

# FINITE ELEMENT APPROXIMATIONS FOR STOCHASTIC CONTROL PROBLEMS WITH UNBOUNDED STATE SPACE

MARTIN G. VIETEN,\*\*\* AND RICHARD H. STOCKBRIDGE,\*\*\*\* University of Wisconsin-Milwaukee

#### Abstract

A numerical method is proposed for a class of one-dimensional stochastic control problems with unbounded state space. This method solves an infinite-dimensional linear program, equivalent to the original formulation based on a stochastic differential equation, using a finite element approximation. The discretization scheme itself and the necessary assumptions are discussed, and a convergence argument for the method is presented. Its performance is illustrated by examples featuring long-term average and infinite horizon discounted costs, and additional optimization constraints.

*Keywords:* Stochastic control; finite element method; linear programming; relaxed controls

2020 Mathematics Subject Classification: Primary 93E20; 65C20 Secondary 90C05

# 1. Introduction

# 1.1. Motivation and literature

This paper considers a class of stochastic control problems for a process *X* whose dynamics are initially specified by the stochastic differential equation (SDE)

$$dX_t = b(X_t, u_t) dt + \sigma(X_t, u_t) dW_t, \qquad X_0 = x_0,$$
(1.1)

where W is a Brownian motion process and the process u represents the control influencing the evolution of X. Both X and u are assumed to be real-valued, and in particular X is allowed to have an unbounded state space E. Given a cost function  $\tilde{c}$ , the control u has to be chosen from a set of admissible controls in such a way that it minimizes either a long-term average or an infinite horizon discounted cost criterion, defined respectively by

$$\limsup_{t \to \infty} \frac{1}{t} \mathbb{E} \left[ \int_0^t \tilde{c}(X_s, u_s) \, \mathrm{d}s \right] \quad \text{and} \quad \mathbb{E} \left[ \int_0^\infty \mathrm{e}^{-\alpha s} \tilde{c}(X_s, u_s) \, \mathrm{d}s \right], \tag{1.2}$$

for some discounting rate  $\alpha > 0$ . Additional constraints on the evolution of X and u can be imposed using criteria which are formulated in a similar fashion. Such control problems

Received 9 February 2023; accepted 22 July 2024.

<sup>\*</sup> Postal address: Department of Mathematical Sciences, University of Wisconsin-Milwaukee, P.O. Box 413, Milwaukee, WI 53201-0413, USA.

<sup>\*\*</sup> Email address: mvieten@uwm.edu

<sup>\*\*\*</sup> Email address: stockbri@uwm.edu

<sup>©</sup> The Author(s), 2024. Published by Cambridge University Press on behalf of Applied Probability Trust. This is an Open Access article, distributed under the terms of the Creative Commons Attribution licence (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted re-use, distribution, and reproduction in any medium, provided the original work is properly cited.

are considered in a relaxed sense by using a martingale problem formulation involving the infinitesimal generator of X, and an equivalent infinite-dimensional linear program for the expected occupation measure of both the process X and the control u. With a certain set of assumptions, approximate solutions to this linear program are attained by restricting the state space to a bounded subset of E, discretizing the infinite-dimensional constraint space of the linear program using a finite element approach, and introducing discrete approximations of the expected occupation measure.

Previous work [27, 28] has focused on processes X which are kept within a bounded state space by either reflection or jump boundaries, or by specific assumptions on the coefficient functions b and  $\sigma$ . The analysis presented in this paper generalizes the approach to processes with an unbounded state space, under specific assumptions on the cost function  $\tilde{c}$ . On top of that, we present how additional optimization constraints can easily be integrated into the linear program structure, which is applicable to both the class of problems presented herein and those problems discussed in [27].

The classical approach to stochastic control problems is given by methods based on the dynamic programming principle, as presented in [7, 8].

Central to these methods is the solution of the Hamilton–Jacobi–Bellman (HJB) equation, taking the form of a second-order, non-linear differential equation. The survey article [23] elaborates on this approach, while also contrasting it with Pontryagin's maximum principle and more recent developments in backward stochastic differential equations.

Linear programming approaches have also been instrumental in the analytic treatment of various stochastic control problems. A first example is given in [21], where an ergodic Markov chain for an inventory problem under long-term average costs is analyzed. Both [2, 16] investigate the linear programming approach for solutions of controlled martingale problems using long-term average and discounted cost criteria for infinite horizon problems, as well as finite horizon and first exit problems for absolutely continuous control. A multi-dimensional diffusion with singular control is considered in [26], while [25] investigates jump diffusions of Lévy–Itô type.

Both the dynamic programming principle and the linear programming approach can be used to establish numerical schemes that solve control problems approximately. An example of the former is given in [18], where the HJB equation is solved for a discrete Markov chain approximating the continuous process. Other such approaches rely on employing numerical techniques for differential equations solving the HJB equation itself. In certain instances, finite difference approximations are equivalent to discrete Markov chain approximations, as illustrated in [4]. Finite element approximations are used in [13, 15]. Another numerical technique using dynamic programming was analyzed in [1]. A survey of numerical techniques for the HJB equation and their theoretical foundations was compiled in [9].

A very general setting for numerical schemes based on the linear programming approach can be found in [22]. Moment-based approaches have been used, e.g. in [11, 19]. A technique employing the finite element method is described in [14, 24], where the infinite-dimensional constraints of the linear program are discretized using a set of finite basis functions. This technique is investigated further in [27], where its convergence for problems with bounded state spaces is analyzed. Additional attention was brought towards theoretic foundations again in [28].

For their convergence arguments, [27, 28] rely on the state space *E* of the considered process being bounded. This assumption allows the use of standard arguments regarding continuous functions over compact spaces when analyzing convergence. However, it imposes a strong

restriction on the types of models to which the method is applicable, as many practical problems typically have an unbounded state space. This article provides a numerical scheme that places additional restrictions on the cost function rather than on the state space, and therefore loosens this restriction. Its contribution lies in identifying the correct assumptions on the cost function that in turn allow the assumption of a bounded state space to be dropped. Furthermore, it demonstrates the generalization of the method established in [27], defining a new set of finite element basis functions suited for unbounded state spaces, and providing an extended convergence argument. In addition, we illustrate how additional optimization constraints which are expressed in the same fashion as the cost criterion can naturally be integrated into the linear program formulation.

Approximate HJB-based approaches to problems with an unbounded state space can be found in [20, 29], where wealth and asset allocation problems are investigated. The former focuses on a finite-time setup and an optimality criterion that incorporates both expected wealth and risk. The latter considers stock build-up with limited fund availability, on an infinite time horizon. The optimization constraint of limited fund availability is expressed therein as an additional criterion on the process and its control, and is incorporated using a Lagrange multiplier approach.

For bounded state spaces, the finite element approach has shown performance that is competitive with other methods, in particular with HJB approaches. By generalizing this approach to unbounded state spaces, we provide a viable alternative to established approaches, especially HJB-based approaches, for a large class of problems. A distinct feature of the proposed method is that the inclusion of optimization constraints is natural, and requires neither additional discretization nor large adaptions of the convergence proof.

In order to illustrate the proposed method, this paper is structured as follows. In the remainder of this section, we introduce the linear programming approach, present formulations for additional constraints, and review results pertinent to the subsequent analysis. Section 2 introduces the discretization stages of the method, and the main results. Section 3 gives some details on the implementation, and illustrates the applicability of the method on three examples. Section 4 discusses the technical details of the convergence proofs. Section 5 concludes the article with a discussion of possible future research directions.

#### 1.2. Notation and formalism

We denote the set of non-negative real numbers by  $\mathbb{R}_+$ , and the set of positive real numbers by  $\mathbb{R}_+ - \{0\}$ . The complement of a set *S* is denoted by  $S^{\sim}$ . For  $S \subset \mathbb{R}$ , let C(S) be the set of continuous functions on *S*,  $C^2(S)$  the set of twice continuously differentiable functions on *S*, and  $C_c^2(S)$  be its subset of functions whose support is contained in a compact subset of *S*. The set of functions on *S* which are continuous except for finitely many points, but either leftor right-continuous at these points, is denoted by  $\dot{C}(S)$ . The set of continuous and bounded functions on *S* is denoted by  $C_b(S)$ . The uniform norm of a bounded function *f* is denoted by  $\|f\|_{\infty}$ . We use  $\mathscr{B}(S)$  to refer to the Borel  $\sigma$ -algebra on *S*.  $\mathcal{P}(S)$  denotes the set of probability measures, and  $\mathcal{M}(S)$  the set of finite Borel measures on  $(S, \mathscr{B}(S))$ . The sign function is denoted by  $\operatorname{sgn}(x)$ . For a function *f* on *S*,  $[f]^+(x) := \max(f(x), 0)$  denotes its positive part. We use 'a.e.' to abbreviate 'almost everywhere'.

Consider the SDE given by (1.1). We assume that  $X_t \in E$ , where *E* is an interval in  $\mathbb{R}$  (or  $\mathbb{R}$  itself), and  $u_t \in U = [u_1, u_r]$  with  $-\infty < u_1 < u_r < \infty$ , for all  $t \ge 0$ . *E* and *U* are called the state space and the control space, respectively. Throughout this paper, we consider the case where  $E = (-\infty, \infty)$ , although the proposed discretization approach is also applicable to the cases

 $E = [e_1, \infty)$  and  $E = (-\infty, e_r]$ , for  $-\infty < e_1 < e_r < \infty$ . Assume that the coefficient functions  $b: E \times U \mapsto \mathbb{R}$  and  $\sigma: E \times U \mapsto \mathbb{R}_+ - \{0\}$  both lie in  $C(E \times U)$ . They are called the drift and diffusion functions, respectively. The infinitesimal generator  $\tilde{A}$  of a process solving (1.1) is defined, for  $f \in C^2(E)$ , by  $\tilde{A}f(x, u) = b(x, u)f'(x) + \frac{1}{2}\sigma^2(x, u)f''(x)$ . Throughout this paper, we consider the restriction of  $\tilde{A}$  to  $C_c^2(E)$ . A specification of the dynamics that requires

$$f(X_t) - f(x_0) - \int_0^t \tilde{A} f(X_s, u_s) \,\mathrm{d}s$$
 (1.3)

to be a martingale for all  $f \in C_c^2(E)$  is equivalent to (1.1) in terms of weak solutions. Hence, the values of the cost criteria determined by (1.2) remain identical. A relaxed formulation of (1.3) is better suited to stochastic control.

**Definition 1.1.** Let *X* be a stochastic process with state space *E*, and let  $\Lambda$  be a stochastic process taking values in  $\mathcal{P}(U)$ . (*X*,  $\Lambda$ ) is a relaxed solution to the controlled martingale problem for  $\tilde{A}$  if there is a filtration  $\{\mathscr{F}_t\}_{t\geq 0}$  of  $\mathscr{B}(E)$  such that *X* and  $\Lambda$  are  $\mathscr{F}_t$ -progressively measurable and

$$f(X_t) - f(x_0) - \int_0^t \int_U \tilde{A} f(X_s, u) \Lambda_s(\mathrm{d}u) \,\mathrm{d}s \tag{1.4}$$

is an  $\{\mathscr{F}_t\}_{t\geq 0}$ -martingale for all  $f \in C^2_{c}(E)$ .

The relaxation is given by the fact that the control is no longer represented by a process u, but is encoded in the measure-valued process  $\Lambda$ . The cost criteria for a relaxed solution of the controlled martingale problem for  $\tilde{A}$  are given by

$$\limsup_{t \to \infty} \frac{1}{t} \mathbb{E} \left[ \int_0^t \int_U \tilde{c}(X_s, u) \Lambda_s(\mathrm{d}u) \,\mathrm{d}s \right] \quad \text{and} \quad \mathbb{E} \left[ \int_0^\infty \int_U \mathrm{e}^{-\alpha s} \tilde{c}(X_s, u) \Lambda_s(\mathrm{d}u) \,\mathrm{d}s \right] \quad (1.5)$$

for the long-term average cost criterion and the infinite horizon discounted cost criterion, respectively. A stochastic control problem given by (1.4), together with one of the criteria from (1.5), can be reformulated as an infinite-dimensional linear program. To this end, for  $\alpha \ge 0$  and  $x_0 \in E$ , define the function *c* and two operators  $A: C_c^2(E) \mapsto C(E \times U)$  and  $R: C_c^2(E) \mapsto \mathbb{R}$  by

$$c(x, u) = \begin{cases} \tilde{c}(x, u) & \text{if } \alpha = 0, \\ \tilde{c}(x, u)/\alpha & \text{if } \alpha > 0, \end{cases}$$

$$(Af)(x, u) = \tilde{A}f(x, u) - \alpha f(x),$$

$$Rf = -\alpha f(x_0).$$
(1.6)

In this setup,  $\alpha = 0$  represents the long-term average, and  $\alpha > 0$  the infinite horizon discounted cost criterion case.

**Definition 1.2.** The infinite-dimensional linear program for a stochastic control problem with unbounded state space is given by

Minimize 
$$J \equiv \int_{E \times U} c \, d\mu$$
 subject to 
$$\begin{cases} \int_{E \times U} Af \, d\mu = Rf & \text{for all } f \in C_c^2(E), \\ \mu \in \mathcal{P}(E \times U). \end{cases}$$
 (1.7)

The measure  $\mu$  is the so-called expected occupation measure (long-term average case) or expected discounted occupation measure (infinite horizon discounted case) of X and A. It

encapsulates both the evolution of X and the control u. For the long-term average cost criterion,  $\mu$  is independent of the starting point  $x_0$  of X, while for the infinite horizon discounted criterion,  $\mu$  does depend on  $x_0$ , as is evident from *Rf*. To keep the notation simple, the possible dependency of  $\mu$  (and thereby that of J) on  $x_0$  is omitted in the following.

