

MARKOV DECISION PROCESSES UNDER MODEL UNCERTAINTY

ARIEL NEUFELD¹, JULIAN SESTER¹, MARIO ŠIKIĆ²

June 14, 2022

¹*NTU Singapore, Division of Mathematical Sciences,
21 Nanyang Link, Singapore 637371.*

²*University of Zürich, Department of Banking and Finance,
Plattenstr. 14, 8032 Zürich*

ABSTRACT. We introduce a general framework for Markov decision problems under model uncertainty in a discrete-time infinite horizon setting. By providing a dynamic programming principle we obtain a local-to-global paradigm, namely solving a local, i.e., a one time-step robust optimization problem leads to an optimizer of the global (i.e. infinite time-steps) robust stochastic optimal control problem, as well as to a corresponding worst-case measure.

Moreover, we apply this framework to portfolio optimization involving data of the *S&P 500*. We present two different types of ambiguity sets; one is fully data-driven given by a Wasserstein-ball around the empirical measure, the second one is described by a parametric set of multivariate normal distributions, where the corresponding uncertainty sets of the parameters are estimated from the data. It turns out that in scenarios where the market is volatile or bearish, the optimal portfolio strategies from the corresponding robust optimization problem outperforms the ones without model uncertainty, showcasing the importance of taking model uncertainty into account.

Keywords: Markov decision problem, Ambiguity, Dynamic programming principle, Portfolio optimization

1. INTRODUCTION

Suppose that today and at all future times an agent observes the state of the surrounding world, and based on the realization of this state she decides to execute an action that may also influence future states. All actions are rewarded according to a *reward function* not immediately but once the subsequent state is realized. The *Markov decision problem* consists of finding at initial time a policy, i.e., a sequence of state-dependent actions, that optimizes the expected cumulated discounted future rewards, referred to as the *value* of the Markov decision problem. The underlying process of states in a Markov decision problem is a stochastic process $(X_t)_{t \in \mathbb{N}_0}$ and is called *Markov decision process*. This process is usually modelled by a discrete-time time-homogeneous Markov process that follows a pre-specified probability which is influenced by the current state of the process and the agent's current action. The Markov decision problem leads to an infinite horizon stochastic optimal control problem in discrete-time which finds many applications in finance and economics, compare, e.g., [7], [22], or [39] for an overview. It can, among a multitude of other applications, be used to learn the optimal structure of portfolios and the optimal trading behaviour, see, e.g. [9], [13], [19], [24], [40], to learn optimal hedging strategies, see, e.g. [3], [4], [12], [16], [17], [21], [29], [34], or even to study socio-economic systems under the influence of climate change as in [35].

In most applications the choice of the distribution, or more specifically the probability kernel, of the Markov decision process however is a priori unclear and hence ambiguous. For this reason, in practice, the distributions of the process often need to be estimated, compare e.g. [1], [33], [36]. To account for distributional ambiguity we are therefore interested to study an optimization problem respecting uncertainty with respect to the choice of the underlying distribution of $(X_t)_{t \in \mathbb{N}_0}$ by identifying a policy that maximizes the expected future cumulated rewards under the worst case probability measure from an ambiguity set of admissible probability measures. This formulation allows the agent to act optimally even if adverse scenarios are realized, such as for example during financial crises or extremely volatile market periods in financial markets.

The recent works [5], [14], [37], and [41] also consider infinite horizon robust stochastic optimal control problems and follow a similar paradigm but use different underlying frameworks. More precisely, [14] and [41] assume a finite action and state space. The approach from [37] assumes an atomless probability space and is restricted to so called *conditional risk mappings*, whereas [5] assumes the ambiguity set of probability measures to be dominated. To the best of our knowledge, the generality of the approach presented in this paper has not been established so far in the literature.

Our general formulation enables to specify a wide range of different ambiguity sets of probability measures and associated transition kernels, given some mild technical assumptions are fulfilled. More specifically, we require the correspondence that maps a state-action pair to the set of transition probabilities to be non-empty, continuous, compact-valued, and to fulfil a linear growth condition; see Assumption 2.2. As we will show, these requirements are *naturally* satisfied. This is for example the case if the ambiguity set is modelled by a Wasserstein-ball around a transition kernel or if parameter uncertainty with respect to multivariate normal distributions is considered.

To solve the robust optimization problem we establish a dynamic programming principle that involves only a one time-step optimization problem. Via Berge's maximum theorem (see [8]) we obtain the existence of both an optimal action and a worst case transition kernel of this *local* one time-step problem. It turns out that the optimal action that solves this one time-step optimization problem determines also the *global* optimal policy of the infinite time horizon robust stochastic optimal control problem by repeatedly executing this *local* solution. Similarly, the *global* worst case measure can be determined as a product measure given by the infinite product of the worst case transition kernel of the *local* one time-step optimization problem. We refer to Theorem 2.7 for our main result. This local-to-global principle is in line with similar results for non-robust Markov decision problems, compare e.g. [7, Theorem 7.1.7], where the optimal global policy can also be determined locally. Note that the local-to-global paradigm obtained in Theorem 2.7 is noteworthy, since $(X_t)_{t \in \mathbb{N}_0}$ does not need to be a time-homogeneous Markov process under each measure from the ambiguity set, as the corresponding transition kernel might vary with time. However, due to the particular setting that the *set* of transition probabilities is constant in time and only depends on the current state and action, and not on the whole past trajectory, we are able to derive the analogue local-to-global paradigm for Markov decision processes under model uncertainty as for the ones without model uncertainty.

