

# DUALITY-BASED *A POSTERIORI* ERROR ESTIMATES FOR SOME APPROXIMATION SCHEMES FOR OPTIMAL INVESTMENT PROBLEMS

ATHENA PICARELLI AND CHRISTOPH REISINGER

**ABSTRACT.** We consider a Markov chain approximation scheme for utility maximization problems in continuous time, which uses, in turn, a piecewise constant policy approximation, Euler-Maruyama time stepping, and a Gauß-Hermite approximation of the Gaussian increments. The error estimates previously derived in *A. Picarelli and C. Reisinger, Probabilistic error analysis for some approximation schemes to optimal control problems, arXiv:1810.04691* are asymmetric between lower and upper bounds due to the control approximation and improve on known results in the literature in the lower case only. In the present paper, we use duality results to obtain *a posteriori* upper error bounds which are empirically of the same order as the lower bounds. The theoretical results are confirmed by our numerical tests.

## 1. INTRODUCTION

We study the numerical approximation of a class of optimal control problems for diffusion processes arising in financial applications. It is well known that, under suitable assumptions, the associated value function can be characterized as the solution of a second order Hamilton-Jacobi-Bellman (HJB) partial differential equation. To deal with the possible degeneracy of the diffusion component of the dynamics, it is in general necessary to consider solutions in the viscosity sense (see [7] for an overview). Furthermore, explicit solutions for this type of nonlinear equations are rarely available, so that their numerical approximation becomes vital. In the framework of viscosity solutions, the basic theory of convergence for numerical schemes is established in [4]. The fundamental properties required are: monotonicity, consistency, and stability of the scheme. While standard finite difference schemes are in general non-monotone, semi-Lagrangian (SL) schemes (see [22, 6, 10]) are monotone by construction. The basic scheme considered in this paper belongs to this family and has been previously analyzed in [25].

We focus here on computable error bounds for the solution. Many of the published error bounds for this kind of maximisation problem, including those in [25], are asymmetrical in the sense that a more accurate lower bound can be given than the upper bound. In this work, we construct an upper bound which consists of two additive contributions: a term which can be computed *a priori* from the model parameters and is of the same order in the mesh parameters as the known lower bounds; and a term which can be computed *a posteriori* from the solution of the dual problem. The practical value of this decomposition is that the second term is empirically (i.e., from our numerical tests) smaller than the first one, so that in practice we can compute rigorous error bounds *a posteriori* which improve on the ones available *a priori*. We discuss this in more detail below.

The machinery for *a priori* bounds for HJB equations is now well-established. By a technique pioneered by Krylov based on “shaking the coefficients” and mollification to construct smooth sub- and/or super-solutions, [19, 21, 1, 2, 3] prove certain fractional convergence orders significantly lower than one. These results are mainly derived by PDE techniques and strongly rely on the comparison principle between viscosity sub- and super-solutions of the HJB equation and the consistency properties of the scheme. For the scheme considered in the present paper, the probabilistic proof in [25] exploits the fact that the numerical scheme is based on a discrete approximation of the optimal control problem, specifically by a piecewise constant policy approximation, Euler-Maruyama time stepping, and a Gauß-Hermite approximation of the Gaussian increments. This yields the desired

error bounds by a direct comparison between two value functions and leads to an improvement of the error contribution of the second and third of these approximations by avoiding the use of the truncation error. The piecewise constant policy approximation, however, introduces an asymmetry between the upper and the lower bound of the error and, as a result, the bounds in [25] give only a partial improvement of the classical PDE-based results.

For the class of convex optimal control problems studied here, namely typical utility maximization problems arising in financial applications, we propose to overcome this issue using information coming from a dual problem. Indeed, an important part of the classical literature dealing with financial applications of optimal control theory (see the seminal work of Kramkov and Schachermayer [18]) applies duality techniques to solve utility maximization problems under suitable convexity assumptions. The basic idea of this method is to write the optimal control problem as a constrained optimization problem with respect to the state variable and then solve it by convex analysis techniques. A systematic approach to utility maximization problems admitting a dual formulation is discussed in [26]. Of these, the fairly general set-up of an optimal investment problem involving nonlinear dynamics given in [9] will be explicitly analyzed in this paper.

More specifically, a direct application of the results in [25] to this problem gives one-sided (lower) error bounds for the considered Markov chain approximation of order

$$h^{(M-1)/2M} + \Delta x^{(M-1)/(3M-1)} \quad (1.1)$$

for timestep  $h$ , spatial mesh size  $\Delta x$  and number of Gaussian points  $M$ , for Lipschitz viscosity solutions. They coincide with the two-sided bounds in [10] for the standard linear-interpolation SL scheme, i.e.  $M = 2$ , and improve them for  $M > 2$ . In contrast, the piecewise constant policy approximation introduces an extra term in the upper bound of order  $h^{1/4}$  (from a recent result in [16]), which strictly restricts the order for  $M > 2$ .

The main contribution of this paper is to analyse the error estimates in the case of optimal investment problems. Their special structure has neither been exploited by the classical literature on PDE-based error estimates for HJB equations nor by the analysis in [25]. We prove that for the class of problems analyzed here, two-sided *a posteriori* bounds of the empirical order (1.1) can be obtained. As a side result, we complete the literature by deriving explicit values for the constants appearing in the error estimates in terms of the Lipschitz (resp. Hölder) regularity of the coefficients and the solution in space (resp. time).

The paper is organised as follows. In Section 2, we introduce the problem set-up and state our assumptions. We define the scheme and give *a priori* lower error bounds for the primal problem in Section 3, and both *a priori* and *a posteriori* upper bounds, by way of the dual problem, in Section 4. We illustrate the theoretical results by numerical tests in Section 5, and offer conclusions and extensions in Section 6. In Appendix A, we derive explicit expressions for the constants in the error bounds.

## 2. MAIN ASSUMPTIONS AND PRELIMINARY RESULTS

Let  $(\Omega, \mathbb{F}, \mathbb{P})$  be a probability space with filtration  $\{\mathbb{F}_t, t \geq 0\}$  induced by a  $d$ -dimensional Brownian motion  $B$  and let  $T > 0$ . We consider a controlled (scalar) process governed by a dynamics of the following form, for  $t \in [0, T)$ ,

$$\begin{cases} dX_s = X_s \left( r(s) + \alpha_s^\top (b(s) - r(s)\mathbb{1}) + g(s, \alpha_s) \right) ds + X_s \alpha_s^\top \sigma(s) dB_s, & s \in (t, T) \\ X_t = x \geq 0, \end{cases} \quad (2.1)$$

where  $r, b, g$  and  $\sigma$  take values, respectively, in  $\mathbb{R}, \mathbb{R}^d, \mathbb{R}$  and  $\mathbb{R}^{d \times d}$  and  $\mathbb{1} \equiv (1, \dots, 1)^\top \in \mathbb{R}^d$ . Denote further by  $\mathcal{A}$  the set of control policies, i.e. progressively measurable processes  $\alpha$  taking values in a given set  $A \subseteq \mathbb{R}^d$  such that  $\int_0^T |\alpha_s|^2 ds < +\infty$ . This framework has been introduced and studied in [9], and encompasses a number of important optimal investment problems involving nonlinear dynamics, including the classical Merton problem [23], as special cases. In such models, the state  $X$  typically represents the wealth of an investor with initial endowment  $x$  at time  $t$ . The control

vector  $\alpha \equiv (\alpha_1, \dots, \alpha_d)^\top$  then determines the proportion of wealth the investor puts in each stock. Here, the coefficient  $r$  is the return rate of a bond (riskless asset), while  $b(\cdot) \equiv (b_1(\cdot), \dots, b_d(\cdot))^\top$  is the vector of the appreciation rates of the  $d$  considered stocks with volatility matrix  $\sigma(\cdot)$ . The nonlinearity in the investment strategy introduced by the function  $g$  models the effects of market frictions and trading constraints on the wealth (see [9, 8, 12]). We refer the reader to [26] for an overview of different utility maximization problems, including (2.1) and its special cases. We consider the following assumptions:

(H1)  $A \subseteq \mathbb{R}^d$  is a bounded and convex set such that  $0 \in A$ .

(H2) (i) There exists  $K_0 \geq 0$  such that

$$|r(t) - r(s)| + |b(t) - b(s)| + \|\sigma(t) - \sigma(s)\| \leq K_0 |t - s|^{1/2} \quad \forall t, s \in [0, T].$$

(ii)  $g : [0, T] \times A \rightarrow \mathbb{R}$  satisfies:

- there exists  $K_1 \geq 0$  such that

$$|g(t, a) - g(t, a')| \leq K_1 |a - a'| \quad \forall a, a' \in A, t \in [0, T];$$

$$|g(t, a) - g(s, a)| \leq K_1 |t - s|^{1/2} \quad \forall t, s \in [0, T], a \in A;$$

- for each  $t \in [0, T]$ ,  $a \rightarrow g(t, a)$  is concave;

-  $g(t, 0) = 0$  for all  $t \in [0, T]$ .

(H3)  $\sigma$  satisfies a uniform ellipticity condition, i.e. there exists  $\eta > 0$  such that

$$\xi^\top \sigma \sigma^\top \xi \geq \eta |\xi|^2 \quad \forall \xi \in \mathbb{R}^d.$$

One has the following existence and uniqueness result:

**Lemma 2.1.** *Let assumptions (H1) to (H3) be satisfied. For any choice of the control  $\alpha \in \mathcal{A}$  and  $x \geq 0$  there exists a unique strong solution to equation (2.1).*

*Proof.* For  $x > 0$ , a solution can be defined as  $X_\cdot = \exp(Z_\cdot)$ , where

$$Z_\cdot = z + \int_t^\cdot r(s) + \alpha_s^\top (b(s) - r(s)\mathbb{1}) + g(s, \alpha_s) - \frac{1}{2}(\alpha_s^\top \sigma)^2 ds + \int_t^\cdot \alpha_s^\top \sigma(s) dB_s,$$

for  $z = \log x$ , which is well defined under assumptions (H1)-(H3) for any  $\alpha \in \mathcal{A}$ . Moreover, for  $x = 0$  the process  $X \equiv 0$  is the unique solution to (2.1) for any  $\alpha \in \mathcal{A}$ .  $\square$

We denote by  $X_\cdot^{t,x,\alpha}$  the unique solution of equation (2.1). To simplify the notation, where no ambiguities arise, we will indicate the starting point  $(t, x)$  of the processes involved as a subscript in the expectation, i.e.  $\mathbb{E}_{t,x}[\cdot]$ .

The value function  $v : [0, T] \times [0, +\infty) \rightarrow \mathbb{R}$  of the optimal control problem is defined by

$$v(t, x) := \sup_{\alpha \in \mathcal{A}} \mathbb{E}_{t,x} [U(X_T^\alpha)], \quad (2.2)$$

where  $U : [0, +\infty) \rightarrow \mathbb{R}$  is the so-called utility function of the investor and it is assumed to satisfy the following assumptions:

(H4)  $U \in C^1((0, +\infty); \mathbb{R})$ ;

$U$  is concave and strictly increasing;

$\lim_{x \rightarrow +\infty} U'(x) = 0$ .

For any  $[0, T - t]$ -valued stopping time  $\theta$ ,  $v$  satisfies the Dynamic Programming Principle (DPP)

$$v(t, x) = \sup_{\alpha \in \mathcal{A}} \mathbb{E}_{t,x} [v(t + \theta, X_{t+\theta}^\alpha)], \quad (2.3)$$

from which, at least formally, one can show that the Hamilton-Jacobi-Bellman (HJB) equation associated with the optimal control problem (2.2) is

$$-v_t + \sup_{a \in A} \left( -x (r(t) + a^\top (b(t) - r(t)\mathbb{1}) + g(t, a)) v_x - \frac{1}{2} x^2 \text{Tr}[a(\sigma \sigma^\top)(t) a^\top] v_{xx} \right) = 0 \quad (2.4)$$

for  $t \in [0, T]$ ,  $x \geq 0$ , completed with the terminal condition  $v(T, x) = U(x)$  for  $x \geq 0$  (see [24, Section 3.6.1]). We refer the reader to [27, Section 3, Chapter 4] and the references therein for a complete overview on the dynamic programming approach to optimal control problems.

In the general case,  $v$  is not expected to have sufficient regularity to satisfy the previous equation in the classical sense and even if (2.4) admits a classical solution, it is rarely found explicitly. To handle the problem in its full generality, the notion of viscosity solution is needed (see [7] for an overview). Indeed, under suitable assumptions, it can be proved (see for instance [27, Theorems 5.2 and 6.1]) that  $v$  defined in (2.2) is the unique continuous viscosity solution to (2.4) on  $[0, T] \times [0, +\infty)$ .