The measure on  $(E, \mathscr{B}(E))$  given by  $\mu_E(\cdot) = \mu(\cdot \times U)$  is called the state-space marginal of  $\mu$ . If  $P: E \times U \ni (x, u) \mapsto x \in E$  is the projection map from  $E \times U$  onto  $E, \mu_E$  is the distribution of P under  $\mu$ . The relation of an infinite-dimensional linear program as seen in Definition 1.2 to stochastic control problems is explored in [16]. Therein, a certain class of relaxed controls, which can be considered as regular conditional probabilities, are of central interest.

**Definition 1.3.** Let  $(E \times U, \mathscr{B}(E \times U), \mu)$  be a measure space. A mapping  $\eta : \mathscr{B}(U) \times E \mapsto [0, 1]$  is called a regular conditional probability of  $\mu$  if:

- (i) for each  $x \in E$ ,  $\eta(\cdot, x)$ :  $\mathscr{B}(U) \mapsto [0, 1]$  is a probability measure;
- (ii) for each  $V \in \mathscr{B}(U)$ ,  $\eta(V, \cdot): E \mapsto [0, 1]$  is a measurable function; and
- (iii) for all  $V \in \mathscr{B}(U)$  and all  $F \in \mathscr{B}(E)$ ,  $\mu(F \times V) = \int_{F} \eta(V, x) \mu_{E}(dx)$ .

**Remark 1.1.** By [6, Theorem 8.1 and Remark 8.2], the existence of regular conditional probabilities for any considered solution to (1.7) follows from the fact that *E* is a measurable space and *U* is a complete, separable metric space, with the fact that  $\mu(E \times U) = 1$ .

Using the notion of a regular conditional probability, Theorem 1.1 addresses the equivalence of infinite-dimensional linear programs and controlled martingale problems.

**Theorem 1.1.** The problem of minimizing one of the cost criteria of (1.5) over the set of all relaxed solutions  $(X, \Lambda)$  to the controlled martingale problem for  $\tilde{A}$  is equivalent to the linear program stated in Definition 1.2. Moreover, there exists an optimal solution  $\mu^*$ . Let  $\eta^*$  be the regular conditional probability of  $\mu^*$  with respect to its state-space marginal  $\mu_E^*$ . Then, an optimal relaxed control is given in feedback form by  $\Lambda_t^* = \eta^*(\cdot, X_t^*)$ , where  $(X^*, \Lambda^*)$  is a relaxed solution to the controlled martingale problem for  $\tilde{A}$  having occupation measure  $\mu^*$ .

*Proof.* See [16, Theorems 6.1 and 6.3] and the erratum in [17] for both the proof and the underlying technical conditions.  $\Box$ 

By this result, it suffices to find optimal solutions to (1.7) in order to solve a stochastic control problem for a cost criterion of (1.2), with the dynamics specified by (1.1). Moreover, we can focus our attention on such measures  $\mu(dx, du)$  that can be factored into  $\eta(du, x)\mu_E(dx)$ , as an optimal solution lies within this set of measures.

**Remark 1.2.** By considering relaxed solutions, the solution space is large enough to make convergence arguments in the proof of Theorem 1.1 work. Practically, relaxed controls will only appear at the transition points (of Lebesgue measure 0) of so-called bang-bang controls, or for rather pathological combinations of the coefficient function b and the cost function c. In both cases, there will always be a non-relaxed, 'strict' control (which assigns full mass on a single control value) with equal value of the cost criterion.

Analytic solutions are usually hard to obtain, making approximate treatments of such control problems necessary. Such numerical approaches have to address three challenges. First, the unboundedness of the state space has to be taken into account. Second, the infinite-dimensional constraint space of (1.7), given by  $C_c^2(E)$ , has to be discretized. Finally, the space of measures  $\mathcal{P}(E \times U)$  has to be made computationally tractable. In order to establish a convergent numerical scheme addressing these challenges, several assumptions have to be introduced. To do so, we need the following two definitions.

**Definition 1.4.** A function  $c : \mathbb{R} \times U \mapsto \mathbb{R}_+$  will be called increasing in |x| if, first, for any  $L_1 \ge 0$ , there exist  $L_2, L_3 > L_1$  such that  $c(x, u) > c(-L_2, u)$  for all  $x < -L_2$  and  $c(x, u) > c(L_3, u)$  for all  $x > L_3$ , and, second, for any K > 0 there is an L large enough such that c(x, u) > K for all  $x \in [-L, L]^{\sim}$ , uniformly in u.

**Definition 1.5.** A function  $c: \mathbb{R} \times U \mapsto \mathbb{R}_+$  allows for compactification if, for all K > 0, there exists a continuous function  $u^-: [-K, K] \mapsto U$  satisfying

 $\sup\{c(x, u^{-}(x)) \mid x \in [-K, K]\} \le \inf\{c(x, u) \mid x \in [-K, K]^{\sim}, u \in U\}.$ 

**Remark 1.3.** Intuitively, these two rather technical conditions describe that it is inherently preferable to confine the process X within a compact set. If c is increasing in |x|, it will grow beyond bound outside of any given compact set. If c allows for compactification, its structure is such that there exists at least one control  $u^-$  for which it is advantageous for the process to remain within [-K, K], as compared to remaining outside it (in  $[-K, K]^{\sim}$ ) under any other possible control.

**Remark 1.4.** A broad class of cost functions *c* which are increasing in |x| also allows for compactification. For example, an additive structure of the costs, such as  $(x, u) \mapsto x^2 + u^2$ , satisfies this property, as do more complex functions such as  $(x, u) \mapsto x^2 + |x|u^2$ . Trivially, a function with no cost of control such as  $x \mapsto |x|$  also allows for compactification.

Throughout this paper, we assume the following conditions on (1.7).

# Assumption 1.1

- (i) The coefficient functions b and  $\sigma$  are continuous.
- (ii) The cost function c is continuous and non-negative.
- (iii) The cost function c is increasing in |x|.
- (iv) The cost function c allows for compactification.
- (v) There exists a solution  $\mu$  to (1.7) such that  $\int_{E \times U} c \, d\mu < \infty$ .

**Remark 1.5.** Assumption 1.1(v) requires that there exists some measure for which the cost integral is finite. Should this condition not be satisfied, then every admissible control process will lead to an infinite cost and every control will be optimal so there is no point to the optimization. In this situation, the problem has most likely been misspecified. Since the diffusion processes are Gaussian at each time, the expected occupation measures will also be Gaussian, provided the drift and diffusion coefficients do not send the process to  $\pm \infty$ . As long as *c* does not grow too quickly at  $\pm \infty$ , the natural decay of the Gaussian distribution will dominate the growth of *c*, resulting in Assumption 1.1(v) being satisfied. Thus, a measurable feedback control u(x) such that the drift coefficient *b* strictly pushes towards the origin from some point outwards while the diffusion coefficient remains bounded will result in an occupation measure with Gaussian decay. Then *c* having polynomial growth will lead to a finite cost, satisfying the condition. Typically, this condition is verified on a case-by-case basis.

**Remark 1.6.** Assumptions 1.1(iii) and (iv) guarantee that the problem is computationally attainable by discrete methods by having an optimal solution that confines the process to a compact set; cf. Remark 1.3. Any other problem would naturally lie outside the scope of the proposed method.

**Assumption 1.2.** For any considered solution  $\mu(dx, du) \equiv \eta(du, x) \mu_E(dx)$ , the following must be satisfied:

- (i) For a set  $V \subset \mathscr{B}(U)$  which is either a singleton or an interval,  $x \mapsto \eta(V, x)$  is continuous *a.e.* with respect to Lebesgue measure.
- (ii) The functions  $x \mapsto \int_{U} b(x, u) \eta(du, x)$  and  $x \mapsto \int_{U} \sigma(x, u) \eta(du, x)$  lie in  $\dot{C}(E)$ .

**Assumption 1.3.** For any considered solution  $\mu(dx, du) \equiv \eta(du, x) \mu_E(dx)$ ,  $\mu_E$  must be absolutely continuous with respect to Lebesgue measure.

**Remark 1.7.** Assumption 1.2 is necessary for the convergence argument for an approximation of the regular conditional probability  $\eta$ . As elaborated in [27, Section 2.2] and [28, Section 2.2], a large class of controls satisfies this assumption. In particular, this applies to bang-bang controls and controls  $\eta$  that satisfy  $\eta(\{v(x)\}, x)) = 1$  for some continuous function v.

# 1.3. Additional optimization constraints

The control problem given by (1.1), with one of the criteria from (1.2), can be enhanced by placing constraints on the process X and the control u. To this end, take two non-negative functions  $g_1, g_2 \in C(E \times U)$ , and  $D_1, D_2 \ge 0$ . For the long-term average cost criterion, these constraints take the form

$$\limsup_{t\to\infty}\frac{1}{t}\mathbb{E}\bigg[\int_0^t g_1(X_s, u_s)\,\mathrm{d}s\bigg] \le D_1 \quad \text{and} \quad \limsup_{t\to\infty}\frac{1}{t}\mathbb{E}\bigg[\int_0^t g_2(X_s, u_s)\,\mathrm{d}s\bigg] = D_2,$$

the former representing a linear inequality constraint, and the latter representing a linear equality constraint. Respectively, for the infinite horizon discounted cost criterion, they take the form

$$\mathbb{E}\left[\int_0^\infty e^{-\alpha s} g_1(X_s, u_s) \, \mathrm{d}s\right] \le D_1 \quad \text{and} \quad \mathbb{E}\left[\int_0^\infty e^{-\alpha s} g_2(X_s, u_s) \, \mathrm{d}s\right] = D_2$$

When expressed in terms of the expected occupation measure  $\mu$ , they take the form

$$\int_{E \times U} g_1 \, \mathrm{d}\mu \le D_1 \quad \text{and} \quad \int_{E \times U} g_2 \, \mathrm{d}\mu = D_2 \tag{1.8}$$

and can easily be integrated into the constraints of (1.7) since they are linear in  $\mu$ . The subsequent analysis is carried out without such additional constraints for the sake of presentation. However, the derived results still hold when such constraints are present, as they can be treated equivalently in the convergence proofs.

**Remark 1.8.** This extension of the problem with additional optimization constraints is not limited to problems with unbounded state space as described herein. Constraints in the style of (1.8) can also be introduced to problems with bounded state space, and singular behavior, which are discussed in [27].

## 1.4. Preliminary results

The analysis presented herein relies on established results on the linear programming approach to stochastic control, which we present in the following. We also review results on the weak convergence of measures, and on B-spline basis functions.

Recent research [27, 28] considers stochastic control problems stemming from an SDE similar to (1.1). This SDE takes the form

$$dX_t = b(X_t, u_t) dt + \sigma(X_t, u_t) dW_t + h(X_{t-}) d\xi_t, \qquad X_0 = x_0,$$
(1.9)

which, in contrast to (1.1), allows for singular behavior of X by introducing the term  $h(X_{t-}) d\xi_t$ . This term takes such forms that X is contained in a bounded state space  $E = [e_1, e_r]$ , with  $-\infty < e_1 < e_r < \infty$ . Typically,  $h(X_{t-}) d\xi_t$  models either a reflection at one of the boundary points  $e_1$  and  $e_r$ , or a jump from one of the boundaries into E. Such behavior is characterized by an infinitesimal generator B, which, for the case of a reflection to the right, takes the form Bf(x) = f'(x). For our purposes, it suffices to consider  $e_1 = -K$  and  $e_r = K$ , where  $K \in \mathbb{N}$  is a constant chosen in the discretization steps discussed in Section 2. The cost criteria of (1.2) can be adapted to this situation by introducing a second cost function  $\tilde{c}_1$  and its scaled version  $c_1$ , cf. (1.6), accounting for the cost accrued by the singular behavior. A detailed discussion of such problems, and a derivation of equivalent linear programs analogous to the setup introduced in Section 1.2, is provided in [28]. For our purposes, it suffices to consider such problems without an additional cost function  $\tilde{c}_1$ . The equivalent linear program to (1.9), with one of the criteria introduced in (1.2), reads as follows. As above, we drop the possible dependence of  $\mu_0$  and J on  $x_0$  from the notation.

**Definition 1.6.** The infinite-dimensional linear program for a stochastic control problem with singular boundary behavior is given by

Minimize 
$$J \equiv \int_{[-K,K] \times U} c \, d\mu_0$$
  
subject to 
$$\begin{cases} \int_{[-K,K] \times U} Af \, d\mu_0 + \int_{[-K,K]} Bf \, d\mu_1 = Rf \text{ for all } f \in C_c^2([-K,K]), \\ \mu_0 \in \mathcal{P}([-K,K] \times U), \\ \mu_1 \in \mathcal{M}([-K,K]). \end{cases}$$
(1.10)

**Remark 1.9.** In Remark 2.2, we consider a specific case of this problem where the singular behavior of *X* is given by reflections at the boundaries of the state space -K and *K*. In this case,  $\mu_1$  has full mass on  $\{-K\}$  and  $\{K\}$ , and  $\int_{[-K,K]} Bf d\mu_1 = f'(-K) \cdot \mu_1(\{-K\}) - f'(K) \cdot \mu_1(\{K\})$ .

In Definition 1.6, the expected occupation measure  $\mu_0$  can be considered to be the analogue of the measure  $\mu$  in (1.7), while the expected occupation measure  $\mu_1$  models the singular behavior occurring at the boundaries of *E*. As introduced in Definition 1.3, regular conditional probabilities can be used to model relaxed controls, and the statement of Theorem 1.1 on the equivalence of this linear program with a controlled martingale problem remains true. The solvability of the linear constraints of (1.10) is discussed in [28], providing the following result.