Eventually we show how the discussed robust stochastic optimal control framework can be applied to portfolio optimization with real data, which was already studied extensively in the non-robust case, for example in [6], [31], [42], and [44]. To that end, we show how, based on a time series of realized returns of multiple assets of the *S&P* 500, a data-driven ambiguity set of probability measures can be derived in two cases. The first case is an entirely data-driven approach where ambiguity is described by a Wasserstein-ball around the empirical measure. In the second case a multivariate normal distribution of the considered returns is assumed while the set of parameters for the multivariate normal distribution is estimated from observed data. Hence, this approach can be considered as semi data-driven approach. We then train neural networks to solve the (semi) data-driven robust optimization problem based on the local-to-global paradigm obtained in Theorem 2.7 and compare the trading performance of the two approaches with non-robust approaches. It turns out that under adverse market scenarios both robust approaches outperform comparable non-robust approaches. These results emphasize the importance of taking into account model uncertainty when making decisions that rely on financial assets.

The remainder of the paper is as follows. In Section 2 we present the setting and formulate the underlying distributionally robust stochastic optimal control problem. We also present our main results that include a dynamic programming principle. In Section 3 we discuss different possibilities to define ambiguity sets of probability measures and we show that these specifications meet the requirements of our setting. In Section 4 we apply the robust optimization approach to portfolio optimization using real financial data and compare the different ambiguity sets introduced in Section 3 also with non-robust approaches. The proof of the main results is reported in Section 5, while the proofs of the results from Section 3 and 4 can be found in Section 6 and 7, respectively. Finally, the appendix contains a description of a numerical routine that can be applied to solve the optimization problem using neural networks as well as several useful auxiliary known mathematical results.

2. SETTING, PROBLEM FORMULATION, AND MAIN RESULT

We first present the underlying setting for the considered stochastic process and then formulate an associated distributionally robust optimization problem.

2.1. Setting. We consider a closed subset $\Omega_{\text{loc}} \subseteq \mathbb{R}^d$, equipped with its Borel σ -field \mathcal{F}_{loc} , which we use to define the infinite Cartesian product

$$\Omega := \Omega_{\text{loc}}^{\mathbb{N}} = \Omega_{\text{loc}} \times \Omega_{\text{loc}} \times \dots$$

and the σ -field $\mathcal{F} := \mathcal{F}_{\text{loc}} \otimes \mathcal{F}_{\text{loc}} \otimes \dots$. We denote by $\mathcal{M}_1(\Omega)$ the set of probability measures on (Ω, \mathcal{F}) , by $d \in \mathbb{N}$ the dimension of the state space, and by $m \in \mathbb{N}$ the dimension of the control space.

On this space, we consider an infinite horizon time-discrete stochastic process. To this end, we define on Ω the stochastic process $(X_t)_{t \in \mathbb{N}_0}$ by the canonical process $X_t((\omega_0, \omega_1, \dots, \omega_t, \dots)) := \omega_t$ for $(\omega_0, \omega_1, \dots, \omega_t, \dots) \in \Omega$, $t \in \mathbb{N}_0$. We fix a compact set $A \subseteq \mathbb{R}^m$ and define the set of controls (also called actions) through

$$\begin{aligned} \mathcal{A} &:= \{ \mathbf{a} = (a_t)_{t \in \mathbb{N}_0} \mid (a_t)_{t \in \mathbb{N}_0} : \Omega \rightarrow A; a_t \text{ is } \sigma(X_t)\text{-measurable for all } t \in \mathbb{N}_0 \} \\ &= \{ (a_t(X_t))_{t \in \mathbb{N}_0} \mid a_t : \Omega_{\text{loc}} \rightarrow A \text{ Borel measurable for all } t \in \mathbb{N}_0 \}. \end{aligned}$$

For every $k \in \mathbb{N}$, $X \subseteq \mathbb{R}^k$, and $p \in \mathbb{N}_0$, we define the set of continuous functions $g : X \rightarrow \mathbb{R}$ with polynomial growth at most of degree p via

$$C_p(X, \mathbb{R}) := \left\{ g \in C(X, \mathbb{R}) \mid \sup_{x \in X} \frac{|g(x)|}{1 + \|x\|^p} < \infty \right\},$$

where $C(X, \mathbb{R})$ denotes the set of continuous functions mapping from X to \mathbb{R} and $\|\cdot\|$ denotes the Euclidean norm on \mathbb{R}^k . We define on $C_p(\Omega_{\text{loc}}, \mathbb{R})$ the norm

$$\|g\|_{C_p} := \sup_{x \in \Omega_{\text{loc}}} \frac{|g(x)|}{1 + \|x\|^p}$$

Moreover, recall the Wasserstein p -topology τ_p on $\mathcal{M}_1(\Omega_{\text{loc}})$ induced by the convergence

$$(2.1) \quad \mu_n \xrightarrow{\tau_p} \mu \text{ for } n \rightarrow \infty \Leftrightarrow \lim_{n \rightarrow \infty} \int g d\mu_n = \int g d\mu \text{ for all } g \in C_p(\Omega_{\text{loc}}, \mathbb{R}).$$

Note that for $p = 0$, the topology τ_0 coincides with the topology of weak convergence. To be able to formulate a robust optimization problem, we make use of the theory of set-valued maps, also called *correspondences*, see also [2, Chapter 17] for an extensive introduction to the topic. In the following we clarify how continuity is defined for correspondences, compare also Lemma B.3 and Lemma B.4, where characterizations of upper hemicontinuity and lower hemicontinuity are provided.

Definition 2.1. Let $\varphi : X \rightrightarrows Y$ be a correspondence between two topological spaces.

- (i) φ is called upper hemicontinuous, if $\{x \in X \mid \varphi(x) \subseteq A\}$ is open for all open sets $A \subseteq Y$.
- (ii) φ is called lower hemicontinuous, if $\{x \in X \mid \varphi(x) \cap A \neq \emptyset\}$ is open for all open sets $A \subseteq Y$.
- (iii) We say φ is continuous, if φ is upper and lower hemicontinuous.

Moreover, for a correspondence $\varphi : X \rightrightarrows Y$ its graph is defined as

$$\text{Gr } \varphi := \{(x, y) \in X \times Y \mid y \in \varphi(x)\}.$$

We impose the following standing assumptions on the process $(X_t)_{t \in \mathbb{N}_0}$ and on the set of admissible measures, which are from now on assumed to be valid for the rest of the paper.