### 3. THE NUMERICAL SCHEME

We consider here the scheme analyzed in [25]. It belongs to the family of the so-called semi-Lagrangian (SL) schemes (see [6, 11, 20, 22] for their earlier introduction) which are based on discretization of the control set  $\mathcal{A}$  and a Markov chain approximation of the associated optimal control problem. For completeness, we briefly discuss below the main features of the scheme. We refer the reader to [25] for further details.

**3.1. Description of the scheme.** We start by introducing a discretization in time. Let  $N \geq 1$ ,

$$h = T/N \quad \text{and} \quad t_n = nh,$$

for  $n = 0, \dots, N$ . The first step in our approximation is to introduce a time discretization of the control set. We consider the set  $\mathcal{A}_h$  of controls  $\alpha \in \mathcal{A}$  which are constant in each interval  $[t_n, t_{n+1}]$ , for  $n = 0, \dots, N-1$ , i.e.

$$\mathcal{A}_h := \left\{ \alpha \in \mathcal{A} : \alpha_s(\omega) \equiv \sum_{i=0}^{N-1} a_i \mathbb{1}_{s \in [t_i, t_{i+1})} \quad \forall \omega \in \Omega \text{ s.t. } a_i \in A, \quad i = 0, \dots, N-1 \right\}.$$

In what follows, we identify any element  $\alpha \in \mathcal{A}_h$  by the sequence of random variables  $a_i$  taking values in  $A$  (denoted by  $a_i \in A$  for simplicity) and will write  $\alpha \equiv (a_0, \dots, a_{N-1})$ . We denote by  $v_h$  the value function obtained by restricting the supremum in (2.2) to controls in  $\mathcal{A}_h$ , that is

$$v_h(t, x) := \sup_{\alpha \in \mathcal{A}_h} \mathbb{E}_{t,x} [U(X_T^\alpha)]. \quad (3.1)$$

Clearly, since  $\mathcal{A}_h \subseteq \mathcal{A}$ , one has

$$v(t, x) \geq v_h(t, x), \quad (3.2)$$

for any  $t \in [0, T]$ ,  $x \geq 0$ . An upper bound of order  $1/6$  for the error related to this approximation was first obtained by Krylov in [20]. Recently, this estimate has been improved to the order  $1/4$  in [16], so that one has

$$v(t, x) \leq v_h(t, x) + Ch^{1/4} \quad (3.3)$$

for some constant  $C \geq 0$ . We point out that the results in [20] and [16] require some additional assumptions on the coefficients and do not directly apply to problem (2.1) to (2.2). It is possible that analogous estimates hold also in the setting of the present paper, but since we do not make use of (3.3) here, we did not check this point in detail. Indeed, a main objective of the present paper is to by-pass the estimate (3.3), which turns out to be a bottleneck in the provable approximation order, while still using the piecewise constant policy approximation itself by building an approximation to  $v_h$ . The more important observation from [16] is therefore that a better order than  $1/4$  is not provable in the general case of Lipschitz viscosity solutions. Then no matter how precise the estimates obtained for the error of the final approximation to  $v_h$  are, without any further information the upper error bounds to  $v$  cannot be more accurate than  $O(h^{1/4})$ . Section 4 will show how this term can be replaced by an expression which is computable from the dual problem and provides sharper bounds in our tests (see Section 5).

|         | $\xi_i$                | $\lambda_i$       |
|---------|------------------------|-------------------|
| $M = 2$ | $\pm 1$                | $1/2$             |
| $M = 3$ | $0$                    | $2/3$             |
|         | $\pm\sqrt{3}$          | $1/6$             |
| $M = 4$ | $\pm\sqrt{3-\sqrt{6}}$ | $(3+\sqrt{6})/12$ |
|         | $\pm\sqrt{3+\sqrt{6}}$ | $(3-\sqrt{6})/12$ |

FIGURE 1. Analytical expressions of  $\{(\xi_i, \lambda_i)\}_{i=1,\dots,M}$  for  $M = 2, 3, 4$ . We refer to [5, p. 464] for numerical approximations of  $\{(z_i, \omega_i)\}_{i=1,\dots,M}$  for larger  $M$ .

For any given  $\alpha \equiv (a_0, \dots, a_{N-1}) \in \mathcal{A}_h$ , we consider the Euler-Maruyama approximation of the process  $X^{t,x,\alpha}$  given by the following recursive relation:

$$X_{t_{i+1}} = X_{t_i} + h X_{t_i} (r(t_i) + a_i^\top (b(t_i) - r(t_i)\mathbb{1}) + g(t_i, a_i)) + X_{t_i} a_i^\top \sigma(t_i) \Delta B_i \quad (3.4)$$

for  $i = 0, \dots, N-1$ . The increments  $\Delta B_i := B_{t_{i+1}} - B_{t_i}$  are independent, identically distributed random variables such that

$$\Delta B_i \sim \sqrt{h} \mathcal{N}(0, I_d) \quad \forall i = 0, \dots, N-1. \quad (3.5)$$

We denote by  $\bar{X}^{t_n, x, \alpha}$  the solution to (3.4) with the control  $\alpha \equiv (a_0, \dots, a_{N-1}) \in \mathcal{A}_h$  and such that  $\bar{X}_{t_n}^{t_n, x, \alpha} = x$ . In the next step, we work towards a Markov chain approximation of  $\bar{X}^{t_n, x, \alpha}$ .

Let us start for simplicity with the case  $d = 1$ . Let  $M \geq 2$  and denote by  $\{z_i\}_{i=1,\dots,M}$  the zeros of the Hermite polynomial  $H_M$  of order  $M$  and by  $\{\omega_i\}_{i=1,\dots,M}$  the corresponding weights given by

$$\omega_i = \frac{2^{M-1} M! \sqrt{\pi}}{M^2 [H_{M-1}(z_i)]^2}, \quad i = 1, \dots, M.$$

With the definitions

$$\lambda_i := \frac{\omega_i}{\sqrt{\pi}} \quad \text{and} \quad \xi_i := \sqrt{2} z_i, \quad i = 1, \dots, M,$$

one can make use of the following approximation (see, e.g., [14, p. 395])

$$\int_{-\infty}^{+\infty} f(y) \frac{e^{-\frac{y^2}{2}}}{\sqrt{2\pi}} dy \approx \sum_{i=1}^M \lambda_i f(\xi_i), \quad (3.6)$$

which holds for any smooth real-valued function  $f$  (say  $f$  at least  $C^{2M}$ ). Observing that  $\lambda_i \geq 0, \forall i = 1, \dots, M$ , and  $\sum_{i=1}^M \lambda_i = 1$ , given the sequence  $\{\zeta_n\}_{n=0,\dots,N-1}$  of i.i.d. random variables such that for any  $n = 0, \dots, N-1$

$$\mathbb{P}(\zeta_n = \xi_i) = \lambda_i, \quad i = 1, \dots, M,$$

one has

$$\mathbb{E}[\zeta_n] = 0 \quad \text{and} \quad \text{Var}[\zeta_n] = 1 \quad \forall n = 0, \dots, N-1.$$

For any control  $\alpha \equiv (a_0, \dots, a_{N-1}) \in \mathcal{A}_h$ , we will denote by  $\hat{X}^{t_n, x, \alpha}$  the Markov chain approximation of the process  $\bar{X}^{t_n, x, \alpha}$ , i.e.

$$\begin{cases} \hat{X}_{t_n} = x, \\ \hat{X}_{t_{i+1}} = \hat{X}_{t_i} + h \hat{X}_{t_i} (r(t_i) + a_i (b(t_i) - r(t_i)) + g(t_i, a_i)) + \sqrt{h} \hat{X}_{t_i} a_i \sigma(t_i) \zeta_i, \end{cases} \quad (3.7)$$

for  $i = n, \dots, N-1$ .

Applying to (2.3) with  $\theta = h$  the piecewise control approximation, the Euler-Maruyama discretization and the Gauß-Hermite quadrature formula (3.6), we obtain the following recursive



**Proposition 3.1.** *Let assumptions (H1) to (H3) be satisfied and let the function  $U$  be Lipschitz continuous with Lipschitz constant  $L \geq 0$ . Then, there exists a constant  $C \geq 0$  such that for any  $n = 0, \dots, N$ ,  $x \geq 0$*

$$v(t_n, x) \geq V(t_n, x) - LC(1 + x^{2M})h^{(M-1)/2M}, \quad (3.12)$$

and for any  $n = 0, \dots, N$ ,  $m \in \mathbb{N}$

$$v(t_n, x_m) \geq W(t_n, x_m) - LC(1 + x_m^{2M}) \left( h^{(M-1)/2M} + \Delta x/h \right). \quad (3.13)$$

*Proof.* When only the time discretization is taken into account, the estimate (3.12) directly follows by [25, Section 4.2, equation (4.9)]. Moreover by [25, Section 4.3] for the fully discrete scheme one has (3.13), where  $C$  only depends on  $M$ ,  $T$ , the constants  $K_0$  and  $K_1$  in assumption (H2).  $\square$

**Remark 1.** *Balancing the terms  $h^{(M-1)/2M}$  and  $\Delta x/h$  on the right-hand side of (3.13) by judicious choice of  $\Delta x$  in relation to  $h$  leads to*

$$v(t_n, x_m) - W(t_n, x_m) \geq -LC(1 + x_m^{2M}) \left( h^{(M-1)/2M} + \Delta x^{(M-1)/(3M-1)} \right).$$

The scheme we are considering is monotone, stable and it has order one of consistency (for smooth test functions) for any  $M \geq 2$ . For a scheme of this type, (upper and lower) error bounds of order  $1/4$  in  $h$  have been provided in [3, 10] by PDE techniques. Splitting each contribution to the error, namely the control discretization and the Euler-Maruyama and Gauß-Hermite approximations, the probabilistic proof proposed in [25] gives an improvement to the lower bound of these estimates by increasing the value of  $M > 2$ , i.e. by using a more accurate quadrature formula. For large  $M$ , the order is arbitrarily close to  $1/2$  in  $h$  and  $1/3$  in  $\Delta x$ , improving the corresponding orders  $1/4$  and  $1/5$  obtained for  $M = 2$ . It is hence for  $M > 2$  that the term  $h^{1/4}$  in the upper bound becomes strictly dominant and restricts the order to  $1/4$ , independent of  $M$ . The analysis that follows aims to eliminate this dominant term and replace it by a term which can be computed *a posteriori* from the numerical solution of the original and its dual problem, which we expect to be smaller generally than that from the Gauß-Hermite approximations. This is confirmed in our tests. Hence we provide a computable upper bound which is empirically of the same order as the lower bounds obtained in Proposition 3.1.

#### 4. DUALITY-BASED ERROR ESTIMATES

In this section we discuss how duality theory can be employed to obtain an upper bound of the error associated with our approximation scheme. Assuming to be able to extend either the PDE-based error estimates in [19, 21, 1, 2, 3] or the probabilistic ones in [25] to the particular problem (2.1) and (2.2), this would result in both cases in an upper bound of order  $1/4$  in  $h$  for any choice of  $M \geq 2$ , as explained at the end of the previous section. We show here that for our class of problems it is possible to pass through the definition of a dual problem to replace these *a priori* estimates by *a posteriori* computable bounds, which are empirically significantly smaller.