**Theorem 1.2.** Let  $v: E \mapsto U$  be a continuous function, and assume that  $\eta_0$  is a regular conditional probability such that, for all  $x \in E$ ,  $\eta_0(v(x), x) = 1$ . Then, there exists a unique solution  $(\hat{\mu}_0, \hat{\mu}_1)$  to the linear constraints of (1.10) such that, for all  $F \in \mathscr{B}(E)$  and  $V \in \mathscr{B}(U)$ ,  $\hat{\mu}_0(F \times V) = \int_F \eta_0(V, x) d\hat{\mu}_{0,E}$ , where  $\hat{\mu}_{0,E}$  denotes the state-space marginal of  $\hat{\mu}_0$ .

*Proof.* This a consequence of [28, Theorem 2.5].  $\Box$ 

To introduce the notion of weak convergence of measures, let S be a measurable space equipped with a topology.

**Definition 1.7.** Consider a sequence of probability measures  $\{\mu_n\}_{n\in\mathbb{N}}$  and another probability measure  $\mu$  on S. We say that  $\mu_n$  converges weakly to  $\mu$ , denoted  $\mu_n \Rightarrow \mu$ , if, for all  $f \in C_b(S)$ ,  $\int_S f(x) \mu_n(dx) \Rightarrow \int_S f(x) \mu(dx)$  as  $n \to \infty$ .

An extensive discussion of weak convergence of measures can be found in [3, Volume 2, Chapter 8]. Central to our purposes is Theorem 1.3, which states sufficient conditions for the existence of weakly convergent subsequences of sequences of probability measures, based on the following concept.

**Definition 1.8.** A sequence of probability measures  $\{\mu_n\}_{n \in \mathbb{N}}$  on S is called tight if, for each  $\varepsilon > 0$ , there is a compact set  $K_{\varepsilon}$  in S such that  $\mu_n(K_{\varepsilon}^{\sim}) < \varepsilon$  for all  $n \in \mathbb{N}$ .

**Remark 1.10.** If S is compact, any sequence of probability measures on S is tight.

**Theorem 1.3.** Let  $\{\mu_n\}_{n \in \mathbb{N}}$  be a sequence of probability measures on *S*. Then the following are equivalent:

- (i)  $\{\mu_n\}_{n \in \mathbb{N}}$  contains a weakly convergent subsequence.
- (ii)  $\{\mu_n\}_{n\in\mathbb{N}}$  is tight.

*Proof.* [3, Theorem 8.6.2] considers this statement for the more general case of a sequence of *finite* measures, for which the sequence of measures  $\{\mu_n\}$  is required to be uniformly bounded. This requirement is trivially satisfied when considering a sequence of probability measures.

While discretizing the constraints of (1.7) and (1.10), we need to embrace the space  $C_c^2(E)$ as a normed space. To this end, for  $f \in C_c^2(E)$ , set  $||f||_{\mathscr{D}} = ||f||_{\infty} + ||f'||_{\infty} + ||f''||_{\infty}$  and define  $\mathscr{D}_{\infty}(E) = (C_c^2(E), ||\cdot||_{\mathscr{D}})$ . Of special importance is the collection of spaces  $\mathscr{D}_{\infty}((-K, K))$  for  $K \in \mathbb{N}$ . Note that for  $f \in \mathscr{D}_{\infty}((-K, K))$ , the support of f is always a *proper* subset of (-K, K). Furthermore,  $\mathscr{D}_{\infty}(E) \supset \mathscr{D}_{\infty}((-K, K))$  and  $\mathscr{D}_{\infty}(E) = \bigcup_{K \in \mathbb{N}} \mathscr{D}_{\infty}((-K, K))$ . Less obvious is the fact that, for any  $K \in \mathbb{N}$ ,  $\mathscr{D}_{\infty}((-K, K))$  is separable. To show this, we rely on cubic B-spline basis functions. These functions form a basis for the space of all cubic splines on a closed interval, which in turn show specific convergence properties crucial to our claim. To construct these basis functions, fix  $q, K \in \mathbb{N}$  and consider the set of grid points

$$E^{(q,K)} = \left\{ e_r^{(q)} = -K + \frac{2K}{2^q} \cdot r, \ r = -3, -2, \dots, 2^q + 3 \right\}.$$

**Definition 1.9.** For fixed  $q, K \in \mathbb{N}$ , the set of cubic B-spline basis functions on  $E^{(q,K)}$  is defined by

$$\left\{s_r^{(q)} = (e_{r+4}^{(q)} - e_r^{(q)}) \sum_{i=r}^{r+4} \frac{\left[(e_i^{(q)} - x)^3\right]^+}{\Psi'_r(e_i^{(q)})}, \ r = -3, -2, \dots, 2^q - 1\right\},\$$

where  $\Psi_r(x) = \prod_{i=r}^{r+4} (x - e_i^{(q)})$  for  $r = -3, -2, \dots, 2^q - 1$ .

An analysis of these basis function is given in [5]. Provided that

$$\max_{r=-3,\dots,2^{q}+2} \left( e_{r+1}^{(q)} - e_{r}^{(q)} \right) \to 0 \quad \text{and} \quad \frac{\max_{r=-3,\dots,2^{q}+2} \left( e_{r+1}^{(q)} - e_{r}^{(q)} \right)}{\min_{r=-3,\dots,2^{q}+2} \left( e_{r+1}^{(q)} - e_{r}^{(q)} \right)} \to 1$$

as  $q \to \infty$ , [10, Theorem 1] holds and the following statement can be shown.

**Proposition 1.1.** The set  $\bigcup_{q\in\mathbb{N}} \{s_r^{(q)}\}_{r=-3}^{2^q-1}$  is dense in  $C^2([-K, K])$  with respect to  $\|\cdot\|_{\mathscr{D}}$ .

Using this result, we can approximate  $f \in \mathscr{D}_{\infty}((-K, K))$  arbitrarily closely with respect to  $\|\cdot\|_{\mathscr{D}}$  by a linear combination of B-spline basis functions. However, these linear combinations are not guaranteed to have their support in a compact subset of (-K, K). To remedy this, we need to alter the basis functions.

**Definition 1.10.** For  $q \in \mathbb{N}$ , let  $h^{(q)}$  be a piecewise polynomial function in  $C_c^2((-K, K))$  with the following properties.

- (i)  $h^{(q)}(x) = 0$  for  $x \le e_1^{(q)}$ .
- (ii)  $h^{(q)}(x) = 1$  for  $e_2^{(q)} \le x \le e_{2^q-2}^{(q)}$ .
- (iii)  $h^{(q)}(x) = 0$  for  $e_{2^q 1}^{(q)} \le x$ .

For  $q_1, q_2 \in \mathbb{N}$ , consider two sets of grid points  $E^{(q_1,K)}$  and  $E^{(q_2,K)}$ . The set of altered B-spline basis functions for  $q_1$  and  $q_2$  is defined by  $\{h^{(q_1)} \cdot s_r^{(q_2)}\}_{r=-3}^{2q_2-1}$ .

**Proposition 1.2.** The set  $\bigcup_{q_1 \in \mathbb{N}} \bigcup_{q_2 \in \mathbb{N}} \{h^{(q_1)} \cdot s_r^{(q_2)}\}_{r=-3}^{2^{q_2}-1}$  is dense in  $\mathscr{D}_{\infty}((-K, K))$ .

*Proof.* Take  $f \in \mathscr{D}_{\infty}((-K, K))$ , and  $\varepsilon > 0$ . Choose  $q_1 \in \mathbb{N}$  large enough that  $\operatorname{supp}(f) \subset (e_2^{(q_1)}, e_{2^{q_1}-2}^{(q_1)})$ . Interpolate f on [-K, K] by a linear combination  $s^{(q_2)}$  of B-spline basis functions on  $E^{(q_2, K)}$ , with  $q_2 \ge q_1$  such that  $||f - s^{(q_2)}||_{\mathscr{D}} < \varepsilon/(4 \cdot \max\{1, ||h^{(q_1)}||_{\mathscr{D}}\})$ . Set  $\tilde{s} = h^{(q_1)} \cdot s^{(q_2)}$ , and  $I_1 = [e_1^{(q_1)}, e_{2^{q_1}-1}^{(q_1)}]$ ,  $I_2 = [e_2^{(q_1)}, e_{2^{q_1}-2}^{(q_1)}]$ ,  $I_3 = [e_1^{(q_1)}, e_2^{(q_1)}]$  for notation

Set  $\tilde{s} = h^{(q_1)} \cdot s^{(q_2)}$ , and  $I_1 = \lfloor e_1^{(q_1)}, e_{2^{q_1}-1} \rfloor$ ,  $I_2 = \lfloor e_2^{(q_1)}, e_{2^{q_1}-2} \rfloor$ ,  $I_3 = \lfloor e_1^{(q_1)}, e_2^{(q_1)} \rfloor$  for notational purposes. By construction,  $\tilde{s}$  and  $s^{(q_2)}$  are identical on  $I_2$ , and thereby  $M_1 := \|(\tilde{s} - f)\|_{\mathscr{D}(I_2)} < \varepsilon$  and  $\|(\tilde{s} - f)\|_{\mathscr{D}(I_1^-)} = 0$  as  $\tilde{s} \equiv f \equiv 0$  on  $I_1^-$ . Furthermore,  $\|s^{(q_2)}\|_{\mathscr{D}(I_3)} < \varepsilon/(4 \cdot \max\{1, \|h^{(q_1)}\|_{\mathscr{D}}\})$  as  $f \equiv 0$  on  $I_3$ . But then,

$$\begin{split} M_{2} &:= \| (\tilde{s} - f) \|_{\mathscr{D}(I_{3})} \\ &= \| (h^{(q_{1})} \cdot s^{(q_{2})}) \|_{\mathscr{D}(I_{3})} \\ &= \| (h^{(q_{1})} \cdot s^{(q_{2})}) \|_{I_{3}} \|_{\infty} + \| (h^{(q_{1})} \cdot s^{(q_{2})})' \|_{I_{3}} \|_{\infty} + \| (h^{(q_{1})} \cdot s^{(q_{2})})'' \|_{I_{3}} \|_{\infty} < \varepsilon \end{split}$$

by the construction of  $h^{(q_1)}$  and  $s^{(q_2)}$ , the properties of the uniform norm, and the rules of differentiation. The same holds for  $M_3 := \|(\tilde{s}-f)\|_{\mathscr{D}(\left[e_{2^{q_1}-2}^{(q_1)}, e_{2^{q_1}-1}^{(q_1)}\right])}$ . Hence,  $\|\tilde{s}-f\|_{\mathscr{D}_{\infty}} = \max\{M_1, M_2, M_3\} < \varepsilon$ .

**Remark 1.11.** By the properties of the norm  $\|\cdot\|_{\mathscr{D}}$ , convergence in  $\mathscr{D}_{\infty}((-K, K))$  in particular means bounded, pointwise convergence in function, first, and second derivative. If.

for a sequence  $\{g_k\}_{k\in\mathbb{N}}$  of functions and another function  $f, g_k \to f$  in  $\mathscr{D}_{\infty}((-K, K))$  as  $k \to \infty$ , then  $\int_{E \times U} Af \, d\mu = \lim_{k \to \infty} \int_{E \times U} Ag_k \, d\mu$  is satisfied for any  $\mu \in \mathcal{P}(E \times U)$ , as is  $Rf = \lim_{k \to \infty} Rg_k$ .

# 2. Discretization

This section presents three discretization stages needed for our numerical scheme, and states the respective results that show its convergence. This is conducted for  $E = (-\infty, \infty)$ . In order to deal with the cases  $E = [e_1, \infty)$  and  $E = (-\infty, e_r]$  for  $-\infty < e_1 < e_r < \infty$ , singular behavior of X, like reflection and jumps at  $e_1$  or  $e_r$ , has to be introduced, or the coefficient functions b and  $\sigma$  have to be chosen in such a way that X remains inside these intervals. In these cases, a hybrid approach of the technique presented in [27] and the technique presented here has to be employed.

# 2.1. Addressing the unboundedness of the state space

Looking at the constraints in (1.7), we denote the set of feasible solutions as

$$\mathscr{M}_{\infty} = \left\{ \mu \in \mathcal{P}(E \times U) \colon \int_{E \times U} Af \, \mathrm{d}\mu = Rf \text{ for all } f \in \mathscr{D}_{\infty}(E) \right\}.$$

We denote the cost criterion by  $J: \mathcal{P}(E \times U) \ni \mu \mapsto J(\mu) = \int_{E \times U} c \, d\mu$ . To make the feasible solutions computationally tractable, for  $K \in \mathbb{N}$  we introduce

$$\mathcal{M}_{\infty,K} = \left\{ \mu \in \mathcal{P}(E \times U) \colon \int_{E \times U} Af \, \mathrm{d}\mu = Rf \text{ for all } f \in \mathcal{D}_{\infty}((-K,K)) \right\},$$
$$\mathcal{M}_{\infty,K} = \left\{ \mu \in \mathcal{P}([-K,K] \times U) \colon \int_{E \times U} Af \, \mathrm{d}\mu = Rf \text{ for all } f \in \mathcal{D}_{\infty}((-K,K)) \right\}.$$

 $\mathcal{M}_{\infty,K}$  has fewer constraints than  $\mathcal{M}_{\infty}$ , as it only considers those constraint functions in  $\mathcal{D}_{\infty}(E)$  with support within (-K, K). Thus,  $\mathcal{M}_{\infty} \subset \mathcal{M}_{\infty,K}$ , while on the other hand,  $\mathcal{M}_{\infty,K} \supset \mathcal{M}_{\infty,K}$  since the latter set is more restrictive by requiring measures to have mass in [-K, K]. Foremost, however, is that  $\mathcal{M}_{\infty,K}$  features both measures and constraints that 'live' in a bounded set, which are accessible to discrete approaches.