Standing Assumption 2.2 (Assumptions on the set of measures). Fix $p \in \{0, 1\}$.

- (i) The set-valued map

$$\begin{aligned} \Omega_{\text{loc}} \times A &\rightarrow (\mathcal{M}_1(\Omega_{\text{loc}}), \tau_p) \\ (x, a) &\mapsto \mathcal{P}(x, a) \end{aligned}$$

is assumed to be nonempty, compact-valued, and continuous.

- (ii) There exists $C_P \geq 1$ such that for all $(x, a) \in \Omega_{\text{loc}} \times A$ and $\mathbb{P} \in \mathcal{P}(x, a)$ it holds

$$(2.2) \quad \int_{\Omega_{\text{loc}}} (1 + \|y\|^p) \mathbb{P}(dy) \leq C_P(1 + \|x\|^p).$$

Under these assumptions we define for every $x \in \Omega_{\text{loc}}$, $\mathbf{a} \in \mathcal{A}$ the set of admissible measures

$$\mathfrak{P}_{x,\mathbf{a}} := \left\{ \delta_x \otimes \mathbb{P}_0 \otimes \mathbb{P}_1 \otimes \cdots \left| \begin{array}{l} \text{for all } t \in \mathbb{N}_0 : \mathbb{P}_t : \Omega_{\text{loc}} \rightarrow \mathcal{M}_1(\Omega_{\text{loc}}) \text{ Borel-measurable,} \\ \text{and } \mathbb{P}_t(\omega_t) \in \mathcal{P}(\omega_t, a_t(\omega_t)) \text{ for all } \omega_t \in \Omega_{\text{loc}} \end{array} \right. \right\},$$

where the notation $\mathbb{P} = \delta_x \otimes \mathbb{P}_0 \otimes \mathbb{P}_1 \otimes \cdots \in \mathfrak{P}_{x,\mathbf{a}}$ abbreviates

$$\mathbb{P}(B) := \int_{\Omega_{\text{loc}}} \cdots \int_{\Omega_{\text{loc}}} \cdots \mathbf{1}_B((\omega_t)_{t \in \mathbb{N}_0}) \cdots \mathbb{P}_{t-1}(\omega_{t-1}; d\omega_t) \cdots \mathbb{P}_0(\omega_0; d\omega_1) \delta_x(d\omega_0), \quad B \in \mathcal{F}.$$

Remark 2.3. *To ensure that the set $\mathfrak{P}_{x,\mathbf{a}}$ is nonempty, one needs to show that \mathcal{P} admits a measurable selector. By Assumption 2.2 the correspondence $\Omega_{\text{loc}} \times A \ni (x, a) \rightarrow \mathcal{P}(x, a)$ is closed-valued and measurable. Hence, by Kuratovski's Theorem (compare, e.g., [2, Theorem 18.13]), there exists a measurable selector $\Omega_{\text{loc}} \times A \ni (x, a) \mapsto \mathbb{P}(x, a) \in \mathcal{M}_1(\Omega_{\text{loc}})$ such that $\mathbb{P}(x, a) \in \mathcal{P}(x, a)$ for all $(x, a) \in \Omega_{\text{loc}} \times A$. Since actions are by definition measurable, we also obtain that for all $(a_t)_{t \in \mathbb{N}_0} \in \mathcal{A}$ and for all $t \in \mathbb{N}_0$ the map $\Omega_{\text{loc}} \ni \omega_t \mapsto \mathbb{P}(\omega_t, a_t(\omega_t)) =: \mathbb{P}_t(\omega_t; d\omega_{t+1})$ is measurable, as required. Then, the non-emptiness of $\mathfrak{P}_{x,\mathbf{a}}$ follows by the Ionescu–Tulcea theorem.*

2.2. Problem Formulation. Let $r : \Omega_{\text{loc}} \times A \times \Omega_{\text{loc}} \rightarrow \mathbb{R}$ be some reward function. We assume from now on that it fulfils the following assumptions.

Standing Assumption 2.4 (Assumptions on the reward function and the discount factor). *Let $p \in \{0, 1\}$ be the number fixed in Assumption 2.2.*

(i) *The map*

$$\Omega_{\text{loc}} \times A \times \Omega_{\text{loc}} \ni (x_0, a, x_1) \mapsto r(x_0, a, x_1)$$

is continuous

(ii) *There exists some $L > 0$ such that for all $x_0, x'_0, x_1 \in \Omega_{\text{loc}}$ and $a, a' \in A$ we have*

$$(2.3) \quad |r(x_0, a, x_1) - r(x'_0, a', x_1)| \leq L \cdot \max\{1, \|x_1\|^p\} (\|x_0 - x'_0\| + \|a - a'\|).$$

(iii) *There exists some $C_r \geq 1$ such that for all $x_0, x_1 \in \Omega_{\text{loc}}$ we have*

$$(2.4) \quad |r(x_0, a, x_1)| \leq C_r(1 + \|x_0\|^p + \|x_1\|^p) \text{ for all } a \in A.$$

(iv) *We fix an associated discount factor $\alpha < 1$ which satisfies*

$$0 < \alpha < \frac{1}{C_r(C_p + 1)C_p}.$$

Remark 2.5. *Note that if r is Lipschitz-continuous, i.e., if there exists some $L > 0$ such that for all $x_0, x'_0, x_1, x'_1 \in \Omega_{\text{loc}}$ and $a, a' \in A$ we have*

$$|r(x_0, a, x_1) - r(x'_0, a', x'_1)| \leq L (\|x_0 - x'_0\| + \|a - a'\| + \|x_1 - x'_1\|),$$

then (2.3) follows directly. Therefore the requirement of Assumption 2.4 (i) and (ii) is weaker than assuming Lipschitz continuity of the reward function. In particular, if $m = d$ holds for the dimensions, then the function of the form

$$(2.5) \quad \Omega_{\text{loc}} \times A \times \Omega_{\text{loc}} \ni (x_0, a, x_1) \mapsto r(x_0, a, x_1) := a \cdot x_1$$

fulfils the requirement imposed in (2.3) but are not Lipschitz continuous, unless Ω_{loc} is bounded. Compare also Section 4, where we apply portfolio optimization while taking into account a reward function of the form (2.5).