**4.1. The dual problem.** The dual problem associated with (2.1) and (2.2) is defined in [9] by

$$\begin{cases} dY_s = -(r(s) + \tilde{g}(s, \nu_s))Y_s ds + (\sigma\sigma^T(s))^{-1}Y_s(r(s)\mathbb{1} - b(s) - \nu_s) \cdot \sigma(s) dB_s, & s \in (t, T), \\ Y_t = y, \end{cases} \quad (4.1)$$

for all  $t \in [0, T)$ , where

$$\tilde{g}(t, \nu) := \sup_{a \in A} \{g(t, a) - a \cdot \nu\}, \quad (4.2)$$

and the dual utility function  $\tilde{U}$  is the convex conjugate of  $U$ , i.e.

$$\tilde{U}(y) := \sup_{x \geq 0} \{U(x) - xy\}.$$

The dual value function is defined by

$$\tilde{v}(t, y) = \inf_{\nu \in \mathcal{V}} \mathbb{E}_{t, y} [\tilde{U}(Y_T^\nu)], \quad (4.3)$$

where  $\mathcal{V}$  is the set of  $\mathbb{R}^d$ -valued progressively measurable processes such that  $\int_0^T |\nu_s|^2 ds + \int_0^T \tilde{g}(s, \nu_s) ds < +\infty$ . One has the following duality result:

**Proposition 4.1** ([9], Theorem 2). *Let assumptions (H1) to (H4) be satisfied. Then for any  $t \in [0, T]$ ,  $x \geq 0$ , the primal and dual value functions,  $v$  and  $\tilde{v}$ , satisfy the conjugate relation*

$$v(t, x) = \inf_{y \geq 0} \{ \tilde{v}(t, y) + xy \}. \quad (4.4)$$

**Remark 2.** *The results in [9] hold also if  $r, b, \sigma$  and  $g$  are stochastic processes. However, as our approximation scheme makes use of the Markovian structure, we would have to add extra variables to the state space to account for this, which is outside the scope of this work.*

**Remark 3.** *As mentioned in [26, Section 6.5], the usual Inada condition*

$$\lim_{x \rightarrow 0^+} U'(x) = +\infty$$

*(requested in [9]) is not necessary for proving the main duality results.*

**4.2. Approximation of the dual problem.** The scheme presented in Section 3.1 can be used to approximate the value function  $\tilde{v}$  associated with the dual problem (4.1)-(4.3). To this end, we define by  $\Gamma \subset \mathbb{R}^d$  a compact set and by  $\mathcal{V}^\Gamma \subset \mathcal{V}$  the set of all a.s.  $\Gamma$ -valued elements of  $\mathcal{V}$ . One clearly has

$$\tilde{v}(t, y) \leq \tilde{v}^\Gamma(t, y) = \inf_{\nu \in \mathcal{V}^\Gamma} \mathbb{E}_{t, y} [\tilde{U}(Y_T^\nu)]. \quad (4.5)$$

If there exists a uniformly bounded optimal control  $\nu^* \in \mathcal{V}$ , one can find a compact set  $\Gamma$  such that  $\tilde{v} = \tilde{v}^\Gamma$ . Otherwise, such an approximation introduced on the set of controls will result in a strictly bigger value function and in a duality gap which does not diminish under mesh refinement and can only be decreased by increasing  $\Gamma$ . Nonetheless, the inequalities stated in this section still hold in this case.

Let further  $\zeta_i$ ,  $i = 0, \dots, N-1$ , be i.i.d. copies of the increments from the definition of the primal approximation. For the discrete time scheme, one can then recursively define

$$\begin{cases} \tilde{V}(t_n, y) = \inf_{\gamma \in \Gamma} \mathbb{E}_{t_n, y} [\tilde{V}(t_{n+1}, \hat{Y}_{t_{n+1}}^\gamma)], & n = N-1, \dots, 0, \\ \tilde{V}(t_N, y) = \tilde{U}(y), \end{cases} \quad (4.6)$$

where  $y \geq 0$ , and  $\hat{Y}_{t_n, y, \gamma}$  is the Markov chain recursively defined by

$$\begin{cases} \hat{Y}_{t_n} = y, \\ \hat{Y}_{t_{i+1}} = \hat{Y}_{t_i} - h(r(t_i) + \tilde{g}(t_i, \gamma_i)) \hat{Y}_{t_i} + \sqrt{h}(\sigma \sigma^T(t_i))^{-1} \hat{Y}_{t_i} (r(t_i) \mathbb{1} - b(t_i) - \gamma_i) \cdot \sigma(t_i) \zeta_i \end{cases} \quad (4.7)$$

for  $i = n, \dots, N-1$ . Denoting by  $\mathcal{V}_h^\Gamma$  the set of all  $\nu \equiv (\gamma_0, \dots, \gamma_{N-1})$  adapted to the filtration generated by  $(\zeta_0, \dots, \zeta_{N-1})$ , with  $\gamma_i$  random variables taking values in  $\Gamma$ ,  $i = 0, \dots, N-1$  one has

$$\tilde{V}(t_n, y) = \inf_{\nu \in \mathcal{V}_h^\Gamma} \mathbb{E}_{t_n, y} [\tilde{U}(\hat{Y}_T^\nu)].$$

The fully discrete version of the scheme is then given by

$$\begin{cases} \tilde{W}(t_n, y_j) = \inf_{\gamma \in \Gamma} \sum_{i=1}^M \lambda_i \mathcal{I}[\tilde{W}](t_{n+1}, y_j + h y_j (r(t_n) + \tilde{g}(t_n, \gamma)) \\ \quad \quad \quad + \sqrt{h} y_j (\sigma \sigma^T(t_n))^{-1} (r(t_n) \mathbb{1} - b(t_n) - \gamma) \cdot \sigma(t_n) \xi_i), \\ \tilde{W}(t_N, y_j) = \tilde{U}(y_j), \end{cases} \quad (4.8)$$

for  $n = N-1, \dots, 0$  and  $j \in \mathbb{N}$ .



**4.3. An *a priori* upper bound for  $\tilde{v}$ .** The approximation scheme we defined for the dual problem is the same we used for the primal one, with the only difference that we have to handle a minimization problem. Therefore, we can use the arguments in [25] to obtain an accurate upper bound for the differences  $\tilde{v} - \tilde{V}$  and  $\tilde{v} - \tilde{W}$ .

**Proposition 4.2.** *Let assumptions (H1) to (H4) be satisfied and let  $\tilde{U}$  be Lipschitz continuous with Lipschitz constant  $\tilde{L} \geq 0$ . Then, there exists a constant  $\tilde{C} \geq 0$ , such that for any  $n = 0, \dots, N$ ,  $y > 0$ ,*

$$\tilde{v}(t_n, y) \leq \tilde{V}(t_n, y) + \tilde{L}\tilde{C}(1 + y^{2M})h^{(M-1)/2M}, \quad (4.9)$$

and for any  $n = 0, \dots, N$ ,  $j \in \mathbb{N}^+$ ,

$$\tilde{v}(t_n, y_j) \leq \tilde{W}(t_n, y_j) + \tilde{L}\tilde{C}(1 + y_j^{2M}) \left( h^{(M-1)/2M} + \Delta x/h \right). \quad (4.10)$$

*Proof.* The result follows by applying the estimates from Section 4.2 in [25], adapted to a minimization problem, to  $\tilde{v}^\Gamma$  in (4.5). Under the assumptions (H1)-(H3), one has

$$\begin{aligned} & \|(\sigma\sigma^T(t))^{-1}(r(t)\mathbb{1} - b(t) - \gamma) - (\sigma\sigma^T(s))^{-1}(r(s)\mathbb{1} - b(s) - \gamma)\| \\ & \quad + |\tilde{g}(t, \gamma) - \tilde{g}(s, \gamma)| + |r(t) - r(s)| \leq C_0|t - s|^{1/2}, \end{aligned}$$

where  $C_0$  depends on  $T$ , the constants  $K_0, K_1$  and  $\eta$  in assumptions (H2)-(H3) and the uniform bounds on the elements of  $\Gamma$ , so that the dynamics (4.1) satisfies the assumptions in [25].  $\square$

**Remark 4.** *A similar truncation strategy as for  $\mathcal{V}$  can be applied to the set of controls  $\mathcal{A}$  if  $A$  is unbounded. For the thus obtained numerical solution, the inequalities in Proposition 3.1 still hold. Again, in this case, the duality gap can only be reduced by increasing  $A$  and not only by letting  $h$  and  $\Delta x$  go to 0 alone.*

**4.4. Using duality in error estimates.** In the sequel, we will use the following notation: for any  $n = 0, \dots, N$ ,  $x > 0$

$$G^h(t_n, x) := \min_{y>0} \left\{ \tilde{V}(t_n, y) + xy \right\} - V(t_n, x) \quad (4.11)$$

$$I^h(t_n, x) := \arg \min_{y>0} \left\{ \tilde{V}(t_n, y) + xy \right\},$$

and for any  $n = 0, \dots, N$ ,  $m \in \mathbb{N}^+$

$$G^{h, \Delta x}(t_n, x_m) := \min_{j \in \mathbb{N}^+} \left\{ \tilde{W}(t_n, y_j) + x_m y_j \right\} - W(t_n, x_m), \quad (4.12)$$

$$I^{h, \Delta x}(t_n, x_m) := \arg \min_{j \in \mathbb{N}^+} \left\{ \tilde{W}(t_n, y_j) + x_m y_j \right\}.$$

We refer to  $G^h$  and  $G^{h, \Delta x}$  as the *numerical duality gap* of the semidiscrete and fully discrete scheme respectively.

One has the following result:

**Theorem 4.1.** *Let assumptions (H1) to (H4) be satisfied and let  $U$  and  $\tilde{U}$  be Lipschitz continuous with Lipschitz constants  $L$  and  $\tilde{L}$ , respectively. Then, there exist some constants  $C, \tilde{C} \geq 0$  such that for any  $n = 0, \dots, N$ ,  $x > 0$*

$$-LC(1 + x^{2M})h^{(M-1)/2M} \leq v(t_n, x) - V(t_n, x) \leq G^h(t_n, x) + \tilde{L}\tilde{C}(1 + (I^h(t_n, x))^{2M})h^{(M-1)/2M} \quad (4.13)$$

and for any  $n = 0, \dots, N$ ,  $m \in \mathbb{N}^+$

$$\begin{aligned} & -LC(1 + x_m^{2M}) \left( h^{(M-1)/2M} + \Delta x/h \right) \leq v(t_n, x_m) - W(t_n, x_m) \\ & \leq G^{h, \Delta x}(t_n, x_m) + \tilde{L}\tilde{C}(1 + (I^{h, \Delta x}(t_n, x_m))^{2M}) \left( h^{(M-1)/2M} + \Delta x/h \right). \end{aligned} \quad (4.14)$$

*Proof.* The first inequalities in (4.13) and (4.14) follow directly by Proposition 3.1. It remains to prove the upper bounds. We prove the result for the semi-discrete scheme, while the proof for the fully discrete scheme follows by similar arguments. Thanks to Proposition 4.1, Proposition 4.2 and the definition of  $I^h(\cdot, \cdot)$  one has

$$\begin{aligned} v(t_n, x) &= \inf_{y>0} \{ \tilde{v}(t_n, y) + xy \} \leq \inf_{y>0} \left\{ \tilde{V}(t_n, y) + xy + \tilde{L}\tilde{C}(1 + y^{2M})h^{(M-1)/2M} \right\} \\ &\leq \tilde{V}(t_n, I^h(t_n, x)) + x I^h(t_n, x) + \tilde{L}\tilde{C} (1 + (I^h(t_n, x))^{2M}) h^{(M-1)/2M} \\ &= \inf_{y>0} \left\{ \tilde{V}(t_n, y) + xy \right\} + \tilde{L}\tilde{C} (1 + (I^h(t_n, x))^{2M}) h^{(M-1)/2M}. \end{aligned}$$

Therefore,

$$v(t_n, x) - V(t_n, x) \leq \inf_{y>0} \left\{ \tilde{V}(t_n, y) + xy \right\} - V(t_n, x) + \tilde{L}\tilde{C} (1 + (I^h(t_n, x))^{2M}) h^{(M-1)/2M},$$

which gives the desired result.  $\square$

Observe that due to the particular convexity feature of the dual problem, the quantity  $I(x)$  typically increases as  $x$  approaches 0.

The duality gap for the fully discrete scheme is computable efficiently, see e.g. [13, Section 3.4], so that (4.14) provides a practically useful *a posteriori* bound.