**Remark 2.1.** The choice of considering measures  $\mu \in \mathcal{P}([-K, K] \times U)$  on the closed interval [-K, K] while considering functions  $f \in \mathcal{D}_{\infty}((-K, K))$  on the open interval (-K, K) is taken purposefully here, as it will ensure that the result of Proposition 1.2 is applicable in the convergence proof of Proposition 2.2,

The convergence result of Proposition 2.1 shows how so-called  $\varepsilon$ -optimal measures in  $\mathcal{M}_{\infty,K}$  relate to  $\varepsilon$ -optimal measures in  $\mathcal{M}_{\infty}$ .

**Definition 2.1.** A measure  $\mu^{\varepsilon} \in \mathscr{M}_{\infty}$  is said to be  $\varepsilon$ -optimal for  $\mathscr{M}_{\infty}$  if  $J(\mu^{\varepsilon}) - \varepsilon < J(\mu)$  for all  $\mu \in \mathscr{M}_{\infty}$ . In identical fashion,  $\varepsilon$ -optimality is defined for  $\mathscr{M}_{\infty,K}$  and  $\mathscr{M}_{\infty,K}$ .

**Proposition 2.1.** For every K > 0, let  $\mu_K^{\varepsilon}$  be an  $\varepsilon$ -optimal solution for  $\mathcal{M}_{\infty,K}$ . Then, given  $\delta > 0$ , there exists a  $K_0(\delta) > 0$  such that, for all  $K \ge K_0$ ,  $\mu_K^{\varepsilon}$  is a  $(2\varepsilon + \delta)$ -optimal solution to  $\mathcal{M}_{\infty}$ .

This result reduces our infinite-dimensional linear program to the task of finding an  $\varepsilon$ optimal measure for  $\mathcal{M}_{\infty,K}$ . It is the main contribution of this article. For its proof in

Section 4.1 we consider arguments concerning both  $\mathcal{M}_{\infty,K}$  and  $\mathcal{M}_{\infty,K}$ . The analysis presented is an adaption of the proof techniques employed in [27].

**Remark 2.2.** Let  $(\mu_0, \mu_1)$  be a solution to the constraints of (1.10), with the singular behavior of *X* given by reflections at -K and *K*; cf. Remark 1.9. For  $f \in \mathcal{D}((-K, K))$ ,  $\int_{[-K,K]\times U} Af \, d\mu_0 = Rf$  as the support of *f* is a proper subset of [-K, K]. Hence,  $\mu_0 \in \mathcal{M}_{\infty,K}$ . Furthermore, if  $x_0 = K$ , Rf = 0 for both the long-term average and infinite horizon discounted cost criterion.

**Remark 2.3.** Under Assumption 1.2, by [28, Theorem 2.5], for any solution  $(\mu_0, \mu_1)$  to the constraints of (1.10), we have that  $\mu_{0,E}$  is absolutely continuous with respect to Lebesgue measure. By Remark 2.2,  $\mathcal{M}_{\infty,K}$  contains solutions  $\mu$  whose state-space marginal  $\mu_E$  is absolutely continuous with respect to Lebesgue measure.

#### 2.2. Discretizing the constraint space

 $\mathscr{M}_{\infty,K}$  remains an analytic construct, and we have to introduce discretizations to make it computationally tractable. Both the function space  $\mathscr{D}_{\infty}((-K, K))$  and the set of measures  $\mathcal{P}([-K, K] \times U)$  are approximated by discrete forms. To approach the first, for given K and some  $q_1, q_2 \in \mathbb{N}$ , consider the set of altered B-spline basis functions from Definition 1.10. Its linear span is denoted by  $\mathscr{D}_{(q_1,q_2)}((-K, K))$ . Following Proposition 1.2,  $\bigcup_{q_1 \in \mathbb{N}, q_2 \in \mathbb{N}} \mathscr{D}_{(q_1,q_2)}((-K, K))$  is dense in  $\mathscr{D}_{\infty}((-K, K))$ . Set

$$\mathscr{M}_{(q_1,q_2),K} = \left\{ \mu \in \mathcal{P}([-K,K] \times U) \colon \int Af \, \mathrm{d}\mu = Rf \text{ for all } f \in \mathscr{D}_{(q_1,q_2)}((-K,K)) \right\}, \quad (2.1)$$

which features a finite set of constraints given by the finite number of B-spline basis functions that span  $\mathcal{D}_{(q_1,q_2)}((-K, K))$ . Set  $J^*_{(q_1,q_2),K} = \inf\{J(\mu) \mid \mu \in \mathcal{M}_{(q_1,q_2),K}\}$ .

**Remark 2.4.** For our purposes, it suffices to consider a sequence of indices  $q_n \equiv (q_{n,1}, q_{n,2})$  with  $q_{n,1} \leq q_{n,2}$  and  $\lim_{n\to\infty} q_{n,1} = \infty$ , and the respective set  $\mathcal{M}_{q_n,K}$ .

Definition 2.1 carries over to  $\mathcal{M}_{q_n,K}$ , and the convergence result is follows.

**Proposition 2.2.** For every  $n \in \mathbb{N}$ , let  $\mu_n^{\varepsilon}$  be an  $\varepsilon$ -optimal solution for  $\mathcal{M}_{q_n,K}$ . Then, given a  $\delta > 0$ , there exists an  $N_0(\delta) \in \mathbb{N}$  such that, for all  $n \ge N_0$ ,  $\mu_n^{\varepsilon}$  is a  $(2\varepsilon + \delta)$ -optimal solution to  $\mathcal{M}_{\infty,K}$ .

This result represents another reduction of the initial problem, as the constraints are now computationally tractable. Its proof is akin to that of Proposition 2.1. Therefore, Section 4.2 merely describes the necessary alterations with respect to Section 4.1.

## 2.3. Discretizing the expected occupation measure

For fixed  $q_1, q_2 \in \mathbb{N}$ , take an expected occupation measure  $\mu \in \mathcal{M}_{(q_1,q_2),K}$  and decompose it into its regular conditional probability and its state-space marginal by setting  $\mu(dx, du) = \eta(du, x)\mu_E(dx)$ . By assumption,  $\mu_E(dx)$  is absolutely continuous with respect to Lebesgue measure. Hence, we can introduce its density p and write  $\mu(dx, du) = \eta(du, x)p(x) dx$ . The following discretization makes  $\eta$  and p tractable. As c, b, and  $\sigma$  are continuous, and we have restricted our attention to the compact set [-K, K], for all  $m \in \mathbb{N}$  there is a  $\delta_m > 0$  such that, for all  $u, v \in U$  with  $|u - v| \le \delta_m$ ,

$$\max\left\{ |c(x, u) - c(x, v)|, |b(x, u) - b(x, v)|, \left| \frac{1}{2} \sigma^2(x, u) - \frac{1}{2} \sigma^2(x, v) \right| \right\} \le \frac{1}{2^{m+1}}$$

holds uniformly in  $x \in [-K, K]$ . Set  $k_m$  to be the smallest integer such that  $(u_r - u_l)/2^{k_m} \le \delta_m$ . The parameter  $k_m$  controls the discretization of the control space U, and the specific choice enables an accurate approximate integration against the relaxed control  $\eta_0$  in the proof of Proposition 4.3. Define

$$U^{(m)} = \left\{ u_j = u_l + \frac{u_r - u_l}{2^{k_m}} \cdot j, \ j = 0, \dots, 2^{k_m} \right\},\$$
$$E^{(m)} = \left\{ e_j = e_l + \frac{e_r - e_l}{2^m} \cdot j, \ j = 0, \dots, 2^m \right\}.$$

The unions of these sets over all  $m \in \mathbb{N}$  are dense in the control space and state space, respectively. They are used to define discretizations for the measure  $\mu$  as follows. First, we approximate the density p of  $\mu_E$ . Denote the countable basis of  $L^1(E)$  given by indicator functions over closed intervals of E by  $\{p_n\}_{n \in \mathbb{N}}$ . Truncate this basis to  $p_0, \ldots, p_{2^m-1}$ , given by the indicator functions of the intervals of length  $1/2^m$ , whose boundary points are given by  $E^{(m)}$ . Then, we approximate the density p by

$$\hat{p}_m(x) = \sum_{k=0}^{2^m - 1} \gamma_k p_k(x),$$
(2.2)

where  $\gamma_k \in \mathbb{R}_+$ ,  $k = 0, ..., 2^m - 1$ , are weights to be chosen such that  $\int_F \hat{p}_m(x) dx = 1$ .

**Remark 2.5.** Other choices of basis functions  $p_0, \ldots, p_{2^m-1}$  can indeed be considered for the discretization of p. As we will see in Section 3.2, analytical solutions to specific control problems may have densities with higher regularity than merely  $L^1(E)$ . In order to keep the analysis as general as possible, we continue using the basis of  $L^1(E)$  given by the indicator functions.

Set 
$$E_j = [e_j, e_{j+1})$$
 for  $j = 0, 1, ..., 2^m - 2$  and  $E_{2^m-1} = [e_{2^m-1}, e_{2^m}]$  to define

$$\hat{\eta}_m(V, x) = \sum_{j=0}^{2^m - 1} \sum_{i=0}^{2^{k_m}} \beta_{j,i} I_{E_j}(x) \delta_{\{u_i\}}(V),$$
(2.3)

where  $\beta_{j,i} \in \mathbb{R}_+$ ,  $j = 0, ..., 2^m - 1$ ,  $i = 0, ..., k_m$ , are weights to be chosen such that  $\sum_{i=0}^{2^{k_m}} \beta_{j,i} = 1$  for  $j = 0, ..., 2^m - 1$ . We approximate  $\eta_0$  using (2.3), which means that this relaxed control is approximated by point masses in the *U*-'direction' and piecewise constants in the *E*-'direction'. Then, we set  $\hat{\mu}_m(du \times dx) = \hat{\eta}_m(du, x)\hat{p}_m(x) dx$ , and introduce the final reduction of our control problem:

$$\mathscr{M}_{(q_1,q_2,m),K} = \{ \hat{\mu}_m \in \mathscr{M}_{(q_1,q_2),K} \colon \mu_m(\mathrm{d}x \times \mathrm{d}u) = \hat{\eta}_m(\mathrm{d}u, x)\hat{p}_m(x)\,\mathrm{d}x \} \}$$

This set features finitely many constraints, inherited from  $\mathcal{M}_{(q_1,q_2),K}$ , but also finitely many unknowns which are the coefficients of  $\hat{\eta}_m$  and  $\hat{p}_m$ . As the constraints are linear, and the objective function J is linear, we are able to find an optimal solution  $\mu_m^*$  in  $\mathcal{M}_{(q_1,q_2,m),K}$  with  $J(\mu^*) \leq J(\mu)$  for all  $\mu \in \mathcal{M}_{(q_1,q_2,m),K}$  by standard (finite-dimensional) linear programming. If additional optimization constraints in the style of (1.8) are present, we can simply evaluate them on  $\hat{\mu}_m$  in the same way we evaluate the cost criterion  $J(\hat{\mu}_m)$ , by integrating  $g_1, g_2$ , or *c* against  $J(\hat{\mu}_m)$ . Doing so, each additional optimization constraint of the original problem introduces merely a single constraint to  $\mathcal{M}_{(q_1,q_2,m),K}$ .

**Proposition 2.3.** For each  $m \in \mathbb{N}$ , assume that  $\mu_m^* \in \mathcal{M}_{(q_1,q_2,m),K}$  and that, for all  $m \in \mathbb{N}$ ,  $\mu_m^*$  is an optimal solution in  $\mathcal{M}_{(q_1,q_2,m),K}$ . Then, the sequence of numbers  $\{J(\mu_m^*)\}_{m \in \mathbb{N}}$  converges to  $J^*_{(q_1,q_2),K}$  as  $m \to \infty$ .

The proof of this result requires a detailed investigation of the approximation properties of (2.2) and (2.3). It is, with some minor modifications, identical to the proof of [27, Corollary 3.3]. Section 4.3 thus only presents a short sketch of the proof.

#### 2.4. Combining the results

The results of the previous three sections combined lead to our main result.

**Theorem 2.1.** For any  $\varepsilon > 0$ , there is a sequence  $\{q_n\}_{n \in \mathbb{N}}$  as in Remark 2.4, and  $K_0, N_0, N_1 \in \mathbb{N}$  such that, for all  $K \ge K_0$ ,  $n \ge N_0$ , and  $m \ge N_1$ , an optimal solution  $\mu_m^*$  in  $\mathcal{M}_{(q_1,q_2,m),K}$  is  $\varepsilon$ -optimal in  $\mathcal{M}_{\infty}$ .

*Proof.* Take  $\varepsilon_1$ ,  $\delta_1 > 0$  satisfying  $2\varepsilon_1 + \delta_1 < \varepsilon$  and, by Proposition 2.1, select  $K_0 \in \mathbb{N}$  such that, for all  $K \ge K_0$ , an  $\varepsilon_1$ -optimal solution to  $\mathscr{M}_{\infty,K}$  is a  $(2\varepsilon_1 + \delta_1)$ -optimal solution to  $\mathscr{M}_{\infty}$ . For any such K, take  $\varepsilon_2$ ,  $\delta_2 > 0$  satisfying  $2\varepsilon_2 + \delta_2 < \varepsilon_1$ . By Proposition 2.2, there is an  $N_0 \in \mathbb{N}$  such that, for all  $n \ge N_0$ , an  $\varepsilon_2$ -optimal solution to  $\mathscr{M}_{q_n,K}$  is  $(2\varepsilon_2 + \delta_2)$ -optimal for  $\mathscr{M}_{\infty,K}$ . For any such n, pick  $N_1 \in \mathbb{N}$  such that, for all  $m \ge N_1$ , an optimal solution  $\mu_m^*$  for  $\mathscr{M}_{(q_{n,1},q_{n,2},m),K}$  is  $\varepsilon_2$ -optimal for  $\mathscr{M}_{(q_{n,1},q_{n,2}),K}$ . By construction,  $\mu_m^*$  is  $\varepsilon$ -optimal in  $\mathscr{M}_{\infty}$ .