Our main problem consists, for every initial value $x \in \Omega_{\text{loc}}$, in maximizing the expected value of $\sum_{t=0}^{\infty} \alpha^t r(X_t, a_t, X_{t+1})$ under the worst case measure from $\mathfrak{P}_{x,\mathbf{a}}$ over all possible actions $\mathbf{a} \in \mathcal{A}$. More precisely, we introduce the value function

$$(2.6) \quad \Omega_{\text{loc}} \ni x \mapsto V(x) := \sup_{\mathbf{a} \in \mathcal{A}} \inf_{\mathbb{P} \in \mathfrak{P}_{x,\mathbf{a}}} \left(\mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{\infty} \alpha^t r(X_t, a_t, X_{t+1}) \right] \right).$$

Definition 2.6. *We call $(X_t)_{t \in \mathbb{N}_0}$ a Markov decision process under model uncertainty on state space $\Omega_{\text{loc}} \subseteq \mathbb{R}^d$ with corresponding set of transition probabilities \mathcal{P} , and we call the problem defined in (2.6) a Markov decision problem under model uncertainty.*

2.3. Main Result: The Dynamic Programming Principle. In this section we provide the main results of the paper which comprise a *dynamic programming principle* which in particular allows to solve the optimization problem (2.6) by solving a related one-step fix point equation.

To this end, we define the space of *one-step actions*

$$\mathcal{A}_{\text{loc}} := \{a_{\text{loc}} : \Omega_{\text{loc}} \rightarrow A \text{ measurable}\},$$

and, we define for every $a_{\text{loc}} \in \mathcal{A}_{\text{loc}}$ the set of kernels

$$\mathbf{P}_{a_{\text{loc}}} := \{\mathbb{P}_0 : \Omega_{\text{loc}} \rightarrow \mathcal{M}_1(\Omega_{\text{loc}}) \text{ measurable} \mid \mathbb{P}_0(x) \in \mathcal{P}(x, a_{\text{loc}}(x)) \text{ for all } x \in \Omega_{\text{loc}}\}.$$

Moreover, we define on $C_p(\Omega_{\text{loc}}, \mathbb{R})$ the operator \mathcal{T} which for every $v \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ is defined by

$$\Omega_{\text{loc}} \ni x \mapsto \mathcal{T}v(x) := \sup_{a \in A} \inf_{\mathbb{P} \in \mathcal{P}(x, a)} \mathbb{E}_{\mathbb{P}} [r(x, a, X_1) + \alpha v(X_1)].$$

Our main findings are collected in the subsequent theorem.

Theorem 2.7. *Assume that Assumption 2.2 and Assumption 2.4 hold true. Then the following holds.*

- (i) *For every $v \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ there exists $\mathbb{P}_0^* : \Omega_{\text{loc}} \times A \rightarrow \mathcal{M}_1(\Omega_{\text{loc}})$ such that for all $(x, a) \in \Omega_{\text{loc}} \times A$ we have $\mathbb{P}_0^*(x, a) \in \mathcal{P}(x, a)$ and*

$$(2.7) \quad \begin{aligned} \mathbb{E}_{\mathbb{P}_0^*(x, a)} [r(x, a, X_1) + \alpha v(X_1)] &:= \int_{\Omega_{\text{loc}}} r(x, a, \omega_1) + \alpha v(\omega_1) \mathbb{P}_0^*(x, a; d\omega_1) \\ &= \inf_{\mathbb{P}_0 \in \mathcal{P}(x, a)} \mathbb{E}_{\mathbb{P}_0} [r(x, a, X_1) + \alpha v(X_1)]. \end{aligned}$$

Moreover, there exists $a_{\text{loc}}^* \in \mathcal{A}_{\text{loc}}$ such that for every $x \in \Omega_{\text{loc}}$ we have

$$(2.8) \quad \begin{aligned} &\inf_{\mathbb{P}_0 \in \mathcal{P}(x, a_{\text{loc}}^*(x))} \mathbb{E}_{\mathbb{P}_0} [r(x, a_{\text{loc}}^*(x), X_1) + \alpha v(X_1)] \\ &= \sup_{a_{\text{loc}} \in \mathcal{A}_{\text{loc}}} \inf_{\mathbb{P}_0 \in \mathcal{P}(x, a_{\text{loc}}(x))} \mathbb{E}_{\mathbb{P}_0} [r(x, a_{\text{loc}}(x), X_1) + \alpha v(X_1)]. \end{aligned}$$

Furthermore, let $\mathbb{P}_{\text{loc}}^* : \Omega_{\text{loc}} \rightarrow \mathcal{M}_1(\Omega_{\text{loc}})$ be defined by

$$\mathbb{P}_{\text{loc}}^*(x) := \mathbb{P}_0^*(x, a_{\text{loc}}^*(x)), \quad x \in \Omega_{\text{loc}}.$$

Then $\mathbb{P}_{\text{loc}}^* \in \mathbf{P}_{a_{\text{loc}}^*}$ and for every $x \in \Omega_{\text{loc}}$ it holds that

$$(2.9) \quad \begin{aligned} \mathcal{T}v(x) &= \sup_{a_{\text{loc}} \in \mathcal{A}_{\text{loc}}} \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}}} \mathbb{E}_{\mathbb{P}_0(x)} [r(x, a_{\text{loc}}(x), X_1) + \alpha v(X_1)] \\ &= \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}^*}} \mathbb{E}_{\mathbb{P}_0(x)} [r(x, a_{\text{loc}}^*(x), X_1) + \alpha v(X_1)] \\ &= \mathbb{E}_{\mathbb{P}_{\text{loc}}^*(x)} [r(x, a_{\text{loc}}^*(x), X_1) + \alpha v(X_1)]. \end{aligned}$$

- (ii) *We have that $\mathcal{T}(C_p(\Omega_{\text{loc}}, \mathbb{R})) \subseteq C_p(\Omega_{\text{loc}}, \mathbb{R})$, i.e., $\mathcal{T}v \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ for all $v \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ and for all $v, w \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ the following inequality holds true*

$$(2.10) \quad \|\mathcal{T}v - \mathcal{T}w\|_{C_p} \leq \alpha C_P \|v - w\|_{C_p}.$$

In particular, there exists a unique $v \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ such that $\mathcal{T}v = v$. Moreover, for every $v_0 \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ we have $v = \lim_{n \rightarrow \infty} \mathcal{T}^n v_0$.