*A priori* bounds could be obtained by proving that the numerical duality gap  $G^h$  (resp.  $G^{h,\Delta x}$ ) decays with order at most  $h^{(M-1)/2M}$  (resp.  $h^{(M-1)/2M} + Dx/h$ ). This requires a proof that  $V$  and  $\tilde{V}$  (resp.  $W$  and  $\tilde{W}$ ) satisfy an approximated duality relation. Indeed, the key feature of dynamics (2.1) and (4.1) leading to the conjugate relation (4.4) is the following so called “polar property”

$$\sup_{\nu \in \mathcal{V}} \mathbb{E} [X_T^{t,x,\alpha} Y_T^{t,y,\nu}] = xy \quad \forall x, y \geq 0, t \in [0, T], \alpha \in \mathcal{A}.$$

For the discrete time dynamics  $\hat{X}$  and  $\hat{Y}$  defined in (3.7) and (4.7), respectively, a straightforward calculation shows that for any  $\alpha \equiv (a_n, \dots, a_{N-1}) \in \mathcal{A}_h$  and  $\nu \equiv (\gamma_n, \dots, \gamma_{N-1}) \in \mathcal{V}_h^\Gamma$

$$\begin{aligned} &\hat{X}_T^{t_n,x,\alpha} \hat{Y}_T^{t_n,y,\nu} - xy \\ &= \sum_{i=n}^{N-1} \hat{X}_{t_i}^{t_n,x,\alpha} \hat{Y}_{t_i}^{t_n,y,\nu} \left\{ h \left( a_i(b(t_i) - r(t_i)) + g(t_i, a_i) - \tilde{g}(t_i, \gamma_i) + a_i(r(t_i) - b(t_i) - \gamma_i) \zeta_i^2 \right) \right. \\ &\quad \left. + h^2 \left( - (r(t_i) + \tilde{g}(t_i, \gamma_i)) (r(t_i) + a_i(b(t_i) - r(t_i)) + g(t_i, a_i)) \right) \right. \\ &\quad \left. + (\dots) \zeta_i + (\dots) (\zeta_i^2 - 1) \right\}. \end{aligned}$$

Taking the expectation in the expression above, thanks to the independence and distribution of the random variables  $\zeta_i$  and the definition of the convex conjugate  $\tilde{g}$ , one gets

$$\mathbb{E} \left[ \hat{X}_T^{t_n,x,\alpha} \hat{Y}_T^{t_n,y,\nu} \right] - xy \leq Ch^2 \sum_{i=n}^{N-1} \mathbb{E} \left[ \hat{X}_{t_i}^{t_n,x,\alpha} \hat{Y}_{t_i}^{t_n,y,\nu} \right]$$

for some constant  $C$  depending on  $T$ , the uniform bounds on  $A$  and  $\Gamma$  and the constants appearing in assumption (H2). For any  $i = n, \dots, N-1$  one can easily prove that

$$\mathbb{E} \left[ \hat{X}_{t_i}^{t_n,x,\alpha} \hat{Y}_{t_i}^{t_n,y,\nu} \right] \leq xy C e^{CT},$$

for some possibly different constant  $C \geq 0$ , so that it is possible to conclude that there exists some  $C \geq 0$  such that

$$\sup_{\nu \in \mathcal{V}_h^\Gamma} \mathbb{E} \left[ \hat{X}_T^{t_n,x,\alpha} \hat{Y}_T^{t_n,y,\nu} \right] \leq xy (1 + Ch) \quad \forall x, y \geq 0, n = 0, \dots, N, \alpha \in \mathcal{A}.$$

We conjecture that a similar approximate lower bound also holds. This finds a confirmation in our numerical tests (see Tables 3 and 5 in Section 5) where at least first order of convergence in  $h$

of the numerical duality gap is observed. However the rigorous prove of the result involves delicate convex analysis arguments and we plan to investigate this point in future work.

**4.5. The case of non Lipschitz utility functions.** We assumed for the results above that the primal (and, where applicable, dual) utility functions  $U$  (and  $\tilde{U}$ ) are Lipschitz continuous (see Propositions 3.1, 4.2, and Theorem 4.1). This is a standard assumption in the numerical literature, including our previous work [25] which we draw on here. This property is, however, not satisfied by commonly used utility functions in finance, such as the power utility  $U(x) = x^p/p$ ,  $x \geq 0$ , with  $p \in (0, 1)$ , or the dual of the exponential utility. To deal with such cases, we introduce a further approximation of the problem and consequently have to estimate an additional error contribution.

We assume first, in addition to (H4), that  $U$  is bounded (from below) at 0. Letting  $\rho, c_0 > 0$  and  $x_\rho = c_0/\rho$ ,  $y_\rho = \rho$ , we define

$$U_\rho(x) := \begin{cases} U(0) + \frac{U(x_\rho) - U(0)}{x_\rho}x & \text{if } 0 \leq x \leq x_\rho, \\ U(x) & \text{if } x_\rho < x \leq y_\rho, \\ U(y_\rho) & \text{if } x > y_\rho, \end{cases} \quad (4.15)$$

so that  $U_\rho$  is Lipschitz with Lipschitz constant  $L_\rho := (U(x_\rho) - U(0))/x_\rho$  and  $U_\rho \rightarrow U$  as  $\rho \rightarrow +\infty$  (uniformly on compact sets).

We denote by  $v_\rho$ ,  $V_\rho$  and  $W_\rho$  the value function and the numerical approximations defined respectively by (2.2), (3.8) and (3.10), replacing  $U$  with  $U_\rho$ . Observe that as  $U$  is concave and therefore  $U_\rho \leq U$ , one has for any  $t \in [0, T]$ ,  $x \geq 0$

$$v_\rho(t, x) \leq v(t, x). \quad (4.16)$$

Let  $\tilde{U}_\rho$  be the convex conjugate of the approximated utility function  $U_\rho$ , i.e.

$$\tilde{U}_\rho(y) := \sup_{x \geq 0} \{U_\rho(x) - xy\}.$$

We denote by  $\tilde{v}_\rho$ ,  $\tilde{V}_\rho$  and  $\tilde{W}_\rho$  the value function and the numerical approximations obtained respectively by (4.3), (4.6) and (4.8), replacing  $\tilde{U}$  with  $\tilde{U}_\rho$ . Observe that  $\tilde{U}_\rho : [0, +\infty) \rightarrow \mathbb{R}$  is decreasing and Lipschitz continuous with constant  $\tilde{L}_\rho := y_\rho$ . Moreover, it follows by the very definition of  $U_\rho$  that  $\tilde{U}_\rho(y) = 0$  for  $y \geq L_\rho$ .

**Remark 5.** *The modified utility function  $U_\rho$  is not of class  $C^1$ , however the discussion in [26, Section 6.5] can be used to show that (4.4) also holds for  $v$  and  $\tilde{v}$  replaced by  $v_\rho$  and  $\tilde{v}_\rho$ , i.e.*

$$v_\rho(t, x) = \inf_{y > 0} \{\tilde{v}_\rho(t, y) + xy\}. \quad (4.17)$$

The following large deviations-type argument is needed to estimate the error of this Lipschitz continuous approximation.

**Lemma 4.2.** *Consider an  $\mathbb{R}$ -valued process  $p$  and an  $\mathbb{R}^d$ -valued process  $q$ , both progressively measurable with  $\int_0^T |p_s| ds \leq \mu T$  and  $\int_0^T |q_s|^2 ds \leq \gamma^2 T$  a.s., respectively, for some constants  $\mu, \gamma \geq 0$ , and let*

$$X_t = x \exp \left( \int_0^t p_s ds + \int_0^t q_s dW_s \right)$$

for  $t \in [0, T]$ . Then

$$\mathbb{P}[X_t \geq \rho] \leq 2 \exp \left( -\frac{3}{8\gamma^2 T} (\log \rho/x - \mu T)^2 \right), \quad (4.18)$$

$$\mathbb{P}[X_t \leq c_0/\rho] \leq 2 \exp \left( -\frac{3}{8\gamma^2 T} (\log \rho/(c_0 x) - \mu T)^2 \right). \quad (4.19)$$

Moreover, for each  $p > 0$ ,  $x > 0$  there exists  $C > 0$  such that

$$\mathbb{E}[X_t \mathbb{1}_{\{X_t \geq \rho\}}] \leq C \rho^{-p} \quad (4.20)$$

for all  $t \in [0, T]$ .

*Proof.* We have, for any  $\lambda > 0$ ,

$$\begin{aligned} \mathbb{P}[X_t \geq \rho] &= \mathbb{P}\left[\exp\left(\frac{\lambda}{2}\left(\log X_t/x - \int_0^t p_s ds\right)^2\right) \geq \exp\left(\frac{\lambda}{2}\left(\log \rho/x - \int_0^t p_s ds\right)^2\right)\right] \\ &\leq \mathbb{P}\left[\exp\left(\frac{\lambda}{2}\left(\log X_t/x - \int_0^t p_s ds\right)^2\right) \geq \exp\left(\frac{\lambda}{2}(\log \rho/x - \mu T)^2\right)\right]. \end{aligned}$$

Following the same steps as in the proof of Lemma 2.6 in [15], we obtain for  $\lambda\gamma^2 T < 1$

$$\mathbb{E}\left[\exp\left(\frac{\lambda}{2}\left(\int_0^t q_s dW_s\right)^2\right)\right] \leq \frac{1}{\sqrt{1 - \lambda\gamma^2 T}},$$

and hence from Markov's inequality

$$\mathbb{P}[X_t \geq \rho] \leq \frac{\exp\left(-\frac{\lambda}{2}(\log \rho/x - \mu T)^2\right)}{\sqrt{1 - \lambda\gamma^2 T}}.$$

Choosing  $\lambda = 3/(4\gamma^2 T)$  we obtain (4.18).

The second statement (4.19) follows immediately by replacing  $X_t$  by  $1/X_t$ ,  $x$  by  $1/x$ ,  $(p, q)$  by  $(-p, -q)$ , and  $\rho$  by  $\rho/c_0$ .

Finally, the last estimate is obtained from

$$\mathbb{E}[X_t \mathbb{1}_{\{X_t \geq \rho\}}] \leq \sum_{k=\lfloor \rho \rfloor}^{\infty} (k+1) \mathbb{P}(X_t \in [k, k+1)) = (\lfloor \rho \rfloor + 1) \mathbb{P}(X_t \geq \lfloor \rho \rfloor) + \sum_{k=\lfloor \rho \rfloor}^{\infty} \mathbb{P}(X_t \geq k),$$

and estimating each term by substituting  $\lfloor \rho \rfloor$  and  $k$  into (4.18).  $\square$

Let  $G_\rho^h, I_\rho^h, G_\rho^{h,\Delta x}, I_\rho^{h,\Delta x}$  denote the quantities defined by (4.11) and (4.12) replacing  $V, \widetilde{V}, W, \widetilde{W}$  by  $V_\rho, \widetilde{V}_\rho, W_\rho, \widetilde{W}_\rho$ . We then obtain the following extension of Theorem 4.1 to the general case of non Lipschitz utility functions.

**Theorem 4.3.** *Let assumptions (H1) to (H4) be satisfied. Then, there exist some constants  $C, \widetilde{C} \geq 0$  and  $\delta : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  with  $\delta(x, \rho) = o(\rho^{-p})$  for all  $x$  as  $\rho \rightarrow \infty$  for all  $p > 0$ , such that for any  $n = 0, \dots, N$ ,  $x > 0$*

$$\begin{aligned} -L_\rho C(1 + x^{2M})h^{(M-1)/2M} &\leq v(t_n, x) - V_\rho(t_n, x) \\ &\leq G_\rho^h(t_n, x) + \widetilde{L}_\rho \widetilde{C}(1 + (I_\rho^h(t_n, x))^{2M})h^{(M-1)/2M} + \delta(x, \rho) \end{aligned} \quad (4.21)$$

and for any  $n = 0, \dots, N$ ,  $m \in \mathbb{N}^+$

$$\begin{aligned} -L_\rho C(1 + x_m^{2M})\left(h^{(M-1)/2M} + \Delta x/h\right) &\leq v(t_n, x_m) - W_\rho(t_n, x_m) \\ &\leq G_\rho^{h,\Delta x}(t_n, x_m) + \widetilde{L}_\rho \widetilde{C}(1 + (I_\rho^{h,\Delta x}(t_n, x_m))^{2M})\left(h^{(M-1)/2M} + \Delta x/h\right) + \delta(x, \rho). \end{aligned} \quad (4.22)$$

*Proof.* Let us consider for simplicity the semi discrete case. The lower bounds follow by (4.16) applying Proposition 3.1 to the value function  $v_\rho$ . As  $A$  is bounded, the definition of  $X^{t,x,\alpha}$  from (2.1) satisfies the assumptions on the coefficients in Lemma 4.2. We therefore get immediately

$$\begin{aligned} 0 \leq v(t, x) - v_\rho(t, x) &= \sup_{\alpha \in \mathcal{A}} \mathbb{E}_{t,x}[U(X_T^\alpha)] - \sup_{\alpha \in \mathcal{A}} \mathbb{E}_{t,x}[U_\rho(X_T^\alpha)] \\ &\leq \sup_{\alpha \in \mathcal{A}} \mathbb{E}_{t,x}[U(X_T^\alpha) - U_\rho(X_T^\alpha)] \\ &\leq U(c_0/\rho) \sup_{\alpha \in \mathcal{A}} \mathbb{P}[X_T^{t,x,\alpha} \leq c_0/\rho] + U'(\rho) \sup_{\alpha \in \mathcal{A}} \mathbb{E}_{t,x}[(X_T^\alpha - \rho) \mathbb{1}_{\{X_T^\alpha \geq \rho\}}] \\ &= o(\rho^{-p}) \end{aligned}$$

for all  $p$  and all  $x$ . Thanks to the duality property (4.17), one has

$$v(t_n, x) \leq \inf_{y>0} \{\tilde{v}_\rho(t_n, y) + xy\} + \delta(x, \rho).$$

Applying Proposition 4.2, one has

$$\tilde{v}_\rho(t_n, y) \leq \tilde{V}_\rho(t_n, y) + xy + \tilde{L}_\rho \tilde{C}(1 + y^{2M})h^{(M-1)/2M},$$

so that arguing as in the proof of Theorem 4.1 we get the upper bounds.  $\square$