# 3. Examples and implementation

We showcase the applicability of the proposed method on three control problems. We start with a very simple problem of a controlled Brownian motion process. This is followed by an Ornstein–Uhlenbeck process with costs of control, which features a more irregular density of the expected occupation measure. A third example also considers an Ornstein–Uhlenbeck process, but introduces a control budget constraint. These examples were selected to showcase the performance of the proposed method for a variety of irregularities typically encountered. As we are dealing with diffusion processes, the density of the expected occupation measure will always be continuous, and infinitely differentiable almost everywhere. Single points where the density is not differentiable occur on 'switching points', where the behavior of the control fundamentally changes, as can be seen in Sections 3.2 and 3.4. Section 3.3 features a non-symmetric density. A more indicative example, although for a bounded state space, can be found in [27].

Before we delve into the examples, we give some insights on how the proposed method is implemented. The main idea of the implementation is to bring the discrete versions of our control problem into the standard form of a linear program, which can be solved with commonly available libraries.

#### 3.1. Details of the implementation

An implementation of the method proposed in this article needs to compute the components of a linear program as required by typical solvers, and then pass them into such a solver for computation. Some short post-processing of the results reveals the optimal control of our discretized problem. For the examples shown in this section, MATLAB<sup>®</sup> was used. Its linprog

function expects the following arguments, where *x* represents a possible solution to the linear program:

- the cost vector f codifying the objective function, which is evaluated by computing  $f^{\top}x$ ;
- the matrix  $A^{(1)}$  and the vector  $b^{(1)}$  codifying the equality constraints  $A^{(1)}x = b^{(1)}$ ;
- the matrix  $A^{(2)}$  and the vector  $b^{(2)}$  codifying possible inequality constraints  $A^{(2)}x \le b^{(2)}$ ;
- two vectors  $l_b$  and  $u_b$  representing the lower and upper bounds on *x*, respectively, such that  $l_b \le x \le u_b$  holds component-wise.

While other packages might demand a slightly different interface, we consider MATLAB's interface general enough to express the ideas of this section in terms of this form of linear program. To compute  $A^{(1)}$  and  $b^{(1)}$ , we need to take a closer look at the discretized measure  $\hat{\mu}_m$ . Recalling the forms of  $\hat{p}_m$  and  $\hat{\eta}_m$  as in (2.2) and (2.3), we can derive the following formula:

$$\hat{\mu}_{m}(F \times V) = \int_{F} \left[ \sum_{k=0}^{2^{m}-1} \left( \sum_{j=0}^{2^{m}-1} \sum_{i=0}^{2^{k_{m}}} \beta_{j,i} I_{E_{j}}(x) \delta_{\{u_{i}\}}(V) \right) \gamma_{k} p_{k}(x) \right] dx$$
$$= \int_{F} \left[ \sum_{k=0}^{2^{m}-1} \sum_{j=0}^{2^{m}-1} \sum_{i=0}^{2^{k_{m}}} \underbrace{\gamma_{k} \beta_{j,i}}_{x_{k,j,i}} I_{E_{j}}(x) \delta_{\{u_{i}\}}(V) p_{k}(x) \right] dx.$$
(3.1)

To form a 'proper' row vector, the vector  $x_{k,j,i} \equiv \gamma_k \beta_{k,j,i}$  can be re-indexed. For simplicity, we will retain the 'triple' index notation with indices k, j, i in the following. Wherever they appear, it is understood that multiplication of matrices and vectors with x can be expressed by the standard dot product, summing over one index. It may appear as if the multiplication of  $\gamma_k$  and  $\beta_{j,i}$  would make the problem non-linear. However, it is possible to solve for the vector x and compute  $\beta_{j,i}$  from it, which is our variable of interest as it encodes the optimal control.

The computation of the cost vector f in (3.1) is done as follows. Observe that

$$\int_{E \times U} c \, \mathrm{d}\hat{\mu}_m = \sum_{k=0}^{2^m - 1} \sum_{j=0}^{2^m - 1} \sum_{i=0}^{2^{k_m}} \underbrace{\gamma_k \beta_{j,i}}_{x_{k,j,i}} \underbrace{\int_{E_j} c(x, u_i) p_k(x) \, \mathrm{d}x}_{f_{k,j,i}} = f^\top x.$$

Hence, finding the components of f comes down to integration against Lebesgue measure. As long as c is a polynomial (as it is in all the examples considered in this article), Gaussian quadrature can be used to obtain exact results. As the support of  $p_k$  is disjoint from  $E_j$  for most indices k and j, only a small number of components of f have to be actually computed, and sparse structures lend themselves to storing f efficiently.

Having considered the computation of f, the computation of  $A^{(1)}$  can be run in similar fashion. Let  $\{f_r\}_{r=-3}^{2^{q_2}-1}$  be the finite basis of  $\mathscr{D}_{(q_1,q_2)}((-K, K))$  as used in (2.1). Then,

$$\int_{E \times U} Af_r \, \mathrm{d}\hat{\mu}_m = \sum_{k=0}^{2^m - 1} \sum_{j=0}^{2^m - 1} \sum_{i=0}^{2^{k_m}} \underbrace{\gamma_k \beta_{j,i}}_{x_{k,j,i}} \underbrace{\int_{E_j} Af_r(x, u_i) p_k(x) \, \mathrm{d}x}_{A_{r,(k,j,i)}^{(1)}}.$$

Using this formula, all components of  $A^{(1)}$  can be computed exactly using numerical integration, as  $Af_r$  is piecewise polynomial. The sparsity of  $A^{(1)}$  can be accounted for for the same reasons mentioned with respect to the computation of f. Also, for  $r = -3, \ldots, 2^{q_2} - 1$ , we set  $b_r^{(1)} = Rf_r$ , which is 0 in the case of the long-term average cost criterion, or a simple function evaluation in the case of the infinite horizon discounted cost criterion. To ensure that  $\hat{\mu}_m$ is indeed a probability measure, we set  $A_{2^{q_2}-1,(k,j,i)}^{(1)} = \int_{E_j} p_k(x) dx$  and  $b_{2^{q_2}-1}^{(1)} = 1$ . Additional rows are appended to  $A^{(1)}$  and  $b^{(1)}$  if an additional equality constraint is present; cf. (1.8). Similarly,  $A^{(2)}$  and  $b^{(2)}$  are defined by the inequality constraint of (1.8), if present. The same integration techniques can be applied for these constraints. The lower bound  $l_b$  on x is set to be zero, ensuring that all coefficients, and, by extension, the considered probability measures and densities, are non-negative. The upper bound  $u_b$  is not needed in this context. Putting all the elements together, we can use a linear programming package to solve

Minimize 
$$f^{\top}x$$
 subject to 
$$\begin{cases} A^{(1)}x = b^{(1)}, \\ A^{(2)}x \le b^{(2)}, \\ 0 \le x. \end{cases}$$

The dimensions of  $A^{(1)}$ ,  $b^{(1)}$ , f, and x depend on the choice of discretization parameters  $q_2$ , m, and  $k_m$ , as seen above. Typically,  $A^{(1)}$  features several thousand rows. As an additional optimization constraint contributes a single line to  $A^{(1)}$  (or  $A^{(2)}$ ), the added complexity of additional optimization constraints is relatively small.

We previously defined  $\hat{p}_m$ , see (2.2), by using indicator functions of the intervals of the dyadic partition of *E*. This approximation is also used in the proofs of Section 4. For the numerical examples, we adapt this approach by splitting several of those intervals to introduce additional degrees of freedom to the problem. This became necessary due to the observation that, for some examples, clean numerical solutions could only be obtained when keeping  $q_2 = m$ . However, this can lead to the number of constraints outnumbering the degrees of freedom due to the structure of the B-spline basis functions and indicator functions. Furthermore, additional constraints are present, either ensuring that  $\hat{p}_m$  is a probability density, or ensuring the fulfillment of additional optimization constraints. Splitting the intervals counters this issue. The theoretical results of Section 4 nevertheless remain valid. Once the linear program is solved for *x*, we can use the equation

$$\hat{\mu}_m(E_j, \{u_i\}) = \sum_{k=0}^{2^m - 1} \int_{E_j} \gamma_k \beta_{j,i} \hat{p}_k(x) \, \mathrm{d}x = \beta_{j,i} \cdot \hat{\mu}_m(E_j, U),$$

which can easily be solved for  $\beta_{j,i}$ , to obtain the coefficients of the regular conditional probability  $\hat{\eta}_m$ ; again, see (2.3). This ultimately gives us the actual behavior of the control we are interested in finding.

#### 3.2. Example 1: Controlled Brownian motion process

Consider a drifted Brownian motion process, with the drift entirely defined by the control u, featuring a constant diffusion coefficient  $\sigma$ . As a cost criterion, we adopt the long-term average criterion with  $\tilde{c}(x, u) = x^2$ . The resulting SDE and cost criterion J are

$$dX_t = u_t dt + \sigma dW_t, \quad X_0 = 0; \qquad J \equiv \limsup_{t \to \infty} \frac{1}{t} \mathbb{E} \bigg[ \int_0^t X_s^2 ds \bigg].$$



TABLE 1. Configuration of problem, discretization parameters, and cost criterion values for a controlled Brownian motion process.

FIGURE 1. Density of state-space marginal, controlled Brownian motion process.

We have  $E = (-\infty, \infty)$ , and set U = [-1, 1]. Note that the assumptions of our method, especially in the light of Remark 1.5, are satisfied, as we are considering a simple drifted Brownian motion with polynomial costs. Basic principles of optimal control, in particular the linear influence of the control on both SDE and cost criterion, dictate that the optimal control is  $u_t = u(X_t) = -\text{sgn}(X_t) \cdot I_{X_t \neq 0}$ , a so-called bang-bang control. It is a 'strict' control a.e., meaning that it assigns full mass on single points, as opposed to using the flexibility of a relaxed control. We can compute that the respective density is  $p(x) = 1/\sigma^2 \cdot \exp(-|x| \cdot 2/\sigma^2)$  by using the method in the proof of [12, Proposition 3.1]. The configuration of the problem and the numerical solution are shown in Table 1. For this example, we can compare the exact value of the cost criterion for the optimal solution  $J^*$  with the numeric approximation  $\hat{J}$ . Note that even though the expected optimal control is a bang-bang control, we allow the solver to choose from  $2^{10}$  different control values in order to not implicitly assume the form of the optimal control. For this configuration, we have  $|J^* - \hat{J}| \approx 4.8 \cdot 10^{-5}$ , while errors on the magnitude of  $10^{-7}$  can be achieved by increasing  $q_1, q_2$ , and m.

The approximations for p and the average control value  $x \mapsto \int_U \eta(u, x) du$  are displayed in Figures 1 and 2, respectively. Note that a plot of the average control value is preferable over plotting a three-dimensional representation of  $\eta$  with  $2^{10} \cdot 2^{10}$  points, for the sake of presentation. Concerning the state-space density, we can clearly see that it is concentrated around the origin, as the process is pushed towards it by the control. The average control value indeed takes the form of the expected bang-bang control, with some numerical artifacts at the boundary of the state space. These boundaries are imposed by the discretization, and do not exist in the original problem. Given the very small values of the state-space density at the fringes of the state space, the influence of these artifacts on the optimal value  $\hat{J}$  can be neglected. However,



FIGURE 2. Average of optimal relaxed control, controlled Brownian motion process.

they indicate that special treatment of the discretization at the borders of the state space might be necessary. This is the subject of current research.

#### 3.3. Example 2: Ornstein–Uhlenbeck process with costs of control

Consider a controlled Ornstein–Uhlenbeck process. Ornstein–Uhlenbeck processes feature a 'mean-reverting' drift that automatically pushes the process X back to its mean v. The strength of that push is determined by a coefficient  $\rho$ . The diffusion coefficient  $\sigma$  is constant. As a cost criterion, we use the infinite horizon discounted criterion with  $\tilde{c}(x, u) = x^2 + u^2$ . In contrast to Section 3.2, this introduces a cost for using the control, which in turn prevents the optimal control from being of bang-bang type. This is because it balances the costs induced by the position of X with the costs induced by using the control. The linear influence of the control on both SDE and cost criterion implies that it is a 'strict' control a.e., meaning that it assigns full mass on single points, as opposed to using the flexibility of a relaxed control. The SDE and cost criterion J considered are

$$dX_t = [\rho(\nu - X_t) + u_t] dt + \sigma dW_t, \quad X_0 = x_0; \qquad J \equiv \mathbb{E}\bigg[\int_0^\infty e^{-\alpha s} ((X_s - \nu)^2 + u_s^2) ds\bigg].$$

We have  $E = (-\infty, \infty)$ , and set U = [-1, 1]. The assumptions of our method, especially in the light of Remark 1.5 (given the mean-reverting behavior of the model), are satisfied. The configuration of the problem and the numerical solution are shown in Table 2. We have chosen a rather high discounting factor  $\alpha$  whose influence on the results will be visible. The proposed method also works well with smaller discounting factors.

The approximations for the state space density p and the average control value  $x \mapsto \int_U \eta(u, x) du$  are displayed in Figures 3 and 4, respectively. p has its mode at the origin, which it is pushed towards by both the mean-reverting drift and the control. The influence of the discounting factor and the starting point  $x_0 = 1.5$  are clearly visible. As 'earlier' values of X are weighted higher by the discounted expected occupation measure when considering a discounted cost criterion, its state-space density p shows a feature at  $x_0$ . Again, the average control value shows numerical artifacts at the boundary of the state space.

<i>x</i> <sub>0</sub>	α	ρ	ν	σ	K	$q_1$	$q_2$	т	k <sub>m</sub>	$\hat{J}$
1.5	0.2	0.1	0	$\sqrt{2}/2$	4	12	10	10	10	0.393 46
		C	<sup>).8</sup> [						-	
		С	0.6			$\wedge$				
		ity.	).4			/				
		dens	0.2			1				
							ſ			
			-4	-2		<u>.</u>	<u>.</u>			

TABLE 2. Configuration of problem and discretization parameters for an Ornstein–Uhlenbeck process with costs of control.