- (iii) *Let $v \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ satisfy $\mathcal{T}v = v$ and let $a_{\text{loc}}^* \in \mathcal{A}_{\text{loc}}$, $\mathbb{P}_{\text{loc}}^* \in \mathbf{P}_{a_{\text{loc}}^*}$ be defined as in (i). Define $\mathbf{a}^* := (a_{\text{loc}}^*(X_0), a_{\text{loc}}^*(X_1), \dots) \in \mathcal{A}$ and for all $x \in \Omega_{\text{loc}}$, $\mathbb{P}_x^* := \delta_x \otimes \mathbb{P}_{\text{loc}}^* \otimes \mathbb{P}_{\text{loc}}^* \otimes \dots \in \mathfrak{P}_{x, \mathbf{a}^*}$. Then, for all $x \in \Omega_{\text{loc}}$ we have that*

$$(2.11) \quad \begin{aligned} \mathbb{E}_{\mathbb{P}_x^*} \left[\sum_{t=0}^{\infty} \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) \right] &= \inf_{\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}^*}} \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{\infty} \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) \right] \\ &= V(x) \\ &= v(x). \end{aligned}$$

Remark 2.8. Note that the local-to-global paradigm obtained in Theorem 2.7 is noteworthy, since $(X_t)_{t \in \mathbb{N}_0}$ does not need to be a (time homogeneous) Markov process under each $\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}^*}$, as the corresponding transition kernel might vary with time. However, due to the particular setting that the set of transition probabilities $(x, a) \mapsto \mathcal{P}(x, a)$ is constant in time and only depends on the current state and action, and not on the whole past trajectory, we are able to derive the analogue local-to-global paradigm for Markov decision processes under model uncertainty as for the ones without model uncertainty. Moreover, if $(x, a) \mapsto \mathcal{P}(x, a)$ is single-valued, then $(X_t)_{t \in \mathbb{N}_0}$ is a Markov decision process in the classical sense, compare, e.g., [7]. This justifies to call $(X_t)_{t \in \mathbb{N}_0}$ a Markov decision process under model uncertainty on state space $\Omega_{\text{loc}} \subseteq \mathbb{R}^d$ with respect to \mathcal{P} . Moreover, note that a posteriori, we see that $(X_t)_{t \in \mathbb{N}_0}$ is a time-homogeneous Markov process under the worst-case measure, and the optimal strategy only depends on the current state of the process, and not on time, as observed for classical Markov decision problems.

3. CAPTURING DISTRIBUTIONAL UNCERTAINTY

In this section we present different approaches that enable to capture uncertainty with respect to the choice of the underlying probability measure. We show that all of the presented approaches fulfil the requirements of the setting presented in Section 2.

3.1. Uncertainty expressed through the Wasserstein distance. The first example involves the case when distributional uncertainty is captured through the q -Wasserstein-distance $W_q(\cdot, \cdot)$ for some $q \in \mathbb{N}$. For any $\mathbb{P}_1, \mathbb{P}_2 \in \mathcal{M}_1(\Omega_{\text{loc}})$ let $W_q(\mathbb{P}_1, \mathbb{P}_2)$ be defined as

$$W_q(\mathbb{P}_1, \mathbb{P}_2) := \left(\inf_{\pi \in \Pi(\mathbb{P}_1, \mathbb{P}_2)} \int_{\Omega_{\text{loc}} \times \Omega_{\text{loc}}} \|x - y\|^q d\pi(x, y) \right)^{1/q},$$

where $\|\cdot\|$ denotes the Euclidean norm on \mathbb{R}^d , and where $\Pi(\mathbb{P}_1, \mathbb{P}_2)$ denotes the set of joint distributions of \mathbb{P}_1 and \mathbb{P}_2 , compare also for example [38, Definition 6.1.].

We fix some $q \in \mathbb{N}$ and specify $p := 0$ in Assumption 2.2 and 2.4. Further, we assume that there exists a continuous map

$$(3.1) \quad \begin{aligned} \Omega_{\text{loc}} \times A &\rightarrow (\mathcal{M}_1(\Omega_{\text{loc}}), \tau_q) \\ (x, a) &\mapsto \widehat{\mathbb{P}}(x, a) \end{aligned}$$

such that $\widehat{\mathbb{P}}(x, a)$ has finite q -th moments for all $(x, a) \in \Omega_{\text{loc}} \times A$. Then, we define for any $\varepsilon > 0$ the set-valued map

$$(3.2) \quad \Omega_{\text{loc}} \times A \ni (x, a) \mapsto \mathcal{P}(x, a) := \mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a)) := \left\{ \mathbb{P} \in \mathcal{M}_1(\Omega_{\text{loc}}) \mid W_q(\mathbb{P}, \widehat{\mathbb{P}}(x, a)) \leq \varepsilon \right\},$$

where $\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a))$ denotes the q -Wasserstein-ball (or Wasserstein-ball of order q) with ε -radius and center $\widehat{\mathbb{P}}(x, a)$.