**Remark 6.** *The previous error estimates clearly depend on the parameter  $\rho$  and the utility function  $U$  via the Lipschitz constant  $L_\rho$ . In the case of power utility, we have  $L_\rho = x_\rho^{p-1}/p = c_0^{p-1}/p \rho^{1-p}$ . As  $\delta(x, \rho)$  goes to zero faster than any power of  $1/\rho$ , we can choose  $\rho = h^{-r}$  for arbitrarily small positive  $r$  and therefore obtain an order of  $L_\rho h^{(M-1)/2M} + \delta(x, \rho)$  arbitrarily close to the Lipschitz case, i.e.  $(M-1)/2M$ .*

**Remark 7.** *The above result can also be extended to cases where  $\lim_{x \downarrow 0} U(x) = -\infty$ , by considering  $U_\rho(x) = U(x_\rho) + U'(x_\rho)(x - x_\rho)$  for  $x \in [0, x_\rho]$ . Then we can estimate  $\mathbb{E}[(U_\rho(X_t) - U(X_t))\mathbb{1}_{\{X_t \geq \rho\}}]$  similar to the proof of Lemma 4.2, as long as  $U$  does not grow more than, e.g., exponentially in  $-1/x$  as  $x \rightarrow 0$ . This is in particular satisfied by the commonly used log-utility.*

## 5. NUMERICAL TESTS

We test our theoretical results on some concrete examples numerically. We consider  $d = 1$  and the computational domain  $[0, x_{\max}]$ . We denote by  $N$  and  $J$  the number of time and space steps, respectively, i.e.

$$h = \frac{T}{N} \quad \text{and} \quad \Delta x = \frac{x_{\max}}{J}.$$

We study the case of a power utility function:

$$U(x) = \frac{x^p}{p} \quad \text{for some } p \in (0, 1). \quad (5.1)$$

We consider the modification  $U_\rho$  of the utility function obtained in (4.15), for  $\rho = 18$  and  $c_0 = 8$ . The utility function  $U$  for  $p = 0.5$  and its conjugate  $\tilde{U}$ , as well as its Lipschitz continuous approximation  $U_\rho$  and its conjugate  $\tilde{U}_\rho$  are shown in Figure 2.

In our tests, we take  $M = 4$  with  $\Delta x \sim h^{11/8}$  obtained from (4.22) balancing the error terms, more specifically  $J \sim \lceil N^{11/8} \rceil$ . Taking  $M > 2$  has only (theoretical) advantages for non-smooth solutions, while we would observe order of convergence at most one for any choice of  $M \geq 2$ , even in the smooth case. This is due to the fact that, even in the case of smooth solutions, the use of the Euler-Maruyama scheme reduces the order of consistency of the overall scheme to one (noting that a modified proof utilising the higher weak order 1 of the Euler-Maruyama scheme, compared to the strong order  $1/2$ , can be used in the smooth case), regardless of the value of  $M$ . An improvement of the order of consistency might be achieved by combining higher values of  $M$  with the use of higher order time-stepping schemes, for instance the higher order Taylor schemes of [17].

For the optimization over the controls in our computations, we truncate  $A$  and  $\Gamma$  first to a finite interval, if necessary, and then discretise the interval by  $N_a$  and  $N_\gamma$  equally spaced mesh points, respectively. As already pointed out in Section 4.2 and Remark 4, this further approximation decreases the value of the discrete primal (maximisation) problem and increases the value of the discrete dual (minimisation) problem, in the same way as the piecewise constant (in time) control approximation does. This implies that this component of the error is captured in the numerical duality gap which we compute *a posteriori*. The approximation can generally only be improved by increasing the size of the control intervals and decreasing the control mesh spacing, concurrently with decreasing  $h$  and  $\Delta x$ .

As the optimal control in our examples is bounded, the error of the control truncation is zero if the interval is chosen large enough. It is seen from the computations that the contribution of

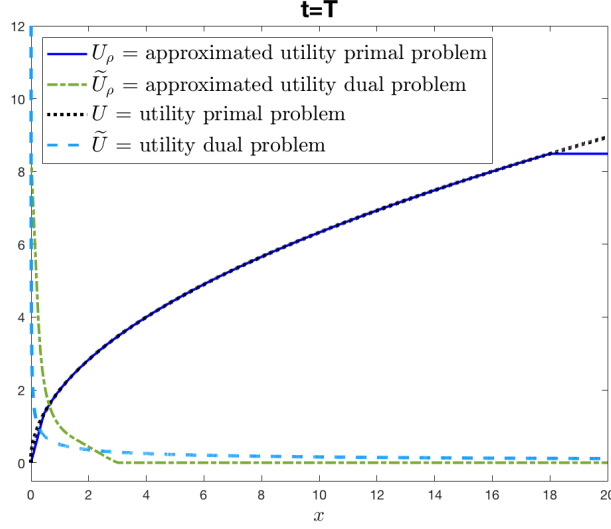


FIGURE 2. The power utility function  $U$  (in dotted black) with its conjugate  $\tilde{U}$  (in dashed cyan) together with the Lipschitz continuous approximation  $U_\rho$  (in solid blue) and its conjugate  $\tilde{U}_\rho$  (in dash-dotted green). Here,  $x_{\max} = 20$ ,  $\rho = 18$  and  $c_0 = 8$ .

the control discretisation error is small, decreasing quadratically in  $N_a^{-1}$  and  $N_\gamma^{-1}$  since we have a smooth dependence of the Hamiltonian on the control. In our tests, we take  $N_a \sim N_\gamma \sim N$ , such that the control discretisation error becomes eventually negligible.

As the point  $x_m$  approaches 0 or  $x_{\max}$ , it may happen that  $\hat{X}_{t_{n+1}}^{t_n, x_m}$  oversteps the domain  $(0, x_{\max})$ . In this case, we use linear extrapolation in order to define  $W_\rho$  and  $\tilde{W}_\rho$  outside the computational mesh. More precisely, one can verify that, due to the boundedness of the control and coefficients, the process  $X_\cdot$  from (2.1) never reaches 0 for  $x > 0$  and equation (2.4) holds up to the left boundary. From equation (3.10), it is clear that for  $m = 0$ , the argument of the expression on the right-hand side is  $(t_{n+1}, 0)$ , so that  $W_\rho(t_n, 0) = W_\rho(t_{n+1}, 0)$  for all  $n$ , at the boundary point. For  $m > 0$  and  $h$  small enough, the argument is  $(t_{n+1}, x)$  for some  $x > 0$ . If  $x < x_{\max}$ , i.e.  $x$  in some interval  $(x_k, x_{k+1}]$ ,  $k \geq 0$ , in the interior of the domain, the value can be obtained by linear interpolation from  $W_\rho(t_{n+1}, x_k)$  and  $W_\rho(t_{n+1}, x_{k+1})$ . In the rare case that  $x < 0$  (for larger  $h$ ) we extend  $W_\rho(t_{n+1}, \cdot)$  in (3.10) by linear extrapolation from  $[x_0, x_1]$  to negative  $x$ . As this is only needed for  $h$  above a certain threshold, it does not affect our estimates. In the case  $x > x_{\max} > \rho$  (where we choose  $x_{\max}$  and  $\rho$  so that the second inequality holds), we can set  $W_\rho(t_{n+1}, x) = U_\rho(\rho)$ , where we have exactly  $v_\rho(t_{n+1}, x) = U_\rho(\rho)$  if  $x_{\max}$  is large enough because of the constancy of the solution for large  $x$ . A similar argument holds for the dual problem.

**Test 1: Merton problem.** We first study the classical Merton problem. This corresponds to the dynamics (2.1) with  $g \equiv 0$ , constant coefficients  $b, r, \sigma$  and  $A = \mathbb{R}$ . It is well known that for this problem there exists a closed-form solution given by (see, e.g. [24])

$$v(t, x) = \exp \left\{ t \left( a^*(b - r) + r - \frac{1}{2} (a^*)^2 (1 - p) \sigma^2 \right) \right\} U(x),$$

where  $U$  and  $p$  are given in (5.1), and

$$a^* := \frac{(b - r)}{\sigma^2(1 - p)}$$

is the optimal control. We recall that in this case the dual problem is linear and no optimisation is necessary since  $\Gamma = \{0\}$ . The values of the coefficients used in the test is given in Table 1. For these values, setting  $A = [-1, 1]$  is sufficient to have  $a^* \in A$ .

| $p$ | $r$ | $b$ | $\sigma$ | $T$ | $x_{\max}$ |
|-----|-----|-----|----------|-----|------------|
| 0.5 | 0.8 | 1.2 | 1        | 0.5 | 20         |

TABLE 1. Test 1: Parameters used in numerical experiments.

Table 2 reports the error and the estimated convergence rate of  $W_\rho$  to the exact solution  $v$  of the primal problem. As expected, the order of convergence is around 1. It is important to notice that continuing to refine the mesh without increasing  $\rho$ , we cannot get convergence to  $v$ . In fact, the probability in (4.18) and (4.19), even if small at points  $x$  far from the boundaries of the domain, is different from zero everywhere (see also Figure 3, left). To reduce the contribution to the error coming from the term in  $\rho$  we compute the error locally, away from the boundary of the computational domain.

In Table 3, we report the numerical duality gap, i.e. the quantity  $G_\rho^{h,\Delta x}(T, x)$ . This quantity also decreases with order 1 or even slightly higher. In this case, the duality gap is bigger than the error, but of the same order. In Figure 3 (right) we show the numerical solutions  $W_\rho$  and  $\widetilde{W}_\rho$  of the primal and the dual problem, together with the convex conjugate of  $\widetilde{W}_\rho$ .

| $J$  | $N$ | Error $L^1$ | Order $L^1$ | Error $L^2$ | Order $L^2$ | Error $L^\infty$ | Order $L^\infty$ | CPU (s) |
|------|-----|-------------|-------------|-------------|-------------|------------------|------------------|---------|
| 18   | 8   | 1.96E-01    | -           | 1.86E-01    | -           | 1.77E-01         | -                | 0.30    |
| 46   | 16  | 1.44E-01    | 0.44        | 1.12E-01    | 0.74        | 1.05E-01         | 0.75             | 1.05    |
| 118  | 32  | 5.85E-02    | 1.30        | 4.54E-02    | 1.30        | 5.86E-02         | 0.84             | 3.91    |
| 305  | 64  | 1.52E-02    | 1.94        | 1.14E-02    | 2.00        | 1.52E-02         | 1.95             | 15.54   |
| 790  | 128 | 5.70E-03    | 1.42        | 4.11E-03    | 1.47        | 4.76E-03         | 1.67             | 61.95   |
| 2048 | 256 | 2.35E-03    | 1.28        | 1.68E-03    | 1.29        | 1.74E-03         | 1.45             | 467.54  |
| 5312 | 512 | 1.12E-03    | 1.07        | 8.14E-04    | 1.04        | 9.18E-04         | 0.92             | 2169.45 |

TABLE 2. Test 1: Local ( $x \in [1, 2]$ ) errors and convergence order comparing  $W_\rho$  with the exact solution  $v$ , for  $M = 4$  (Gauß-Hermite quadrature points),  $N = 4 \cdot 2^k$  (time steps),  $J = \lceil N^{11/8} \rceil$  (space steps),  $N_a = 2^k + 1$  (discrete controls), for  $k = 1, 2, \dots, 8$ .

From the results in Table 3 we deduce that (given the choice of  $\Delta x$  in relation to  $h$ )

$$|G_\rho^{h,\Delta x}(t, x)| \leq C \left( h + \Delta x^{8/11} \right),$$

which, combined with (4.14) and taking  $M = 4$ , gives the *a posteriori* bounds

$$-C \left( h^{3/8} + \Delta x^{3/11} \right) \leq v(t_n, x_m) - W_\rho(t_n, x_m) \leq C \left( h + \Delta x^{8/11} + h^{3/8} + \Delta x^{3/11} \right), \quad (5.2)$$

which in conclusion is a symmetric bound of order 3/8 in time and 3/11 in space.