FIGURE 3. Density of state-space marginal, Ornstein-Uhlenbeck process with costs of control.



FIGURE 4. Average of optimal relaxed control, Ornstein-Uhlenbeck process with costs of control.

# 3.4. Example 3: Ornstein–Uhlenbeck process with control budget constraint

Again consider a controlled Ornstein–Uhlenbeck process. In contrast to Section 3.3, we will use the absolute value of the deviation of X from  $\nu$  as the cost function with the long-term average cost criterion and introduce a constraint in the style of (1.8) to the problem:

TABLE 3. Configuration of problem, discretization parameters, and approximate cost criterion value for an Ornstein–Uhlenbeck process with a control budget.

<i>x</i> <sub>0</sub>	ρ	ν	σ	$D_1$	K	$q_1$	$q_2$	т	k <sub>m</sub>	$\hat{J}_{ m c}$	$\hat{J}_{\mathrm{u}}$
0	0.1	0	$\sqrt{2}/2$	0.15	5	14	13	13	5	0.632 62	0.238 83



FIGURE 5. Density of state-space marginal, Ornstein–Uhlenbeck process with control budget.

$$dX_{t} = \left[\rho(\nu - X_{t}) + u_{t}\right] dt + \sigma \ dW_{t}, \quad X_{0} = x_{0};$$

$$\lim_{t \to \infty} \sup_{t} \frac{1}{t} \mathbb{E} \left[ \int_{0}^{t} |u_{s}| \ ds \right] \leq D_{1};$$

$$J \equiv \limsup_{t \to \infty} \frac{1}{t} \mathbb{E} \left[ \int_{0}^{t} |X_{s} - \nu| \ ds \right].$$
(3.2)

In this setting, the constraint in (3.2) can be interpreted in such a way that there are costs associated with using the control (given by the absolute value of the control), and these costs must not exceed a certain budget stream of  $D_1$  per unit of time in the long-term average. The assumptions of our method are satisfied, in particular as U is compact and the cost criterion will therefore be finite. To show the influence of the additional constraint on the optimal solution, we compare the solution for this problem to the same problem without this constraint, i.e. where (3.2) is not present.

As in Section 3.4, we have  $E = (-\infty, \infty)$  and set U = [-1, 1]. The configuration of the problem and the numerical solution are shown in Table 3.  $\hat{J}_c$  denotes the approximate solution to the constrained problem, while  $\hat{J}_u$  denotes the approximate solution to the unconstrained problem. In order to attain a smooth solution,  $q_1$  and  $q_2$  must be larger than in the previous examples. To make up for the additional memory needed by this discretization, we choose a smaller value of 5 for  $k_m$ , which is nevertheless large enough to allow the solver to consider non-bang-bang controls. We thereby avoid an implicit assumption on the form of the optimal control.

The approximation for *p* (with the budget constraint present) is shown in Figure 5. Figure 6 shows the average control value  $x \mapsto \int_{U} \eta(u, x) du$  for the case where the budget constraint is



FIGURE 6. Average of optimal relaxed control, Ornstein–Uhlenbeck process with control budget.

present, in black, and for the case when it is not present, in gray. With regard to the optimal control value, we can see that due to the limited budget available, the solver chooses not to act when the process X is rather close to the origin, accruing lower costs, but acts at 'full force' once the process is more than 1.075 units away from v. Again, we see some numerical artifacts towards the borders of the computed state space. If the budget constraint is not present, the control acts at 'full force' at any position of the state space, as expected. As a result of the specific behavior of the control, the associated state space density is rather 'flat' in the interval where the control is not acting, and rather steep outside it. The numerical values of  $\hat{J}_c$  and  $\hat{J}_u$  show the expected behavior. Due to the limited control usage when the budget constraint is present,  $\hat{J}_c$  is larger than its counterpart  $\hat{J}_u$ .

### 4. Proofs

# 4.1. Addressing the unboundedness of the state space

This section discusses the proof of Proposition 2.1. We begin by considering weakly convergent sequences of measures  $\{\mu_K\}_{K \in \mathbb{N}}$  that lie in  $\mathcal{M}_{\infty,K}$  for each  $K \in \mathbb{N}$ .

**Lemma 4.1.** Let  $\{\mu_K\}_{K\in\mathbb{N}}$  satisfy  $\mu_K \in \mathcal{M}_{\infty,K}$  for all  $K \in \mathbb{N}$ . Assume that  $\mu_K \Rightarrow \hat{\mu}$  for some  $\hat{\mu} \in \mathcal{P}(E \times U)$  as  $K \to \infty$ . Then  $\hat{\mu} \in \mathcal{M}_{\infty}$ .

*Proof.* Take  $f \in \mathscr{D}_{\infty}(E)$ . Then there is a  $K_0 \in \mathbb{N}$  such that  $\operatorname{supp}(f) \subset (-K, K)$  for all  $K \ge K_0$ . Also, the function  $(x, u) \mapsto Af(x, u)$  lies in  $C_b^u(E \times U)$ , given the form of A, cf. (1.6), and since f, f', and f'' have compact support. By the weak convergence of measures,  $\int_{E \times U} Af d\hat{\mu} = \lim_{K \to \infty} \int_{E \times U} Af d\mu_K = \lim_{K \to \infty} Rf = Rf$ .

**Remark 4.1.** If constraints in the form of (1.8) are present, Lemma 4.1 still holds since  $g_1$  is non-negative and continuous (hence bounded over any compact interval), and

$$\int_{E \times U} g_1 \, \mathrm{d}\mu = \lim_{L \to \infty} \int_{[-L,L] \times U} g_1 \, \mathrm{d}\mu = \lim_{L \to \infty} \lim_{K \to \infty} \int_{[-L,L] \times U} g_1 \, \mathrm{d}\mu_K \le D_1$$

holds by monotone convergence and, again, the weak convergence of measures. The same arguments hold for the equality constraint of (1.8).

A crucial part of the convergence analysis is to consider, given a sequence of measures  $\{\mu_K\}_{K\in\mathbb{N}}$ , how the sequence of values given by  $\{J(\mu_K)\}_{K\in\mathbb{N}}$  evolves. This is discussed next. Note that we use the notation  $J(\mu) = \int_{E\times U} c \, d\mu$  interchangeably.

**Lemma 4.2.** Let  $\{\mu_K\}_{K \in \mathbb{N}}$  satisfy  $\mu_K \in \mathcal{M}_{\infty,K}$  for all  $K \in \mathbb{N}$ . Assume that  $\mu_K \Rightarrow \hat{\mu}$  for some  $\hat{\mu} \in \mathcal{M}_{\infty}$  as  $K \to \infty$ . Let  $\mu \in \mathcal{M}_{\infty}$  be another measure, and assume that, for some  $\varepsilon > 0$ ,  $J(\hat{\mu}) > J(\mu) + \varepsilon$  holds. Then, there is a  $K_0 \in \mathbb{N}$  large enough such that, for all  $K \ge K_0$ ,  $J(\mu_K) > J(\mu) + \varepsilon$ .

*Proof.* By monotone convergence, there exists an  $L_1$  large enough such that, for all  $L_2, L_3 \ge L_1$ ,

$$\int_{E \times U} c \, \mathrm{d}\hat{\mu} \ge \int_{[-L_2, L_3] \times U} c \, \mathrm{d}\hat{\mu} > \int_{E \times U} c \, \mathrm{d}\mu + \varepsilon$$

is true. Find  $L_2$ ,  $L_3 \ge L_1$  such that  $c(x, u) > c(-L_2, u)$  for all  $x < -L_2$  and  $c(x, u) > c(L_3, u)$  for all  $x > L_3$ , which is possible since *c* is increasing in |x|; cf. Assumption 1.1(iii). Define

$$\bar{c}(x) = \begin{cases} c(-L_2, u) & \text{if } x < -L_2, \\ c(x, u) & \text{if } x \in [-L_2, L_3], \\ c(L_3, u) & \text{if } x > L_3. \end{cases}$$

Observe that  $\bar{c}$  is uniformly continuous and bounded, but also

$$\int_{E \times U} c \, \mathrm{d}\hat{\mu} > \int_{E \times U} \bar{c} \, \mathrm{d}\hat{\mu} \ge \int_{[-L_2, L_3] \times U} c \, \mathrm{d}\hat{\mu} > \int_{E \times U} c \, \mathrm{d}\mu + \varepsilon.$$
(4.1)

Clearly, for any  $K \in \mathbb{N}$  we have  $\int_{E \times U} c \, d\mu_K \ge \int_{E \times U} \bar{c} \, d\mu_K$ . On the other hand, by weak convergence,  $\int_{E \times U} \bar{c} \, d\mu_K \to \int_{E \times U} \bar{c} \, d\hat{\mu}$  as  $K \to \infty$ , and with (4.1) there is a  $K_0$  large enough such that, for all  $K \ge K_0$ ,

$$\int_{E \times U} c \, \mathrm{d}\mu_K \ge \int_{E \times U} \bar{c} \, \mathrm{d}\mu_K > \int_{E \times U} c \, \mathrm{d}\mu + \varepsilon.$$

We turn to weakly convergent sequences which are  $\varepsilon$ -optimal for some arbitrary  $\varepsilon > 0$ . In other words, we consider  $\{\mu_K^{\varepsilon}\}_{K \in \mathbb{N}}$ , with  $\mu_K^{\varepsilon} \in \mathcal{M}_{\infty,K}$  for all  $K \in \mathbb{N}$ , and, for each  $K \in \mathbb{N}$ ,  $J(\mu_K^{\varepsilon}) < J(\mu_K) + \varepsilon$  holding for any  $\mu_K \in \mathcal{M}_{\infty,K}$ , in accordance with Definition 2.1. Note that, due to Assumption 1.1(iv) and the fact that  $J(\mu) \ge 0$  for any  $\mu$ , we can assume the existence of an  $\varepsilon$ -optimal solution in  $\mathcal{M}_{\infty}$  with finite cost. Since  $\mathcal{M}_{\infty} \subset \mathcal{M}_{\infty,K}$ , the existence of an  $\varepsilon$ -optimal solution in  $\mathcal{M}_{\infty,K}$  follows, and it is valid to consider sequences of measures with the properties described. The following argument considers the tightness of such sequences. As U is a bounded set, our analysis solely focuses on the state-space behavior of the measures.

**Lemma 4.3.** For each  $K \in \mathbb{N}$ , assume that  $\mu_K^{\varepsilon} \in \mathcal{M}_{\infty,K}$  and that  $\mu_K^{\varepsilon}$  is an  $\varepsilon$ -optimal solution. Then,  $\{\mu_K^{\varepsilon}\}_{K \in \mathbb{N}}$  is tight.

*Proof.* Assume the opposite. Then there exists a  $\delta > 0$  such that, for all  $L \in \mathbb{N}$ , there exists a  $K \in \mathbb{N}$  with  $\mu_K^{\varepsilon}([-L, L]^{\sim} \times U) \ge \delta$ . Using that *c* is increasing in |x|, choose *L* large enough such that  $c(x, u) > (\int_{E \times U} c \, d\mu + \varepsilon) \cdot (1/\delta)$  for all  $x \in [-L, L]^{\sim}$ , uniformly in *u*, where  $\mu$  is a measure in  $\mathscr{M}_{\infty,K}$  with finite costs. By the hypotheses,  $\mu_K^{\varepsilon}([-L, L]^{\sim} \times U) \ge \delta$  for some  $K \in \mathbb{N}$ , and hence

Finite element approximations for stochastic control in unbounded state space

$$\int_{E \times U} c \, \mathrm{d}\mu_K^{\varepsilon} \ge \int_{[-L,L]^{\sim} \times U} c \, \mathrm{d}\mu_K^{\varepsilon} > \delta \cdot \left(\int_{E \times U} c \, \mathrm{d}\mu + \varepsilon\right) \cdot \frac{1}{\delta} = \int_{E \times U} c \, \mathrm{d}\mu + \varepsilon,$$

which is a contradiction to  $\mu_K^{\varepsilon}$  being  $\varepsilon$ -optimal.

By Theorem 1.3, a tight sequence of measures contains a convergent subsequence. The limit of a convergent sequence of  $\varepsilon$ -optimal measures is  $\varepsilon$ -optimal itself.

**Lemma 4.4.** For each  $K \in \mathbb{N}$ , assume that  $\mu_K^{\varepsilon} \in \mathcal{M}_{\infty,K}$  and that  $\mu_K^{\varepsilon}$  is an  $\varepsilon$ -optimal solution in  $\mathcal{M}_{\infty,K}$ . Assume that  $\mu_K^{\varepsilon} \Rightarrow \hat{\mu}$  for some  $\hat{\mu} \in \mathcal{P}(E \times U)$ . Then  $\hat{\mu}$  is  $\varepsilon$ -optimal in  $\mathcal{M}_{\infty}$ .

*Proof.* Assume that  $\hat{\mu}$  is not  $\varepsilon$ -optimal. Then there exists a  $\mu \in \mathscr{M}_{\infty}$  such that  $\int_{E \times U} c \, d\hat{\mu} > \int_{E \times U} c \, d\mu + \varepsilon$ . By Lemma 4.2, there exists a  $K \in \mathbb{N}$  such that  $\int_{E \times U} c \, d\mu_{K}^{\varepsilon} > \int_{E \times U} c \, d\mu + \varepsilon$ , which contradicts the assumption that  $\mu_{K}^{\varepsilon}$  is  $\varepsilon$ -optimal, as  $\mu \in \mathscr{M}_{\infty,K}$  for all K.

In general, weak convergence cannot be assumed. The following result investigates how sequences of  $\varepsilon$ -optimal measures behave without such an assumption.