Proposition 3.1. Let $\Omega_{\text{loc}} \times A \ni (x, a) \mapsto \widehat{\mathbb{P}}(x, a) \in (\mathcal{M}_1(\Omega_{\text{loc}}), \tau_q)$ be continuous with finite q -th moments. Then, the set-valued map $\Omega_{\text{loc}} \times A \ni (x, a) \mapsto \mathcal{P}(x, a)$ defined as in (3.2) fulfils the requirements of Assumption 2.2 with $p = 0$.

3.2. Knightian uncertainty in parametric models. Next, we consider a parametric approach, taking into account the so called *Knightian uncertainty* (see [27]). To this end, we consider a set-valued map of the form

$$(3.3) \quad \Omega_{\text{loc}} \times A \ni (x, a) \mapsto \Theta(x, a) \subseteq \mathbb{R}^{\mathfrak{D}}, \quad \text{for some } \mathfrak{D} \in \mathbb{N}.$$

The set $\Theta(x, a)$ refers to the set of parameters that are admissible in dependence of $(x, a) \in \Omega_{\text{loc}} \times A$. The underlying parametric probability distribution is described by

$$(3.4) \quad \begin{aligned} \{(x, a, \theta) \mid (x, a) \in \Omega_{\text{loc}} \times A, \theta \in \Theta(x, a)\} &\rightarrow (\mathcal{M}_1(\Omega_{\text{loc}}), \tau_p) \\ (x, a, \theta) &\mapsto \widehat{\mathbb{P}}(x, a, \theta), \end{aligned}$$

which enables us to define the ambiguity set of probability measures by

$$(3.5) \quad \Omega_{\text{loc}} \times A \ni (x, a) \mapsto \mathcal{P}(x, a) := \left\{ \widehat{\mathbb{P}}(x, a, \theta) \mid \theta \in \Theta(x, a) \right\} \subseteq (\mathcal{M}_1(\Omega_{\text{loc}}), \tau_p).$$

Proposition 3.2. *Let $(x, a) \mapsto \Theta(x, a)$, as defined in (3.3), be nonempty, compact-valued, and continuous, let $(x, a, \theta) \mapsto \widehat{\mathbb{P}}(x, a, \theta)$, as defined in (3.4), be continuous. Then $(x, a) \mapsto \mathcal{P}(x, a)$, as defined in (3.5), is nonempty, compact-valued, and continuous.*

3.3. Uncertainty in autocorrelated time series. Next, we consider the case where the state process $(X_t)_{t \in \mathbb{N}_0}$ is given by an autocorrelated time series. More precisely, we assume that at time $t \in \mathbb{N}_0$ the past $m \in \mathbb{N}$ observations (Y_{t-m+1}, \dots, Y_t) of a time series $(Y_t)_{t \in \{-m, \dots, -1, 0, 1, \dots\}}$ may have an influence on the next value of the state process. In this case we have for all $t \in \mathbb{N}$ a representation of the form

$$X_t := (Y_{t-m+1}, \dots, Y_t) \in \Omega_{\text{loc}} := T^m \subseteq \mathbb{R}^{D \cdot m}, \text{ with } T \subseteq \mathbb{R}^D \text{ closed, for some } D \in \mathbb{N}.$$

To define the ambiguity set of measures, we first consider a set-valued map of the form

$$(3.6) \quad \Omega_{\text{loc}} \times A \ni (x, a) \mapsto \widetilde{\mathcal{P}}(x, a) \subseteq (\mathcal{M}_1(T), \tau_p).$$

We consider the projection $\Omega_{\text{loc}} \ni (x_1, \dots, x_m) \mapsto \pi((x_1, \dots, x_m)) := (x_2, \dots, x_m) \in T^{m-1}$ that projects onto the last $m-1$ components, and define a set-valued map \mathcal{P} , in dependence of $\widetilde{\mathcal{P}}$, by

$$(3.7) \quad \Omega_{\text{loc}} \times A \ni (x, a) \mapsto \mathcal{P}(x, a) := \left\{ \delta_{\pi(x)} \otimes \mathbb{P} \mid \mathbb{P} \in \widetilde{\mathcal{P}}(x, a) \right\} \subseteq (\mathcal{M}_1(\Omega_{\text{loc}}), \tau_p).$$

This means, by considering \mathcal{P} , we take into account uncertainty with respect to the evolution of the next value of the time series. However, we do not want to consider uncertainty with respect to the $m-1$ preceding values of the time series, as they constitute of the already observed realizations.

Proposition 3.3. *Let $(x, a) \mapsto \widetilde{\mathcal{P}}(x, a)$, as defined in (3.6), be nonempty, compact-valued, and continuous. Then $(x, a) \mapsto \mathcal{P}(x, a)$, as defined in (3.7) is nonempty, compact-valued, and continuous.*

4. APPLICATION TO PORTFOLIO OPTIMIZATION

In this section we discuss a finance-related application of the presented robust stochastic optimal control problem of Section 2. In particular, we compare different specifications to measure uncertainty with respect to the choice of the underlying probability measure.

Setting. We present a setting that can be applied to the robust optimization of financial portfolios. Compare among many others also [11], [16, Chapter 10], and [18], where alternative approaches to portfolio optimization relying on the optimal control of Markov decision processes are discussed.

Let $D \in \mathbb{N}$ denote the number of assets that are taken into account for portfolio optimization. Then, the underlying asset returns in the time period between $t-1$ and t are given by

$$\mathcal{R}_t := (\mathcal{R}_t^i)_{i=1, \dots, D} := \left(\frac{S_t^i - S_{t-1}^i}{S_{t-1}^i} \right)_{i=1, \dots, D} \in T \subseteq \mathbb{R}^D, \quad t \in \{-m+1, \dots, 0, 1, \dots\},$$

where $S_t^i \in (0, \infty)$ denotes the time t -value of asset $i \in \{1, \dots, D\}$, $m \in \mathbb{N}$, and $T \subseteq \mathbb{R}^D$ closed.

