x$ , we have

$$\overline{\lim}_{n \rightarrow \infty} \varepsilon_1^n \leq \frac{C}{\sqrt{\rho_1}} + \frac{\rho_1 - \rho}{\sqrt{\rho_1}} \lim_{n \rightarrow \infty} \varepsilon_0^n = \frac{C}{\sqrt{\rho_1}}.$$

Since  $\rho_1$  is arbitrary, we obtain

$$\lim_{n \rightarrow \infty} \varepsilon_1^n = 0. \quad (4.24)$$

Moreover, recall (4.12) and similarly we have

$$v_{xx}(x) = \frac{1}{H_q(x, v_x, v_{xx})} \left[ \rho v - H_z(x, v_x, v_{xx}) v_x - h(x, v_x, v_{xx}) \right].$$

By (4.22) and (4.23), (4.24), for  $x \geq R$  we have

$$|\Delta v_{xx}^n(x)| \leq \left| \frac{2}{\bar{\sigma}^2} [\rho v^{n-1} - \bar{b} v_x^{n-1} + \bar{v}] - \frac{2}{\bar{\sigma}^2} [\rho v - \bar{b} v_x + \bar{v}] \right| (x) + C \delta_R \leq C \delta_R.$$

That is,  $\lim_{x \rightarrow \infty} \sup_n |\Delta v_{xx}^n(x)| = 0$ . Similarly  $\lim_{x \rightarrow -\infty} \sup_n |\Delta v_{xx}^n(x)| = 0$ . Thus, it follows from Theorem 4.5 that

$$\lim_{n \rightarrow \infty} \varepsilon_2^n = 0. \quad (4.25)$$

Combining (4.25), (4.25), (4.25), we obtain (4.19).

**Step 3.** We now derive the rate of convergence when  $\rho$  is large. Denote

$$\sigma_0(x) := \sqrt{2H_q(x, v_x, v_{xx})}, \quad X_t^x = x + \int_0^t \sigma_0(X_s^x) dW_s.$$

Note that

$$\rho \Delta v^n = \frac{1}{2} \sigma_0^2(x) \Delta v_{xx}^n + F^n(x), \quad (4.26)$$

where, by (4.11),

$$\begin{aligned} F^n(x) &:= H_q(x, v_x^{n-1}, v_{xx}^{n-1})(v_x^n - v_x^{n-1}) + H_z(x, v_x^{n-1}, v_{xx}^{n-1})(v_x^n - v_x^{n-1}) \\ &\quad + H(x, v_x^{n-1}, v_{xx}^{n-1}) - H(x, v_x, v_{xx}) - H_q(x, v_x, v_{xx}) \Delta v_{xx}^n. \end{aligned}$$

Since  $H$  is jointly convex in  $(z, q)$ , we have

$$\begin{aligned} 0 &\leq H(x, v_x, v_{xx}) - H(x, v_x^{n-1}, v_{xx}^{n-1}) + H_q(x, v_x^{n-1}, v_{xx}^{n-1}) \Delta v_{xx}^{n-1} \\ &\quad + H_z(x, v_x^{n-1}, v_{xx}^{n-1}) \Delta v_x^{n-1} \leq C [|\Delta v_{xx}^{n-1}|^2 + |\Delta v_x^{n-1}|^2]. \end{aligned}$$

Then

$$\begin{aligned} |F^n(x)| &\leq \left| H_q(x, v_x^{n-1}, v_{xx}^{n-1}) \Delta v_{xx}^n + H_z(x, v_x^{n-1}, v_{xx}^{n-1}) \Delta v_x^n - H_q(x, v_x, v_{xx}) \Delta v_{xx}^n \right| \\ &\quad + C [|\Delta v_{xx}^{n-1}|^2 + |\Delta v_x^{n-1}|^2] \\ &\leq C \left[ (|\Delta v_x^{n-1}| + |\Delta v_{xx}^{n-1}|) |\Delta v_{xx}^n| + |\Delta v_x^n| + |\Delta v_{xx}^{n-1}|^2 + |\Delta v_x^{n-1}|^2 \right] \\ &\leq C \left[ \varepsilon_2^{n-1} [\varepsilon_2^n + \varepsilon_2^{n-1}] + \varepsilon_1^n + \varepsilon_1^{n-1} \right]. \end{aligned} \quad (4.27)$$