**Lemma 4.5.** For each  $K \in \mathbb{N}$ , assume that  $\mu_K^{\varepsilon} \in \mathscr{M}_{\infty,K}$  and that  $\mu_K^{\varepsilon}$  is an  $\varepsilon$ -optimal solution in  $\mathscr{M}_{\infty,K}$ . Then, for any  $\delta > 0$ , there is a  $z \in \mathbb{R}$  and a  $K_0(\delta) \in \mathbb{N}$  such that  $J(\mu_K^{\varepsilon}) \in (z - \varepsilon/2 - \delta, z + \varepsilon/2 + \delta)$  for all  $K \ge K_0(\delta)$ .

*Proof.* Consider two convergent subsequences of  $\{\mu_K^{\varepsilon}\}_{K\in\mathbb{N}}$ , denoted  $\{\mu_{k_j}^{\varepsilon}\}_{j\in\mathbb{N}}$  and  $\{\mu_{l_j}^{\varepsilon}\}_{j\in\mathbb{N}}$ . Let  $\hat{\mu}$  and  $\tilde{\mu}$  be their respective limits, and assume that  $\hat{\mu} \neq \tilde{\mu}$ . Assume that  $\int_{E\times U} c \, d\hat{\mu} > \int_{E\times U} c \, d\tilde{\mu} + \varepsilon$ . By Lemma 4.2, there exists an  $N \in \mathbb{N}$  large enough such that, for all  $j \ge N$ ,  $\int_{E\times U} c \, d\mu_{k_j} > \int_{E\times U} c \, d\tilde{\mu} + \varepsilon$ , contradicting that  $\{\mu_{k_j}^{\varepsilon}\}_{j\in\mathbb{N}}$  is a sequence of  $\varepsilon$ -optimal measures. Hence,  $\int_{E\times U} c \, d\hat{\mu} \le \int_{E\times U} c \, d\tilde{\mu} + \varepsilon$  has to be true. Applying the same argument to  $\{\mu_{l_j}^{\varepsilon}\}_{j\in\mathbb{N}}$  and  $\hat{\mu}$ , we can readily conclude that  $\int_{E\times U} c \, d\tilde{\mu} \le \int_{E\times U} c \, d\hat{\mu} + \varepsilon$ . Both results put together reveal that  $\left|\int_{E\times U} c \, d\hat{\mu} - \int_{E\times U} c \, d\tilde{\mu}\right| \le \varepsilon$ , and hence there is a  $z \in \mathbb{R}$  such that  $J(\mu) \in [z - \varepsilon/2, z + \varepsilon/2]$  for any limit  $\mu$  of a convergent subsequence of  $\{\mu_K^{\varepsilon}\}_{K\in\mathbb{N}}$ .

Fix  $\delta > 0$ , and assume there is a non-convergent subsequence  $\{\mu_{m_j}^{\varepsilon}\}_{j \in \mathbb{N}}$  of  $\{\mu_K^{\varepsilon}\}_{K \in \mathbb{N}}$  such that, for any N, there is a  $j \ge N$  with  $\int_{E \times U} c \, d\mu_{m_j} \notin (z - \varepsilon/2 - \delta, z + \varepsilon/2 + \delta)$ . Then, there exists a sub-subsequence  $\{\mu_{m'_j}^{\varepsilon}\}_{j \in \mathbb{N}}$  with  $\int_{E \times U} c \, d\mu_{m'_j} \notin (z - \varepsilon/2 - \delta, z + \varepsilon/2 + \delta)$  for all  $j \in \mathbb{N}$ . This sub-subsequence, however, remains a sequence of  $\varepsilon$ -optimal measures, as it is a subsequence of  $\{\mu_K^{\varepsilon}\}_{K \in \mathbb{N}}$ , and hence is tight by Lemma 4.3. Theorem 1.3 implies the existence of a convergent 'sub-sub'-subsequence  $\{\mu_{m'_j}^{\varepsilon}\}_{j \in \mathbb{N}}$ , for which  $\lim_{j \to \infty} J(\mu_{m''_j}) \in [z - \varepsilon/2, z + \varepsilon/2]$ , contradicting the existence of  $\{\mu_{m_j}^{\varepsilon}\}_{j \in \mathbb{N}}$ , and proving the claim.

Lemma 4.5 not only reveals the location of limiting values when considering  $\varepsilon$ -optimal sequences, it also lets us form a conclusion about the optimal value for J.

**Lemma 4.6.** Set  $J^* = \inf\{J(\mu): \mu \in \mathcal{M}_{\infty}\}$ , and consider the sequence  $\{\mu_K^{\varepsilon}\}_{K \in \mathbb{N}}$  of  $\varepsilon$ -optimal measures, as well as  $z \in \mathbb{E}$ , from Lemma 4.5. Then  $z - 3\varepsilon/2 \le J^* \le z + \varepsilon/2$ .

*Proof.* Let  $\hat{\mu}$  be the limit of a convergent subsequence of  $\{\mu_K^{\varepsilon}\}_{K \in \mathbb{N}}$ . As seen in the proof of Lemma 4.6,  $J(\hat{\mu}) \in [z - \varepsilon/2, z + \varepsilon/2]$ , and thereby  $J^* \leq J(\hat{\mu}) \leq z + \varepsilon/2$ . But on the other hand,  $\hat{\mu}$  is  $\varepsilon$ -optimal by Lemma 4.4, and thereby

$$J^* + \varepsilon \ge J(\hat{\mu}) \ge z - \varepsilon/2 \iff J^* \ge z - 3\varepsilon/2.$$

**Proposition 4.1.** For each  $K \in \mathbb{N}$ , assume that  $\mu_K^{\varepsilon} \in \mathcal{M}_{\infty,[-K,K]}$  and that  $\mu_K^{\varepsilon}$  is an  $\varepsilon$ -optimal solution in  $\mathcal{M}_{\infty,K}$ . Then, for  $\delta > 0$ , there exists a  $K_0(\delta)$  such that  $|J(\mu_K^{\varepsilon}) - J^*| \le 2\varepsilon + \delta$  for all  $K \ge K_0(\delta)$ .

*Proof.* Fix  $\delta > 0$ , and choose  $K_0$  large enough such that, for all  $K \ge K_0$ ,  $J(\mu_K) \in (z - \varepsilon/2 - \delta, z + \varepsilon/2 + \delta)$  holds by Lemma 4.5. Using Lemma 4.6, we deduce that

$$z - \frac{\varepsilon}{2} - \delta - \left(z + \frac{\varepsilon}{2}\right) \le J(\mu_K) - J^* \le z + \frac{\varepsilon}{2} + \delta - \left(z - \frac{3\varepsilon}{2}\right).$$

Proposition 4.1 reduces the problem of finding  $\varepsilon$ -optimal solutions in  $\mathcal{M}_{\infty}$  to finding optimal solutions in  $\mathcal{M}_{\infty,K}$ . While this is significant progress in terms of attaining a computable formulation, the fact that measures in  $\mathcal{M}_{\infty,K}$  are allowed to have mass anywhere in E still has to be addressed. The following analysis relies on our assumption that c allows for compactification, cf. Assumption 1.1, and the existence of solutions for problems with bounded state space, cf. Theorem 1.2.

**Lemma 4.7.** Let  $\mu_K$  be a measure in  $\mathcal{M}_{\infty,K}$ . There exists a measure  $\tilde{\mu}_K \in \mathcal{M}_{\infty,K}$  with  $J(\tilde{\mu}_K) \leq J(\mu_K)$ .

*Proof.* If  $\mu_K([-K, K] \times U) = 1$ , the statement is trivially true as  $\mu_K \in \mathcal{M}_{\infty,K}$  in this case. So, assume that  $\tau := \mu_K([-K, K] \times U) < 1$ , and let  $u^-$  be a continuous function such that  $\sup\{c(x, u^-(x)): x \in [-K, K]\} \le \inf\{c(x, u): x \in [-K, K]^{\sim}, u \in U\}$ . By Theorem 1.2 we can take a solution  $(\hat{\mu}_0, \hat{\mu}_1)$  to the constraints of the singular linear program with a bounded state space E = [-K, K], cf. (1.10), with  $x_0 = K$  and reflections at both ends of the state space  $\{-K\}$  and  $\{K\}$ , under a control satisfying  $\eta_0(u^-(x), x) = 1$ . By Remark 2.2,  $\int A_{E \times U} f d\hat{\mu}_0 = 0$  for all  $f \in \mathcal{D}_{\infty}((-K, K)) \equiv C_c^2((-K, K))$ . Set  $\tilde{\mu}_K^A = (1 - \tau)\hat{\mu}_0$  and, for  $F \times V \in \mathcal{B}(E \times U)$ , define  $\tilde{\mu}_K(F \times V) = \mu_K(F \cap [-K, K] \times V) + \tilde{\mu}_K^A(F \times V)$ . As  $\hat{\mu}_0([-K, K] \times U) = 1$ , we have  $\tilde{\mu}_K([-K, K] \times U) = 1$ . For  $f \in \mathcal{D}_{\infty}((-K, K))$ ,

$$\int_{E \times U} Af \, \mathrm{d}\tilde{\mu}_K = \int_{E \times U} Af \, \mathrm{d}\tilde{\mu}_K^A + \int_{E \times U} Af \, \mathrm{d}\hat{\mu}_K = (1 - \tau) \cdot 0 + Rf = Rf$$

follows, so  $\tilde{\mu}_K \in \mathscr{M}_{\infty,K}$ . Note that  $c(x, u) = c(x, u^-(x))$  holds  $\tilde{\mu}_K^A$ -a.e. in [-K, K] as, by construction,  $\tilde{\mu}_K^A(dx, \cdot)$  has full mass on  $\{(x, u^-(x)), x \in [-K, K]\}$ . Therefore,

$$\int_{E \times U} c \, d\tilde{\mu}_K = \int_{[-K,K] \times U} c \, d\tilde{\mu}_K^A + \int_{[-K,K] \times U} c \, d\mu_K$$
$$\leq \int_{[-K,K] \times U} \sup\{c(x, u^-(x)) \colon x \in [-K,K]\} \, d\tilde{\mu}_K^A + \int_{[-K,K] \times U} c \, d\mu_K.$$

But c allows for compactification, proving that

$$\begin{split} \int_{[-K,K]\times U} \sup\{c(x, u^{-}(x)) \colon x \in [-K, K]\} \, \mathrm{d}\tilde{\mu}_{K}^{A} + \int_{[-K,K]\times U} c \, \mathrm{d}\mu_{K} \\ &\leq \tilde{\mu}_{K}^{A}([-K, K] \times U) \cdot \inf\{c(x, u) \colon x \in [-K, K]^{\sim}, \ u \in U\} + \int_{[-K,K]\times U} c \, \mathrm{d}\mu_{K} \\ &\leq (1-\tau) \cdot \inf\{c(x, u) \colon x \in [-K, K]^{\sim}, \ u \in U\} + \int_{[-K,K]\times U} c \, \mathrm{d}\mu_{K} \end{split}$$

$$\leq \mu_K([-K,K]^{\sim} \times U) \cdot \inf\{c(x,u) \colon x \in [-K,K]^{\sim}, \ u \in U\} + \int_{[-K,K] \times U} c \, \mathrm{d}\mu_K$$
$$\leq \int_{[-K,K]^{\sim} \times U} c(x,u) \, \mathrm{d}\mu_K + \int_{[-K,K] \times U} c \, \mathrm{d}\mu_K = \int_{E \times U} c \, \mathrm{d}\mu_K.$$

Using Lemma 4.7, we conclude that we can restrict ourselves to finding  $\varepsilon$ -optimal solutions in  $\mathcal{M}_{\infty,K}$ , as shown in the following result.

**Proposition 4.2.** An  $\varepsilon$ -optimal solution  $\mu_K^{\varepsilon}$  for  $\mathscr{M}_{\infty,K}$  is  $\varepsilon$ -optimal for  $\dot{\mathscr{M}}_{\infty,K}$ .

*Proof.* Assume the existence of  $\mu_K \in \mathcal{M}_{\infty,[-K,K]}$  with  $J(\mu_K) < J(\mu_K^{\varepsilon}) + \varepsilon$ . By Lemma 4.7, there is a measure  $\tilde{\mu}_K \in \mathcal{M}_{\infty,K}$  with  $J(\tilde{\mu}_K) \le J(\mu_K)$ , contradicting the  $\varepsilon$ -optimality of  $\mu_K^{\varepsilon}$ .  $\Box$ 

By combining Propositions 4.1 and 4.2, we conclude our effort in reducing the initial problem of finding an  $\varepsilon$ -optimal solution in  $\mathcal{M}_{\infty,K}$ .

Proof of Proposition 2.1. According to Proposition 4.2,  $\mu_K^{\varepsilon}$  is  $\varepsilon$ -optimal in  $\mathscr{M}_{\infty,K}$  for all K > 0. By Proposition 4.1, however, there is a  $K_0(\delta)$  such that, for all  $K \ge K_0(\delta)$ ,  $\mu_K^{\varepsilon}$  is  $(2\varepsilon + \delta)$ -optimal in  $\mathscr{M}_{\infty}$ .

#### 4.2. Discretizing the constraint space

The analysis needed to prove Proposition 2.2 bears a strong resemblance to that needed to prove Proposition 4.1. In both cases, we are considering  $\varepsilon$ -optimal solutions for a nested structure of sets,  $\mathcal{M}_{\infty} \subset \mathcal{M}_{\infty,K}, K \in \mathbb{N}$ , in the former case, and  $\mathcal{M}_{\infty,K} \subset \mathcal{M}_{q_n,K}, n \in \mathbb{N}$ , in the latter case. For that reason, we merely present the steps in which the proof differs, and refer to Section 4.1 for the remaining steps.