For using our error estimates, it is necessary to solve numerically both the primal and the dual problem. The computational cost for the solution of the dual problem is comparable to that of the primal one, which has the same structure and uses the same scheme. This can be partially observed comparing the CPU times in Table 2 and 3 (however, in this case the dual problem is linear and the computational cost is less than double that of the primal one).

We illustrate the different contributions to the error, together with the actual error, in Figure 4. The figure shows the order (at least) one for the empirical error and for the numerical duality gap, as one would have expected from the first order error of the scheme for sufficiently smooth

| $J$   | $N$  | Gap $L^1$ | Order $L^1$ | Gap $L^2$ | Order $L^2$ | Gap $L^\infty$ | Order $L^\infty$ | CPU (s)  |
|-------|------|-----------|-------------|-----------|-------------|----------------|------------------|----------|
| 18    | 8    | 2.17E+01  | -           | 7.17E+00  | -           | 3.22E+00       | -                | 0.56     |
| 46    | 16   | 1.24E+01  | 0.80        | 4.04E+00  | 0.83        | 1.65E+00       | 0.96             | 1.41     |
| 118   | 32   | 7.24E+00  | 0.78        | 2.31E+00  | 0.80        | 9.24E-01       | 0.88             | 4.70     |
| 305   | 64   | 3.92E+00  | 0.89        | 1.26E+00  | 0.88        | 5.06E-01       | 0.87             | 17.98    |
| 790   | 128  | 1.87E+00  | 1.07        | 6.03E-01  | 1.06        | 2.43E-01       | 1.06             | 110.56   |
| 2048  | 256  | 7.16E-01  | 1.38        | 2.37E-01  | 1.35        | 1.00E-01       | 1.28             | 656.69   |
| 5312  | 512  | 1.72E-01  | 2.05        | 5.53E-02  | 2.10        | 2.20E-02       | 2.19             | 2813.47  |
| 13778 | 1024 | 5.97E-02  | 1.53        | 1.94E-02  | 1.51        | 8.05E-03       | 1.45             | 17059.00 |

TABLE 3. Test 1: Global ( $x \in [0, x_{\max}]$ ) duality gap  $G_\rho^{h, \Delta x}$  from (4.12) and related convergence order, for  $M = 4$  (Gauß-Hermite quadrature points),  $N = 4 \cdot 2^k$  (time steps),  $J = \lceil N^{11/8} \rceil$  (space steps),  $N_a = 2^k + 1$  (discrete controls), for  $k = 1, 2, \dots, 8$ .

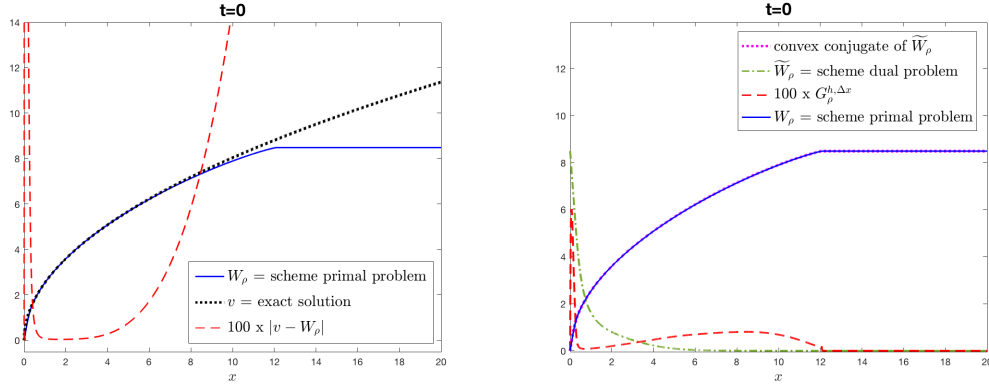


FIGURE 3. Test 1: Numerical solution  $W_\rho$  (in solid blue) compared with the exact solution (dotted black, left) and the convex conjugate of  $\widetilde{W}_\rho$  (dotted magenta, right). The dashed red line represents the error (left) and the numerical duality gap (right), multiplied by a factor 100. The dash-dotted green line on the right is the numerical approximation of the dual problem  $\widetilde{W}_\rho$ .

solutions. We also plot the theoretical error bounds, which hold in the general non-smooth case, for the Euler-Maruyama scheme, given by the expression (A.1) in the Appendix, of order  $1/2$ , and for the Gauß-Hermite approximation, from (A.2), of order  $3/8$ . The big constants appearing in the theoretical *a priori* bounds, which are not sharp, put the magnitude of these theoretical errors far from that of the empirical one.

For this problem, the optimal control is constant over time, so there is no error coming from the piecewise control approximation and theoretical bounds as those provided by (A.1) and (A.2) can be used for both the upper and lower bound. The numerical duality gap in this case contains the sum of the numerical approximation errors for the primal and the dual problem as well as the error coming from the approximation in  $\rho$  and the computation of the numerical convex conjugate.

**Test 2: Cuoco and Liu example.** This example is taken from [9]. In this paper, the authors consider the nonlinear dynamics in (2.1) (i.e.  $g \neq 0$ ) and portfolio constraints (i.e.  $A \subsetneq \mathbb{R}$ ). We still consider a power utility and  $d = 1$ . Let  $A$  be defined by

$$A = \left\{ a \in \mathbb{R} : \max(0, -a)\lambda_- + \max(0, a)\lambda_+ \leq 1 \right\}$$



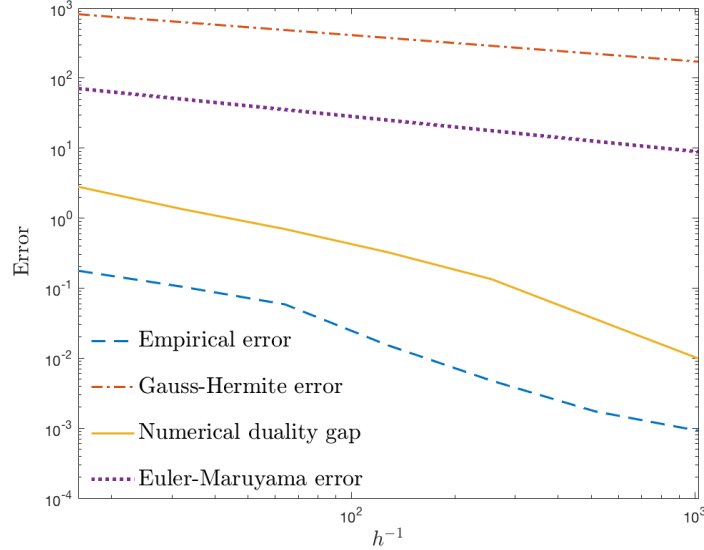


FIGURE 4. Test 1. Local ( $x \in [1, 2]$ ) empirical error  $|v(0, x) - W_\rho(0, x)|$  as reported in Table 3, global numerical duality gap  $G_\rho^{h, \Delta x}$  reported in Table 2, theoretical error estimate for the Euler-Maruyama and Gauß-Hermite approximation given by (A.1) and (A.2), respectively.

| $p$ | $r$ | $R$ | $b$ | $\sigma$ | $T$ | $x_{\max}$ | $\iota$ | $\lambda_+$ | $\lambda_-$ |
|-----|-----|-----|-----|----------|-----|------------|---------|-------------|-------------|
| 0.5 | 0.8 | 1   | 1.2 | 0.5      | 0.5 | 20         | 0.5     | 1           | 1           |

TABLE 4. Test 2: Parameters used in numerical experiments.

for some  $\lambda_- \geq 0$  and  $\lambda_+ \in [0, 1]$ . The function  $g$  is defined by

$$g(a) = -r(1 + \iota\lambda_-) \max(0, -a) - (R - r)(1 - \max(0, a) - \iota\lambda_- \max(0, -a)),$$

where  $R \geq r$  and  $\iota \in [0, 1]$ . The values used in our numerical simulation are reported in Table 4.

Observe that the choice  $\lambda_+ = \lambda_- = 1$  corresponds to  $A = [-1, 1]$ . In order to define  $\Gamma$ , we use the explicit expression given in [9, Section 5.2] for the optimal control. For the data in Table 4, we can take  $\Gamma = [-1, 1]$  to guarantee  $\nu_t^* \in \Gamma$  for any  $t \in [0, T]$ . Table 5 reports the numerical duality gap and the corresponding convergence order. The numerical solutions  $W_\rho$  and  $\widetilde{W}_\rho$  of the primal and the dual problem, together with the convex conjugate of  $\widetilde{W}_\rho$  are shown in Figure 5.

The results in Table 5 give once again an estimate of the form

$$|G_\rho^{h, \Delta x}(t, x)| \leq C \left( h + \Delta x^{8/11} \right)$$

for the duality gap, leading to the *a posteriori* bounds (5.2).

## 6. CONCLUSION AND PERSPECTIVES

For a suitable class of convex optimal control problems, we obtained in this paper *a posteriori* error bounds using the numerical approximation of a dual problem.

Our numerical tests confirm the results given by the theoretical analysis and suggest a convergence to zero with order one of the numerical duality gap. Establishing rigorously a duality relation between the numerical approximations of the primal and the dual problem seems to us an interesting direction of research that we would like to pursue. Beyond the independent theoretical interest,

| $J$  | $N$ | Gap $L^1$ | Order $L^1$ | Gap $L^2$ | Order $L^2$ | Gap $L^\infty$ | Order $L^\infty$ | CPU (s)  |
|------|-----|-----------|-------------|-----------|-------------|----------------|------------------|----------|
| 18   | 8   | 2.26E+01  | -           | 7.44E+00  | -           | 3.59E+00       | -                | 0.79     |
| 46   | 16  | 1.09E+01  | 1.05        | 3.48E+00  | 1.10        | 1.47E+00       | 1.29             | 2.51     |
| 118  | 32  | 5.59E+00  | 0.96        | 1.74E+00  | 1.00        | 6.87E-01       | 1.10             | 9.83     |
| 305  | 64  | 2.82E+00  | 0.99        | 8.79E-01  | 0.99        | 3.47E-01       | 0.98             | 45.94    |
| 790  | 128 | 1.38E+00  | 1.03        | 4.35E-01  | 1.01        | 1.77E-01       | 0.97             | 552.49   |
| 2048 | 256 | 5.75E-01  | 1.26        | 1.83E-01  | 1.25        | 7.49E-02       | 1.24             | 6305.33  |
| 5312 | 512 | 1.56E-01  | 1.88        | 5.00E-02  | 1.87        | 2.08E-02       | 1.85             | 54006.70 |

TABLE 5. Test 2: Global ( $x \in [0, x_{\max}]$ ) duality gap  $G_\rho^{h, \Delta x}$  from (4.12 and related convergence order, for  $M = 4$  (Gauß-Hermite quadrature points),  $N = 4 \cdot 2^k$  (time steps),  $J = \lceil N^{11/8} \rceil$  (space steps),  $N_a = 2^k + 1$  (discrete controls), for  $k = 1, 2, \dots, 8$ .

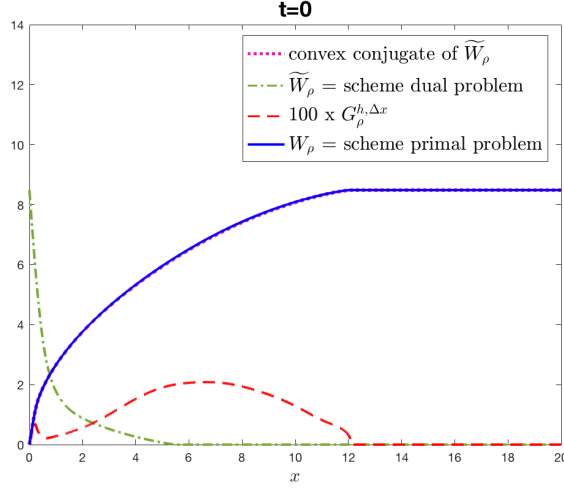


FIGURE 5. Test 2: Numerical solution  $W_\rho$  (in solid blue) compared with the convex conjugate of  $\widetilde{W}_\rho$  (dotted magenta). The dashed red line represents the numerical duality gap multiplied by a factor 100 and the dash-dotted green line the numerical approximation of the dual problem  $\widetilde{W}_\rho$ .

this would also allow us to obtain an *a priori* upper bound for the numerical error. The possibility of improving the order by higher order time stepping is also left for future research.