Note that  $v_{xx} = H^{-1}(x, v_x, \rho v)$ , where  $H^{-1}$  is the inverse function with respect to  $q$ . It is clear that  $v \in C_b^3(\mathbb{R})$  and hence  $\sigma_0 \in C_b^1(\mathbb{R})$ . Then, for the  $N_t^x$  corresponding to  $\sigma_0$ , we have

$$\begin{aligned} |\Delta v^n(x)| &\leq \mathbb{E} \left[ \int_0^\infty e^{-\rho t} |F^n(X_t^x)| dt \right] \leq \frac{C}{\rho} \left[ \varepsilon_2^{n-1} [\varepsilon_2^n + \varepsilon_2^{n-1}] + \varepsilon_1^n + \varepsilon_1^{n-1} \right]; \\ |\Delta v_x^n(x)| &\leq \mathbb{E} \left[ \int_0^\infty e^{-\rho t} |F^n(X_t^x) N_t^x| dt \right] \leq \frac{C}{\sqrt{\rho - C_1}} \left[ \varepsilon_2^{n-1} [\varepsilon_2^n + \varepsilon_2^{n-1}] + \varepsilon_1^n + \varepsilon_1^{n-1} \right]; \end{aligned}$$

for some appropriate  $C_1$  and for  $\rho \geq \rho_0 > 2C_1$ . Then we can easily get

$$\rho \varepsilon_0^n + \sqrt{\rho} \varepsilon_1^n \leq C \left[ \varepsilon_2^{n-1} [\varepsilon_2^n + \varepsilon_2^{n-1}] + \varepsilon_1^n + \varepsilon_1^{n-1} \right]. \quad (4.28)$$

Moreover, by (4.26) and (4.27) we have

$$\varepsilon_2^n \leq C \rho \varepsilon_0^n + C \left[ \varepsilon_2^{n-1} [\varepsilon_2^n + \varepsilon_2^{n-1}] + \varepsilon_1^n + \varepsilon_1^{n-1} \right] \leq C \left[ \varepsilon_2^{n-1} [\varepsilon_2^n + \varepsilon_2^{n-1}] + \varepsilon_1^n + \varepsilon_1^{n-1} \right].$$

Combined with (4.28), this leads to

$$\sqrt{\rho} \varepsilon_1^n + \varepsilon_2^n \leq C_2 \left[ \varepsilon_2^{n-1} [\varepsilon_2^n + \varepsilon_2^{n-1}] + \varepsilon_1^n + \varepsilon_1^{n-1} \right]. \quad (4.29)$$

By (4.25), there exists  $n_0$  such that  $\varepsilon_2^n \leq \frac{1}{3C_2}$  for all  $n \geq n_0$ . Moreover, assume  $\rho \geq \rho_0 \geq 9C_2^2$ . Then, for  $n > n_0$ ,

$$\sqrt{\rho} \varepsilon_1^n + \varepsilon_2^n \leq \frac{1}{3} [\varepsilon_2^n + \varepsilon_2^{n-1}] + \frac{\sqrt{\rho}}{3} [\varepsilon_1^n + \varepsilon_1^{n-1}].$$

This implies that

$$\sqrt{\rho} \varepsilon_1^n + \varepsilon_2^n \leq \frac{1}{2} [\varepsilon_2^{n-1} + \sqrt{\rho} \varepsilon_1^{n-1}], \quad n \geq n_0.$$