**Lemma 4.8.** Let  $\{\mu_n\}_{n\in\mathbb{N}}$  satisfy  $\mu_n \in \mathcal{M}_{q_n,K}$  for all  $n \in \mathbb{N}$ . Assume that  $\mu_n \Rightarrow \hat{\mu}$  for some  $\hat{\mu} \in \mathcal{P}(E \times U)$  as  $n \to \infty$ . Then,  $\hat{\mu} \in \mathcal{M}_{\infty,K}$ .

*Proof.* Take  $f \in \mathscr{D}_{\infty}((-K, K))$ . By Proposition 1.2, there is a sequence  $\{g_k\}_{k \in \mathbb{N}}$  with  $g_k \in \mathscr{D}_{q_k}((-K, K))$  for all k such that  $g_k \to f$  in  $\mathscr{D}_{\infty}((-K, K))$  as  $k \to \infty$ . Following Remark 1.11, we have

$$\int_{E \times U} Af \, \mathrm{d}\hat{\mu} = \lim_{k \to \infty} \int_{E \times U} Ag_k \, \mathrm{d}\hat{\mu} = \lim_{k \to \infty} \lim_{n \to \infty} \int_{E \times U} Ag_k \, \mathrm{d}\mu_n = \lim_{k \to \infty} Rg_k = Rf. \qquad \Box$$

**Lemma 4.9.** Let  $\{\mu_n\}_{n\in\mathbb{N}}$  satisfy  $\mu_n \in \mathcal{M}_{q_n,K}$  for all  $n \in \mathbb{N}$ . Assume that  $\mu_n \Rightarrow \hat{\mu}$  for some  $\hat{\mu} \in \mathcal{M}_{\infty,K}$  as  $n \to \infty$ . Let  $\mu \in \mathcal{M}_{\infty,K}$  be another measure, and assume that, for some  $\varepsilon > 0$ ,  $J(\hat{\mu}) > J(\mu) + \varepsilon$ . Then there is an  $N_0 \in \mathbb{N}$  large enough that, for all  $n \ge N_0$ ,  $J(\mu_n) > J(\mu) + \varepsilon$ .

*Proof.* This is an easy consequence of the fact that *c* is continuous, and hence uniformly continuous and bounded on the compact interval [-K, K].

By Remark 1.10, a sequence  $\{\mu_n\}_{n\in\mathbb{N}}$  of measures with  $\mu_n \in \mathcal{M}_{q_n,K}$  for all  $n \in \mathbb{N}$  is tight, as [-K, K] is compact. From here on, the statements of Lemmas 4.4, 4.5, and 4.6 can be proven for  $\varepsilon$ -optimal sequences  $\{\mu_n^\varepsilon\}_{n\in\mathbb{N}}$  of measures in  $\mathcal{M}_{q_n,K}$  with few modifications. The same holds for the final result of Proposition 4.1, which yields the proof of Proposition 2.2.

## 4.3. Discretizing the expected occupation measure

The proof for Proposition 2.3 is given by a simplified version of the proofs for [27, Corollary 3.3], where discretized problems with bounded state space are treated. The following result on the approximation quality of the discretizations introduced in (2.2) and (2.3) is its first important part. It is equivalent to [27, Proposition 3.2].

**Proposition 4.3.** For every  $\mu \in \mathcal{M}_{(q_1,q_2),K}$  and each  $\varepsilon > 0$ , there is an  $M_0(\varepsilon) \in \mathbb{N}$  such that, for all  $m \ge M_0$ , there exists a  $\hat{\mu}_m \in \mathcal{M}_{(q_1,q_2,m),K}$  with  $|J(\mu) - J(\hat{\mu}_m)| < \varepsilon$ .

*Proof.* The reductions described in Sections 2.1, 2.2, and 2.3 were specifically designed to obtain a problem structure that is akin to that of the discretized problem with bounded state space. In particular, we have achieved that  $\mathcal{M}_{(q_1,q_2,m),K}$  features a finite number of constraints and variables, and that both constraint functions and measures have their support contained in the compact interval [-K, K]. Although we have introduced an altered set  $\{h^{(q_1)} \cdot s_k^{(q_2)}\}_{k=-3}^{2^{q_2}-1}$  of constraint functions, contrary to the usage of 'plain' B-spline basis functions in [27], these functions remain twice continuously differentiable. Furthermore, *c*, *b*, and  $\sigma$  remain continuous functions. This allows for the analysis presented in [27, Section 3.2], relying on arguments of uniform bounds on continuous functions over a compact interval, to be carried out in an identical manner. Any analysis dealing with the singular behavior of *X* can be disregarded.

Continuing from this, we again need to consider a sequence of measures in  $\mathcal{M}_{(q_1,q_2,m),K}$  that converges weakly. To prove the following result, we can again employ the technique of the proofs of Lemmas 4.1 and 4.8.

**Lemma 4.10.** Let  $\{\mu_m\}_{m\in\mathbb{N}}$  satisfy  $\mu_m \in \mathcal{M}_{(q_1,q_2,m),K}$  for all  $m \in \mathbb{N}$ . Assume that  $\mu_m \Rightarrow \hat{\mu}$  for some  $\hat{\mu} \in \mathcal{P}(E \times U)$  as  $m \to \infty$ . Then  $\hat{\mu} \in \mathcal{M}_{(q_1,q_2),K}$ .

The final proof of Proposition 2.3 is now identical to that of [27, Corollary 3.3]. The first step is using Proposition 4.3 and Lemma 4.10 to show that if a sequence of optimal measures  $\{\mu_m^*\}_{m \in \mathbb{N}}$  with  $\mu_m^* \in \mathcal{M}_{(q_1,q_2,m),K}$  for all  $m \in \mathbb{N}$  converges weakly to a measure  $\mu^* \in \mathcal{M}_{(q_1,q_2),K}$ , then  $J(\mu^*) = J^*_{(q_1,q_2),K}$ . Then, we show that the real-valued sequence  $\{J(\mu_m^*)\}_{m \in \mathbb{N}}$  converges. Naturally, any of its subsequences converge to the same limit. As  $\{\mu_m^*\}_{m \in \mathbb{N}}$  is a tight sequence of measures, weakly convergent subsequences  $\{\mu_{m_i}^*\}_{j \in \mathbb{N}}$  exist, and  $\lim_{j\to\infty} J(\mu_{m_i}^*) = J^*_{(q_1,q_2),K}$ .

# 5. Outlook

This paper introduced a discretization scheme for stochastic control problems with unbounded state space, presented a convergence argument, and demonstrated its performance on three examples. Research on this topic can be continued in several directions. More attention has to be given to the assumptions needed for the convergence argument, in particular Assumption 1.3, where we required the state-space marginal of any solution considered to be absolutely continuous with respect to Lebesgue measure. Research on problems with bounded state space [28] has shown that Assumption 1.3 comes as a consequence of Assumption 1.2, with the latter posing rather soft, albeit technical, restrictions. Such results could possibly be extended to unbounded state spaces.

The behavior of the approximations to the relaxed control  $\eta$  at the boundaries of the computed state space [-K, K] has to be investigated further, in order to remove the numerical artifacts. A natural way of pursuing this would be by analyzing the basis functions spanning the

discrete constraint space. In particular, the function h used to ensure that these basis functions have compact support, cf. Definition 1.10, could be investigated further.

On top of that, the applicability of the proposed method to additional examples could be explored. As also pointed out in [27], it would be of great interest to consider the behavior of the proposed method when using higher-order approximations to the density p of  $\mu_E$ ; cf. (2.2). Typically, analytic solutions feature piecewise-smooth densities p, which could be hinting towards the possibility of better performance when using higher-order approximations.

Considering that this paper introduced a generalization to previous research in [27] by allowing for unbounded *state* spaces, it would be natural to investigate whether the method could also be generalized to unbounded *control* spaces. As a first step, we would have to establish conditions which guarantee the existence of optimal solutions. As can be deduced from the form of the optimal solution in Section 3.2, which is a bang-bang control, an optimal solution does not exist if no costs of control are present. Indeed, these controls always take the extreme values of the state space, and these extreme values do not exist in the unbounded case. Contrarily, Section 3.3 indicates that optimal solutions could exist if costs of controls are present. However, such an optimal solution could grow beyond any boundary, and would hence need to be approximated by a control that takes values in a bounded set to make them computationally accessible. Therefore, in a second step, the convergence of these approximations would have to be analyzed. Crucially, the existing convergence arguments might have to be adapted to accommodate the fact that the solution space no longer necessarily includes the optimal solution.

#### **Funding information**

There are no funding bodies to thank relating to the creation of this article.

#### **Competing interests**

There were no competing interests to declare which arose during the preparation or publication process of this article.

#### References

- ANSELMI, J., DUFOUR, F. AND PRIETO-RUMEAU, T. (2016). Computable approximations for continuous-time Markov decision processes on Borel spaces based on empirical measures. J. Math. Anal. Appl. 443, 1323–1361.
- [2] BHATT, A. G. AND BORKAR, V. S. (1996). Occupation measures for controlled Markov processes: Characterization and optimality. Ann. Prob. 24, 1531–1562.
- [3] BOGACHEV, V. I. (2007). Measure Theory, Vols. I and II. Springer, Berlin.
- [4] BONNANS, J. F. AND ZIDANI, H. (2004). Consistency of generalized finite difference schemes for the stochastic HJB equation. SIAM J. Numer. Anal. 41, 1008–1021.
- [5] DE BOOR, C. (2001). A Practical Guide to Splines, revised edn (Appl. Math. Sci. 27). Springer, New York.
- [6] ETHIER, S. N. AND KURTZ, T. G. (1986). Markov Processes: Characterization and Convergence John Wiley, New York.
- [7] FLEMING, W. H. AND RISHEL, R. W. (1975). Deterministic and Stochastic Optimal Control (Appl. Math. 1). Springer, New York.
- [8] FLEMING, W. H. AND SONER, H. M. (2006). Controlled Markov Processes and Viscosity Solutions, 2nd edn (Stoch. Modelling Appl. Prob. 25). Springer, New York.
- [9] GREIF, C. (2017). Numerical methods for Hamilton–Jacobi–Bellman equations. Master's thesis, University of Wisconsin-Milwaukee.
- [10] HALL, C. A. AND MEYER, W. W. (1976). Optimal error bounds for cubic spline interpolation. J. Approx. Theory 16, 105–122.
- [11] HELMES, K., RÖHL, S. AND STOCKBRIDGE, R. H. (2001). Computing moments of the exit time distribution for Markov processes by linear programming. *Operat. Res.* 49, 516–530.

- [12] HELMES, K. L., STOCKBRIDGE, R. H. AND ZHU, C. (2017). Continuous inventory models of diffusion type: Long-term average cost criterion. Ann. Appl. Prob. 27, 1831–1885.
- [13] JENSEN, M. AND SMEARS, I. (2013). On the convergence of finite element methods for Hamilton–Jacobi– Bellman equations. SIAM J. Numer. Anal. 51, 137–162.
- [14] KACZMAREK, P., KENT, S. T., RUS, G. A., STOCKBRIDGE, R. H. AND WADE, B. A. (2007). Numerical solution of a long-term average control problem for singular stochastic processes. *Math. Meth. Operat. Res.* 66, 451–473.
- [15] KUMAR, S. AND MUTHURAMAN, K. (2004). A numerical method for solving singular stochastic control problems. Operat. Res. 52, 563–582.
- [16] KURTZ, T. G. AND STOCKBRIDGE, R. H. (1998). Existence of Markov controls and characterization of optimal Markov controls. SIAM J. Control Optim. 36, 609–653.
- [17] KURTZ, T. G. AND STOCKBRIDGE, R. H. (1999). Erratum: 'Existence of Markov controls and characterization of optimal Markov controls'. SIAM J. Control Optim. 37, 1310–1311.
- [18] KUSHNER, H. J. AND DUPUIS, P. (2001). Numerical Methods for Stochastic Control Problems in Continuous Time, 2nd edn (Appl. Math. (New York) 24). Springer, New York.
- [19] LASSERRE, J.-B. AND PRIETO-RUMEAU, T. (2004). SDP vs. LP relaxations for the moment approach in some performance evaluation problems. *Stoch. Models* 20, 439–456.
- [20] LU, X., YIN, G., ZHANG, Q., ZHANG, C. AND GUO, X. (2017). Building up an illiquid stock position subject to expected fund availability: Optimal controls and numerical methods. *Appl. Math. Optim.* 76, 501–533.
- [21] MANNE, A. S. (1960). Linear programming and sequential decisions. Manag. Sci. 6, 259–267.
- [22] MENDIONDO, M. S. AND STOCKBRIDGE, R. H. (1998). Approximation of infinite-dimensional linear programming problems which arise in stochastic control. SIAM J. Control Optim. 36, 1448–1472.
- [23] PHAM, H. (2005). On some recent aspects of stochastic control and their applications. Prob. Surv. 2, 506-549.
- [24] RUS, G. A. (2009). Finite element methods for control of singular stochastic processes. PhD thesis, The University of Wisconsin-Milwaukee.
- [25] SERRANO, R. (2015). On the LP formulation in measure spaces of optimal control problems for jumpdiffusions. Systems Control Lett. 85, 33–36.
- [26] TAKSAR, M. I. (1997). Infinite-dimensional linear programming approach to singular stochastic control. SIAM J. Control Optim. 35, 604–625.
- [27] VIETEN, M. G. AND STOCKBRIDGE, R. H. (2019). Convergence of finite element methods for singular stochastic control. SIAM J. Control Optim. 56, 4336–4364.
- [28] VIETEN, M. G. AND STOCKBRIDGE, R. H. (2020). On the solution structure of infinite-dimensional linear problems stemming from singular stochastic control problems. *SIAM J. Control Optim.* 58, 3363–3388.
- [29] WANG, J. AND FORSYTH, P. (2010). Numerical solution of the Hamilton–Jacobi–Bellman formulation for continuous time mean variance asset allocation. J. Econom. Dynam. Control 34, 207–230.