To take into account the autocorrelation of the time series, we want to base our portfolio allocation decisions not only on the current portfolio allocation and the present state of the financial market, but also on the past $m \in \mathbb{N}$ observed returns. Thus, we consider at every time $t \in \mathbb{N}_0$ realized returns $(\mathcal{R}_{t-m+1}, \dots, \mathcal{R}_t) \in \mathbb{R}^{D \cdot m}$. Then, the underlying stochastic process $(X_t)_{t \in \mathbb{N}_0}$ is modelled as

$$(4.1) \quad X_t := (\mathcal{R}_{t-m+1}, \dots, \mathcal{R}_t) \in \Omega_{\text{loc}}, \quad t \in \mathbb{N}_0,$$

with

$$\Omega_{\text{loc}} := T^m \subseteq \mathbb{R}^{D \cdot m}.$$

Next, we introduce the compact set

$$A := \{a = (a^i)_{i=1, \dots, D} \in [-C, C]^D\},$$

for the possible values of the controls, which corresponds to the monetary investment in the D stocks, where $C > 0$ relates to a budget constraint when investing. Then, we define the reward function by

$$(4.2) \quad \Omega_{\text{loc}} \times A \times \Omega_{\text{loc}} \ni (X_t, a_t, X_{t+1}) \mapsto r(X_t, a_t, X_{t+1}) := \sum_{i=1}^D a_t^i \cdot \mathcal{R}_{t+1}^i.$$

The reward function in (4.2) expresses the cumulated gain from trading in the period between t and $t+1$.

4.1. Data-driven ambiguity set and Wasserstein-uncertainty. We rely on the setting elaborated above.

As exposed in Section 3.1, we may capture distributional uncertainty by considering a Wasserstein-ball around some kernel

$$\Omega_{\text{loc}} \times A \ni (x, a) \mapsto \widehat{\mathbb{P}}(x, a) \in \mathcal{M}_1(T),$$

for $T \subseteq \mathbb{R}^D$ closed. We consider a time series of past realized returns

$$(4.3) \quad (\mathcal{R}_1, \dots, \mathcal{R}_N) \in T^N, \quad \text{for some } N \in \mathbb{N}.$$

Compare also Figure 1, where we illustrate the relation between this time series and the time series of future returns.

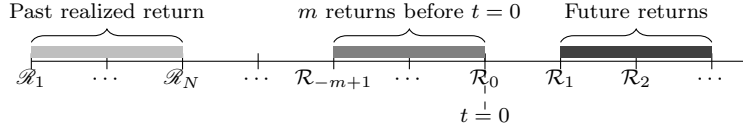


FIGURE 1. Illustration of the observed and already realized return $(\mathcal{R}_t)_{t=1,\dots,N}$ and the future random returns $(\mathcal{R}_t)_{t=-m+1,\dots,0,1,\dots}$

Relying on the time series from (4.3), we aim at constructing an ambiguity set \mathcal{P} . To this end, we define $\widehat{\mathbb{P}}$ through a sum of Dirac-measures given by¹

$$(4.4) \quad \Omega_{\text{loc}} \ni X_t = (\mathcal{R}_{t-m+1}, \dots, \mathcal{R}_t) \mapsto \widehat{\mathbb{P}}(X_t)(dx) := \sum_{s=m}^{N-1} \pi_s(X_t) \cdot \delta_{\mathcal{R}_{s+1}}(dx) \in \mathcal{M}_1(T),$$

where $\pi_s(X_t) \in [0, 1]$, $s = m, \dots, N-1$ with $\sum_{s=m}^{N-1} \pi_s(X_t) = 1$. We want to weight the distance between the past m returns before \mathcal{R}_{t+1} and the m returns before \mathcal{R}_{s+1} , while assigning higher probabilities to more similar sequences of m returns. This means, the measure $\widehat{\mathbb{P}}$ relies its prediction for the next return on the best fitting sequence of m consecutive returns that precede the prediction. To this end, we set for some (small) constant $\tilde{\varepsilon} > 0$ ²

$$\Omega_{\text{loc}} \ni X_t = (\mathcal{R}_{t-m+1}, \dots, \mathcal{R}_t) \mapsto \pi_s(X_t) := \left(\frac{(\text{dist}_s(X_t) + \tilde{\varepsilon})^{-1}}{\sum_{\ell=m}^{N-1} (\text{dist}_\ell(X_t) + \tilde{\varepsilon})^{-1}} \right),$$

with

$$\text{dist}_s(X_t) := \|(\mathcal{R}_{s-m+1}, \dots, \mathcal{R}_s) - X_t\| = \|(\mathcal{R}_{s-m+1}, \dots, \mathcal{R}_s) - (\mathcal{R}_{t-m+1}, \dots, \mathcal{R}_t)\|,$$

for all $s = m, \dots, N-1$. Then, we define for any fixed $\varepsilon > 0$ and $q \in \mathbb{N}$ the ambiguity set of probability measures on $\mathcal{M}_1(\Omega_{\text{loc}})$ via the set-valued map³

$$(4.5) \quad \Omega_{\text{loc}} \ni x \mapsto \mathcal{P}(x) := \left\{ \delta_{\pi(x)} \otimes \mathbb{P} \mid \mathbb{P} \in \mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x)) \right\} \subseteq (\mathcal{M}_1(\Omega_{\text{loc}}), \tau_p)$$

that takes into account a q -Wasserstein-ball around $\widehat{\mathbb{P}}$ for the next future return.

Proposition 4.1. *Let $T \subseteq \mathbb{R}^D$ be compact, and let $p = 0$. Then, the set-valued map \mathcal{P} , defined in (4.5), satisfies Assumption 2.2. Moreover, the reward function r , defined in (4.2), satisfies Assumption 2.4.*

¹Note that $x \mapsto \widehat{\mathbb{P}}(x)$ does not depend on $a \in A$.

²The constant $\tilde{\varepsilon} > 0$ is merely a technical requirement which is considered to avoid division by zero in the case $\text{dist}_s(X_t) = 0$ for some indices $s \in \{1, \dots, N\}, t \in \mathbb{N}_0$, i.e., in the case that a sequence of m random returns equals a sequence of past realized returns. Hence, in practice, $\tilde{\varepsilon}$ can be set to be a negligible small positive real number.