Then it follows from standard arguments that  $\sqrt{\rho} \varepsilon_1^n + \varepsilon_2^n \leq \frac{C}{2^n}$ . Plug this into (4.28), we obtain further that  $\rho \varepsilon_0^n \leq \frac{C}{2^n}$ .  $\blacksquare$

**Remark 4.8.** *Assumption 4.6 is used to prove (4.19), but (4.20) relies only on (4.19), not on Assumption 4.6 directly, as we saw in Step 3 of the proof. In other words, any possible alternative sufficient conditions for (4.19) will imply (4.20) as well when  $\rho$  is large.*

## References

- [1] Bai, L., Gamage, T., Ma, J., and Xie, P., (2023) *Reinforcement Learning for Optimal Dividend Problem under Diffusion Model*, Preprint, arXiv:2309.10242.

- [2] Bismut, J. M., (1984) *Large Deviation and Malliavin Calculus*, Progress in Mathematics **45**. Birkhä user.
- [3] Dong, Y. (2022) *Randomized optimal stopping problem in continuous time and reinforcement learning algorithm*. Preprint, arXiv:2208.02409.
- [4] Elworthy, K. D. and Li, X.M. (1994), *Formulae for the derivatives of heat semigroups*. J. Funct. Anal. **125**, 252–286.
- [5] Guo, X., Xu, R., and Zariphopoulou, T. (2022) *Entropy regularization for mean field games with learning*. *Mathematics of Operations research*, **47**(4), 3239–3260.
- [6] Huang, Y., Wang Z., and Zhou, Z. (2022), *Convergence of Policy Improvement for Entropy-Regularized Stochastic Control Problems*, Preprint, arXiv:2209.07059.
- [7] Jacka, S. and Mijatović, A., (2017) *On the policy improvement algorithm in continuous time*, *Stochastics* **89**(1), 348–359.
- [8] Kerimkulov, B., Šiška, D., and Szpruch, L., (2020) *Exponential convergence and stability of Howard’s policy improvement algorithm for controlled diffusions*, *SIAM J. Control Optim.* **58**(3), 1314–1340.
- [9] Kerimkulov, B., Šiška, D., and Szpruch, L., (2021) *A modified MSA for stochastic control problems*, *Appl. Math. Optim.* **84**(3), 3417–3436.
- [10] Ma, J. and Zhang, J. (2002), *Representation Theorems for Backward Stochastic Differential Equations*. *Annals of Applied Probability*, **12**(4), 1390–1418.
- [11] Nualart, D. (2006). *Malliavin calculus and related topics*. In *Stochastic Processes and Related Topics*, 2nd ed., Springer, Berlin.
- [12] Puterman, M. L., (1981) *On the convergence of policy iteration for controlled diffusions*, *J. Optim. Theory Appl.* **33**(1), 137–144.
- [13] Reisinger, C. and Zhang, Y. (2021). *Regularity and stability of feedback relaxed controls*. *SIAM J. Control Optim.* **59**, 3118–3151.
- [14] Tang, W., Zhang, Y. P., and Zhou, X. Y., (2022) *Exploratory HJB equations and their convergence*. *SIAM Journal on Control and Optimization*, **60**(6), 3191–3216.
- [15] Tran, H. V., Wang, Z., and Zhang, Y. .P., (2024) *Policy Iteration for Exploratory Hamilton-Jacobi-Bellman Equations*, Preprint, arXiv:2406.00612.

- [16] Wang, H., Zariphopoulou, T. and Zhou, X.Y., (2020), *Reinforcement learning in continuous time and space: a stochastic control approach*, *J. Mach. Learn. Res.* **21**(198), 1–34.
- [17] Wang, H. and Zhou, X.Y., (2020), *Continuous-time mean-variance portfolio selection: a reinforcement learning framework*, *Math. Finance* **30**(4), 1273–1308.
- [18] Zhang, J. (2001), *Some fine properties of backward stochastic differential equations*. Ph.D. dissertation, Purdue University.