#### APPENDIX A. EXPLICIT COMPUTATION OF THE CONSTANTS

In this section, we explicitly compute the constant  $C$  which appears in the lower bound of (4.13). Analogous estimates can be used to derive the constant  $\widetilde{C}$  appearing in the upper bound. In what follows we denote for  $t \in [0, T]$ ,  $a \in A$ ,  $x \in \mathbb{R}$ :

$$\mu(t, x, a) := (r(t) + a^\top \cdot (b(t) - r(t)\mathbb{1}) + g(t, a)) x, \quad \psi(t, x, a) := a^\top \sigma(t)x.$$

Let  $C_\mu, C_\psi \geq 0$  such that for  $t, s \in [0, T]$ ,  $a \in A$ ,  $x, y \in \mathbb{R}$ :

$$|\mu(t, x, a) - \mu(s, y, a)| \leq C_\mu \left( |x - y| + (1 + |x|)|t - s|^{1/2} \right),$$

$$|\psi(t, x, a) - \psi(s, y, a)| \leq C_\psi \left( |x - y| + (1 + |x|)|t - s|^{1/2} \right)$$

and

$$|\mu(t, x, a)| \leq C_\mu(1 + |x|), \quad |\psi(t, x, a)| \leq C_\psi(1 + |x|).$$

**A.1. Explicit bounds for the Euler-Maruyama approximation.** We consider the Euler-Maruyama approximation given by (3.4) for  $\alpha \equiv (a_0, \dots, a_{N-1}) \in \mathcal{A}_h$ . This leads to the following expression for  $\bar{X}^{t_n, x, \alpha}$ :

$$\bar{X}_{t_k}^{t_n, x, \alpha} = x + \sum_{i=n}^{k-1} \int_{t_i}^{t_{i+1}} \mu(t_i, \bar{X}_{t_i}^{t_n, x, \alpha}, a_i) ds + \int_{t_i}^{t_{i+1}} \psi(t_i, \bar{X}_{t_i}^{t_n, x, \alpha}, a_i) dW_s.$$

Moreover, by the very definition of  $X^{t_n, x, \alpha}$ :

$$X_{t_k}^{t_n, x, \alpha} = x + \sum_{i=n}^{k-1} \int_{t_i}^{t_{i+1}} \mu(s, X_s^{t_n, x, \alpha}, a_i) ds + \int_{t_i}^{t_{i+1}} \psi(s, X_s^{t_n, x, \alpha}, a_i) dW_s.$$

Therefore, using the Cauchy-Schwartz inequality and Itô isometry together with classical estimates, one has

$$\begin{aligned} \mathbb{E} \left[ |\bar{X}_{t_k}^{t_n, x, \alpha} - X_{t_k}^{t_n, x, \alpha}|^2 \right] &\leq 2T \sum_{i=n}^{k-1} \mathbb{E} \left[ \int_{t_i}^{t_{i+1}} \left| \mu(t_i, \bar{X}_{t_i}^{t_n, x, \alpha}, a_i) - \mu(s, X_s^{t_n, x, \alpha}, a_i) \right|^2 ds \right] \\ &\quad + 2 \sum_{i=n}^{k-1} \mathbb{E} \left[ \int_{t_i}^{t_{i+1}} \left| \psi(t_i, \bar{X}_{t_i}^{t_n, x, \alpha}, a_i) - \psi(s, X_s^{t_n, x, \alpha}, a_i) \right|^2 ds \right] \\ &\leq 8K_1 h \sum_{i=n}^{k-1} \left( \mathbb{E} \left[ |\bar{X}_{t_i}^{t_n, x, \alpha} - X_{t_i}^{t_n, x, \alpha}|^2 \right] + h + h \mathbb{E} \left[ \sup_{s \in [t_i, t_{i+1}]} |X_s^{t_n, x, \alpha}|^2 \right] \right. \\ &\quad \left. + \mathbb{E} \left[ \sup_{s \in [t_i, t_{i+1}]} |X_s^{t_n, x, \alpha} - X_{t_i}^{t_n, x, \alpha}|^2 \right] \right), \end{aligned}$$

where we denoted  $K_1 := (C_\mu^2 T + C_\psi^2)$ . By classical estimates on the process  $X^{t_n, x}$  and denoting  $K_2(\xi) := (C_\mu^2 \xi + 4C_\psi^2)$ , one has

$$\begin{aligned} \mathbb{E} \left[ \sup_{s \in [t_i, t_{i+1}]} |X_s^{t_n, x, \alpha}|^2 \right] &\leq 3(|x|^2 + 2K_2(T)) e^{6K_2(T)T}, \\ \mathbb{E} \left[ \sup_{s \in [t_i, t_{i+1}]} |X_s^{t_n, x, \alpha} - X_{t_i}^{t_n, x, \alpha}|^2 \right] &\leq 4K_2(h)h \left( 1 + 3(|x|^2 + 4K_2(T)) e^{6K_2(T)T} \right). \end{aligned}$$

Putting these estimates together:

$$\mathbb{E} \left[ |\bar{X}_{t_k}^{t_n, x, \alpha} - X_{t_k}^{t_n, x, \alpha}|^2 \right] \leq 8K_1 h \sum_{i=n}^{k-1} \mathbb{E} \left[ |\bar{X}_{t_i}^{t_n, x, \alpha} - X_{t_i}^{t_n, x, \alpha}|^2 \right] + 8K_1 T h (1 + 2K_2(h))(1 + K_3(x))$$

with  $K_3(x) := 3(|x|^2 + 2K_2(T)T) e^{6K_2(T)T}$ , so that, using Gronwall's lemma, one obtains

$$\begin{aligned} \mathbb{E} \left[ |\bar{X}_{t_k}^{t_n, x, \alpha} - X_{t_k}^{t_n, x, \alpha}|^2 \right] &\leq 8K_1 h (1 + 2K_2(h))(1 + K_3(x)) \left( 1 + e^{\sum_{i=n}^{k-1} 8K_1 h} \left( \sum_{i=n}^{k-1} 8K_1 h \right) \right) \\ &\leq 8K_1 h (1 + 2K_2(h))(1 + K_3(x)) (1 + 8K_1 T e^{8K_1 T}). \end{aligned}$$

Using the Lipschitz continuity of  $U$ , one has

$$\left| \sup_{\alpha \in \mathcal{A}_h} \mathbb{E} \left[ U(\bar{X}_T^{t_n, x, \alpha}) \right] - \sup_{\alpha \in \mathcal{A}_h} \mathbb{E} \left[ U(X_T^{t_n, x, \alpha}) \right] \right| \leq L \sup_{\alpha \in \mathcal{A}_h} \mathbb{E} \left[ |\bar{X}_T^{t_n, x, \alpha} - X_T^{t_n, x, \alpha}| \right].$$

In conclusion, the contribution to the error coming from the Euler-Maruyama approximation can be bounded by

$$L \left( 8K_1(1 + 2K_2(h))(1 + K_3(x))(1 + 8K_1Te^{8K_1T}) \right)^{1/2} h^{1/2}.$$

For a linear (in the state), time independent dynamics as the one considered in Section 5, one simply has

$$|\mu(x, a) - \mu(y, a)| \leq C_\mu |x - y|, \quad |\psi(x, a) - \psi(y, a)| \leq C_\psi |x - y|$$

and

$$|\mu(x, a)| \leq C_\mu |x|, \quad |\psi(x, a)| \leq C_\psi |x|.$$

It is possible to verify that this leads to

$$\begin{aligned} \mathbb{E} \left[ |\bar{X}_{t_k}^{t_n, x, \alpha} - X_{t_k}^{t_n, x, \alpha}|^2 \right] &\leq 4K_1 h \sum_{i=n}^{k-1} \left( \mathbb{E} \left[ |\bar{X}_{t_i}^{t_n, x, \alpha} - X_{t_i}^{t_n, x, \alpha}|^2 \right] + \mathbb{E} \left[ \sup_{s \in [t_i, t_{i+1}]} |X_s^{t_n, x, \alpha} - X_{t_i}^{t_n, x, \alpha}|^2 \right] \right) \\ &\leq 4K_1 h \sum_{i=n}^{k-1} \left( \mathbb{E} \left[ |\bar{X}_{t_i}^{t_n, x, \alpha} - X_{t_i}^{t_n, x, \alpha}|^2 \right] + 2K_2(h) h \mathbb{E} \left[ \sup_{s \in [t_i, t_{i+1}]} |X_s^{t_n, x, \alpha}|^2 \right] \right) \end{aligned}$$

with

$$\mathbb{E} \left[ \sup_{s \in [t_i, t_{i+1}]} |X_s^{t_n, x, \alpha}|^2 \right] \leq 3|x|^2 e^{3K_2(T)T}.$$

Neglecting the infinitesimal terms, one has

$$\mathbb{E} \left[ |\bar{X}_{t_k}^{t_n, x, \alpha} - X_{t_k}^{t_n, x, \alpha}|^2 \right] \leq 4K_1 h \left( \sum_{i=n}^{k-1} \mathbb{E} \left[ |\bar{X}_{t_i}^{t_n, x, \alpha} - X_{t_i}^{t_n, x, \alpha}|^2 \right] + 24TC_\psi^2 |x|^2 e^{3K_2(T)T} \right)$$

which leads to the sharper error estimate

$$L \left( 96K_1TC_\psi^2 |x|^2 e^{3K_2(T)T} (1 + 4K_1Te^{4K_1T}) \right)^{1/2} h^{1/2}.$$

In the estimates plotted in Section 5, we consider

$$L \left( 24K_1TC_\psi^2 |x|^2 (1 + 4K_1Te^{4K_1T}) \right)^{1/2} h^{1/2} \quad (\text{A.1})$$

since we can approximate the second order moment of  $X$  by  $x^2$  for a local error.

**A.2. Explicit bounds for the Gauß-Hermite approximation.** We consider the case of a one-dimensional Brownian motion. Given a function  $f \in C^{2M}(\mathbb{R})$ , the analysis in [25, Proposition 3.2] shows that

$$\begin{aligned} &\left| \mathbb{E}_{t_n, x} \left[ f(\bar{X}_{t_{n+1}}^a) \right] - \mathbb{E}_{t_n, x} \left[ f(\hat{X}_{t_{n+1}}^a) \right] \right| \\ &\leq \left| \int_{-\infty}^{+\infty} \frac{f^{(2M)}(\tilde{z})}{(2M)!} (\sqrt{2h}\psi(t_n, x, a)y)^{2M} \frac{e^{-y^2}}{\sqrt{\pi}} dy - \sum_{i=1}^M \lambda_i \frac{f^{(2M)}(\tilde{z})}{(2M)!} (\sqrt{h}\psi(t_n, x, a)\xi_i)^{2M} \right| \\ &\leq 2\|f^{2M}\|_\infty \frac{(2h)^M}{2M!} (\psi(t_n, x, a))^{2M} \int_{-\infty}^{\infty} y^{2M} \frac{e^{-y^2}}{\sqrt{\pi}} dy \\ &\quad + \|f^{2M}\|_\infty \frac{h^M}{2M!} (\psi(t_n, x, a))^{2M} \left| 2^M \int_{-\infty}^{\infty} y^{2M} \frac{e^{-y^2}}{\sqrt{\pi}} dy - \sum_{i=1}^M \lambda_i \xi_i^{2M} \right| \\ &\leq \|f^{2M}\|_\infty \frac{h^M}{2M!} C_\psi^{2M} (1 + |x|)^{2M} \left( 2(2M-1)!! + \left| (2M-1)!! - \sum_{i=1}^M \lambda_i \xi_i^{2M} \right| \right), \end{aligned}$$

where in the last inequality we have used that

$$2^M \int_{-\infty}^{\infty} y^{2M} \frac{e^{-y^2}}{\sqrt{\pi}} dy = (2M-1)!!$$

The estimate above corresponds to the error associated with the Gauß-Hermite approximation at each time step, i.e. considering the error at time  $t_{n+1}$  starting from  $t_n$ . Our scheme being iterative in time, the overall contribution to the error will be

$$\left\| \frac{\partial f}{\partial x^{2M}} \right\|_{\infty} \frac{h^{M-1}}{2M!} 2^{2M-1} C_{\psi}^{2M} \left( (2M-1)!! + \left| (2M-1)!! - \sum_{i=1}^M \frac{\omega_i}{\sqrt{\pi}} z_i^{2M} \right| \right) \left( 1 + \sup_{\substack{\alpha \in \mathcal{A}_h \\ k=n \dots N}} \mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_k}^{\alpha})^{2M} \right] \right),$$

where we also used the classical inequality  $|a+b|^{2M} \leq 2^{2M-1}(|a|^{2M} + |b|^{2M})$ . It remains to estimate  $\mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_k}^{\alpha})^{2M} \right]$ . By the recursive definition of  $\hat{X}$ , one has for any  $k = n, \dots, N$

$$\begin{aligned} \mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_{k+1}}^{\alpha})^{2M} \right] &= \mathbb{E}_{t_n, x} \left[ \left( \hat{X}_{t_k}^{\alpha} + h\mu(t_k, \hat{X}_k, a_k) + \sqrt{h}\psi(t_k, \hat{X}_k, a_k)\zeta_k \right)^{2M} \right] \\ &= \mathbb{E}_{t_n, x} \left[ \sum_{i=0}^{2M} \sum_{j=0}^i \binom{2M}{i} \binom{i}{j} h^{i-j} (\hat{X}_{t_k}^{\alpha})^j (\mu(t_k, \hat{X}_k, a_k))^{i-j} \left( \sqrt{h}\psi(t_k, \hat{X}_k, a_k)\zeta_k \right)^{2M-i} \right] \\ &= \mathbb{E}_{t_n, x} \left[ \sum_{i=0}^M \sum_{j=0}^{2i} \binom{2M}{2i} \binom{2i}{j} h^{2i-j} (\hat{X}_{t_k}^{\alpha})^j (\mu(t_k, \hat{X}_k, a_k))^{2i-j} \left( \sqrt{h}\psi(t_k, \hat{X}_k, a_k)\zeta_k \right)^{2M-2i} \right], \end{aligned}$$

where the last equality follows observing that  $\mathbb{E}[(\dots)\zeta_k^{2j+1}] = 0$  for  $j = 0, \dots, M-1$  for any quantity, represented by “ $(\dots)$ ”, independent of  $\zeta_k$ . Therefore, thanks to the linear growth of  $\mu$  and  $\psi$  (taking for simplicity  $C_1 := \max(C_{\mu}, C_{\psi})$ ):

$$\begin{aligned} \mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_{k+1}}^{\alpha})^{2M} \right] &= \mathbb{E}_{t_n, x} \left[ \sum_{i=0}^M \sum_{j=0}^{2i} \binom{2M}{2i} \binom{2i}{j} h^{M+i-j} C_1^{2M-j} (\hat{X}_{t_k}^{\alpha})^j (1 + |\hat{X}_{t_k}^{\alpha}|)^{2i-j} \zeta_k^{2M-2i} \right] \\ &\leq \sum_{i=0}^M \sum_{j=0}^{2i} \binom{2M}{2i} \binom{2i}{j} h^{M+i-j} C_1^{2M-j} \mathbb{E}_{t_n, x} \left[ |\hat{X}_{t_k}^{\alpha}|^j (1 + |\hat{X}_{t_k}^{\alpha}|)^{2i-j} \right] \mathbb{E} [\zeta_k^{2M-2i}] \\ &= \mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_k}^{\alpha})^{2M} \right] + \sum_{j=0}^{2M-1} \binom{2M}{j} h^{2M-j} C_1^{2M-j} \mathbb{E}_{t_n, x} \left[ |\hat{X}_{t_k}^{\alpha}|^j (1 + |\hat{X}_{t_k}^{\alpha}|)^{2M-j} \right] \\ &\quad + \sum_{i=0}^{M-1} \sum_{j=0}^{2i} \binom{2M}{2i} \binom{2i}{j} h^{M+i-j} C_1^{2M-j} \mathbb{E}_{t_n, x} \left[ |\hat{X}_{t_k}^{\alpha}|^j (1 + |\hat{X}_{t_k}^{\alpha}|)^{2i-j} \right] \mathbb{E} [\zeta_k^{2M-2i}]. \end{aligned}$$

For  $0 \leq i \leq M$  and  $0 \leq j \leq 2i$ , one has

$$\mathbb{E}_{t_n, x} \left[ |\hat{X}_{t_k}^{\alpha}|^j (1 + |\hat{X}_{t_k}^{\alpha}|)^{2i-j} \right] \leq 2^{2i} \left( 1 + \mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_k}^{\alpha})^{2M} \right] \right).$$

This gives:

$$\begin{aligned} \mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_{k+1}}^{\alpha})^{2M} \right] &\leq \mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_k}^{\alpha})^{2M} \right] + \left( 1 + \mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_k}^{\alpha})^{2M} \right] \right) h \left\{ \sum_{j=0}^{2M-1} \binom{2M}{j} h^{2M-j-1} C_1^{2M-j} 2^{2M} \right. \\ &\quad \left. + \sum_{i=0}^{M-1} \sum_{j=0}^{2i} \binom{2M}{2i} \binom{2i}{j} h^{M+i-j-1} C_1^{2M-j} 2^{2i} \mathbb{E} [\zeta_k^{2M-2i}] \right\}. \end{aligned}$$

Neglecting the infinitesimal terms and denoting

$$K_4 := 2MC_1 2^{2M} + \frac{2M(2M-1)}{2} C_1^2 2^{2M-2},$$

we have

$$\mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_{k+1}}^\alpha)^{2M} \right] \leq (1 + K_4 h) \mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_k}^\alpha)^{2M} \right] + K_4 h.$$

Iterating, this leads to

$$\mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_{k+1}}^\alpha)^{2M} \right] \leq (1 + K_4 h)^k x^{2M} + khK_4$$

for any  $n \leq k \leq N-1$ , with  $K_4$  not depending on  $k$  and  $\alpha \in \mathcal{A}_h$ . Therefore, we can conclude that

$$\sup_{\substack{\alpha \in \mathcal{A}_h \\ k=n \dots N}} \mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_k}^\alpha)^{2M} \right] \leq x^{2M} e^{K_4 T} + K_4 T.$$

To avoid an exponential growth in  $M$  of the constants and motivated by the fact that in Section 5 we empirically computed a local error, we can strongly simplify our estimates by approximating

$$\sup_{\substack{\alpha \in \mathcal{A}_h \\ k=n \dots N}} \mathbb{E}_{t_n, x} \left[ (\hat{X}_{t_k}^\alpha)^{2M} \right] \approx x^{2M}.$$

The presence of the  $2M$ -th derivative in the error bound requires to pass by a mollification of the original value function. For a given regularization parameter  $\varepsilon$  and mollified value function  $v_\varepsilon$  it is possible to show that an estimate of the form

$$\left\| \frac{\partial^{2M} v_\varepsilon}{\partial x^{2M}} \right\|_\infty \leq LK_5 \varepsilon^{1-2M}$$

holds with  $K_5 := (3+9K_1 T e^{3K_1 T})^{1/2}$ . The balancing between the Gauß-Hermite and regularization error (the last one giving an extra error term of order  $\varepsilon$ ) leads to the choice of optimal order  $\varepsilon = h^{(M-1)/2M}$ . Therefore, we get

$$LK_5 h^{(M-1)/2M} \frac{2^{2M-1}}{2M!} C_\psi^{2M} \left( (2M-1)!! + \left| (2M-1)!! - \sum_{i=1}^M \frac{\omega_i}{\sqrt{\pi}} z_i^{2M} \right| \right) (1 + x^{2M}). \quad (\text{A.2})$$

## REFERENCES

- [1] G. Barles and E.R. Jakobsen. On the convergence rate of approximation schemes for Hamilton-Jacobi-Bellman equations. *M2AN Math. Model. Numer. Anal.*, 36:33–54, 2002.
- [2] G. Barles and E.R. Jakobsen. Error bounds for monotone approximation schemes for Hamilton-Jacobi-Bellman equations. *SIAM J. Numer. Anal.*, 43(2):540–558, 2005.
- [3] G. Barles and E.R. Jakobsen. Error bounds for monotone approximation schemes for parabolic Hamilton-Jacobi-Bellman equations. *Math. Comput.*, 74(260):1861–1893, 2007.
- [4] G. Barles and P.E. Souganidis. Convergence of approximation schemes for fully nonlinear second order equations. *Asymptotic Anal.*, 4:271–283, 1991.
- [5] W.H. Beyer. *CRC Standard Mathematical Tables*. CRC Press, 28th edition, 1987.
- [6] F. Camilli and M. Falcone. An approximation scheme for the optimal control of diffusion processes. *RAIRO Modél. Math. Anal. Numér.*, 29(1):97–122, 1995.
- [7] M.G. Crandall, H. Ishii, and P.L. Lions. User's guide to viscosity solutions of second order partial differential equations. *Bull. Amer. Math. Soc.*, 27(1):1–67, 1992.
- [8] D. Cuoco and J. Cvitanic. Optimal consumption choices for a 'large' investor. *J. Econ. Dyn. Control*, 22(3):401–436, 1998.
- [9] D. Cuoco and H. Liu. A martingale characterization of consumption choices and hedging costs with margin requirements. *Math. Finance*, 10:355–385, 2000.
- [10] K. Debrabant and E.R. Jakobsen. Semi-Lagrangian schemes for linear and fully non-linear diffusion equations. *Math. Comp.*, 82(283):1433–1462, 2012.
- [11] K. Debrabant and E.R. Jakobsen. Semi-Lagrangian schemes for parabolic equations. In T. Gerstner and P. Kloeden, editors, *Recent Developments in Computational Finance: Foundations, Algorithms and Applications*, pages 279–297. World Scientific, 2013.

- [12] N. El Karoui, S. Peng, and M. C. Quenez. Backward stochastic differential equations in finance. Math. Finance, 7(1):1–71, 1997.
- [13] M. Falcone and R. Ferretti. Semi-Lagrangian Approximation Schemes for Linear and Hamilton-Jacobi Equations, volume 133. SIAM, Philadelphia, 2014.
- [14] F.B. Hildebrand. Introduction to Numerical Analysis. New York: McGraw-Hill, 1956.
- [15] Ying Hu and Shanjian Tang. Existence of solution to scalar BSDEs with  $l \exp\left(\sqrt{\frac{2}{\lambda \log(1+L)}}\right)$ -integrable terminal values. Electron. Commun. Probab., 23, 2018.
- [16] E.R. Jakobsen, A. Picarelli, and C. Reisinger. Improved order 1/4 convergence for piecewise constant policy approximation of stochastic control problems. Electron. Commun. Probab., 24(59):1–10, 2019.
- [17] P.E. Kloeden and E. Platen. Numerical Solution of Stochastic Differential Equations. Berlin, New York, Springer-Verlag, 1992.
- [18] D. Kramkov and W. Schachermayer. The asymptotic elasticity of utility functions and optimal investment in incomplete markets. Ann. Appl. Probab., 9(3):904–950, 1999.
- [19] N.V. Krylov. On the rate of convergence of finite-difference approximations for Bellman’s equations. St. Petersburg Math. J., 9:639–650, 1997.
- [20] N.V. Krylov. Approximating value functions for controlled degenerate diffusion processes by using piece-wise constant policies. Electron. J. Probab., 4(2):1–19, 1999.
- [21] N.V. Krylov. On the rate of convergence of finite-difference approximations for Bellman’s equations with variable coefficients. Probab. Theory Relat. Fields, 117:1–16, 2000.
- [22] J.L. Menaldi. Some estimates for finite difference approximations. SIAM J. Control Optim., 27:579–607, 1989.
- [23] R.C. Merton. Optimal consumption and portfolio rules in continuous time. J. Economic Theory, 3:373–413, 1971.
- [24] H. Pham. Continuous-time Stochastic Control and Optimization with Financial Applications, volume 61. Series Stochastic Modeling and Applied Probability, Springer, 2009.
- [25] A. Picarelli and C. Reisinger. Probabilistic error analysis for some approximation schemes to optimal control problems. Preprint, arXiv:1810.04691, 2018.
- [26] L.C.G. Rogers. Duality in Constrained Optimal Investment Problems: A Synthesis. Number 1814 in Paris–Princeton Lectures on Mathematical Finance. Springer, 2002.
- [27] J. Yong and X.Y. Zhou. Stochastic Controls: Hamiltonian Systems and HJB Equations, volume 43 of Applications of Mathematics. Springer-Verlag, New York, 1999.

DEPARTMENT OF ECONOMICAL SCIENCES, UNIVERSITY OF VERONA, VIA CANTARANE 24, 37129, VERONA, ITALY  
*E-mail address:* `athena.picarelli@univr.it`

MATHEMATICAL INSTITUTE, UNIVERSITY OF OXFORD, ANDREW WILES BUILDING, OX2 6GG, OXFORD, UK  
*E-mail address:* `christoph.reisinger@maths.ox.ac.uk`