³Note that \mathcal{P} does not depend on $a \in A$, and recall that $\Omega_{\text{loc}} \ni (x_1, \dots, x_m) \mapsto \pi((x_1, \dots, x_m)) := (x_2, \dots, x_m) \in T^{m-1}$ denotes the projection onto the last $m-1$ components.

4.2. Parametric Uncertainty. Next, we introduce a parametric approach in which we assume that the asset returns follow a multivariate normal distribution with unknown parameters.⁴

To this end, we build on the setting exposed in Section 4, where $m > 1$, and where we choose $T = \mathbb{R}^D$, and $p = 1$. Moreover, we consider the following unbiased estimators of mean and covariance

$$(4.6) \quad \begin{aligned} \mathbf{m} : (\mathbb{R}^D)^m &\rightarrow \mathbb{R}^D \\ x = (x_1, \dots, x_m) &\mapsto \frac{1}{m} \sum_{i=1}^m x_i, \end{aligned}$$

and

$$(4.7) \quad \mathbf{c} : (\mathbb{R}^D)^m \rightarrow \mathbb{R}^{D \times D} \\ x = (x_1, \dots, x_m) \mapsto \frac{1}{m-1} \sum_{i=1}^m (x_i - \mathbf{m}(x)) \cdot (x_i - \mathbf{m}(x))^T.$$

Let $\varepsilon > 0$. To define the set of admissible parameters we consider the following set-valued maps

$$\begin{aligned} \Omega_{\text{loc}} \ni x &\rightarrow \hat{\mu}(x) := \{ \mu \in \mathbb{R}^D \mid \|\mu - \mathbf{m}(x)\| \leq \varepsilon \}, \\ \Omega_{\text{loc}} \ni x &\rightarrow \hat{\Sigma}(x) := \{ \Sigma \in \mathbb{R}^{D \times D} \mid \Sigma = \mathbf{c}(y) \text{ for some } y \in \Omega_{\text{loc}} \text{ with } \|y - x\| \leq \varepsilon \}, \\ \Omega_{\text{loc}} \ni x &\rightarrow \Theta(x) := \left\{ (\mu, \Sigma) \in \mathbb{R}^D \times \mathbb{R}^{D \times D} \mid \mu \in \hat{\mu}(x), \Sigma \in \hat{\Sigma}(x) \right\}. \end{aligned}$$

We define an ambiguity set related to D -dimensional multivariate normal distributions by

$$\Omega_{\text{loc}} \ni x \rightarrow \tilde{\mathcal{P}}(x) := \{ \mathcal{N}_D(\mu, \Sigma) \mid (\mu, \Sigma) \in \Theta(x) \} \subseteq (\mathcal{M}_1(\mathbb{R}^D), \tau_1).$$

As in Section 3.3 we denote by $\Omega_{\text{loc}} \ni (x_1, \dots, x_m) \mapsto \pi((x_1, \dots, x_m)) := (x_2, \dots, x_m) \in \mathbb{R}^{D \cdot (m-1)}$ the projection onto the last $m - 1$ components. These definitions allow us to define the ambiguity set on $\mathcal{M}_1(\Omega_{\text{loc}})$ by⁵

$$(4.8) \quad \Omega_{\text{loc}} \times A \ni (x, a) \rightarrow \mathcal{P}(x, a) := \left\{ \delta_{\pi(x)} \otimes \mathbb{P} \mid \mathbb{P} \in \tilde{\mathcal{P}}(x) \right\} \subseteq (\mathcal{M}_1(\Omega_{\text{loc}}), \tau_1).$$

This means, we consider as an ambiguity set for the next return a set of multivariate normal distributions with unknown mean and covariance, where the set of admissible means and covariances is specified by the estimators \mathbf{m} and \mathbf{c} as well as by the degree of ambiguity specified through ε .

Proposition 4.2. *Let $T = \mathbb{R}^D$ and $p = 1$. Then, the ambiguity set \mathcal{P} , as defined in (4.8), fulfils the requirements from Assumption 2.2. Moreover, the reward function r , defined in (4.2), satisfies Assumption 2.4.*

4.3. Numerical Experiments. In the sequel we solve the portfolio optimization problem that was discussed in Section 4 by applying the numerical method based on Theorem 2.7 that is elaborated in Appendix A to real financial data. In particular, we compare the different approaches to capture distributional uncertainty outlined in Section 4.1 and 4.2, respectively, and evaluate how these approaches perform under different market scenarios.

4.3.1. Implementation. To apply the numerical method from Appendix A, we use the following hyperparameters: Number of measures $N_{\mathcal{P}} = 10$; Batch size $B = 2^8$; Monte-Carlo sample size $N_{\text{MC}} = 2^3$; Discount factor $\alpha = 0.45$; Number of epochs $E = 50$; number of iterations for a : $\text{Iter}_a = 10$; number of iterations for v : $\text{Iter}_v = 10$. The neural networks that approximate a and v constitute of 2 layers with 128 neurons each possessing *ReLU* activation functions in each layer, except for the output layer of a which possesses a *tanh* activation function in order to constraint the output. The learning rate used to optimize the networks a and v when applying the *Adam* optimizer ([26]) is 0.001. Further details of the implementation can be found under <https://github.com/juliansester/Robust-Portfolio-Optimization>.

⁴We say that $X \in \mathbb{R}^D$ has a D -dimensional multivariate normal distribution with mean $\mu \in \mathbb{R}^D$ and covariance matrix $\Sigma \in \mathbb{R}^{D \times D}$ which is symmetric and positive semidefinite if the characteristic function of X is of the form $\mathbb{R}^D \ni u \mapsto \varphi_X(u) := \exp(iu^T \mu - \frac{1}{2}u^T \Sigma u)$, compare e.g. [20, p. 124]. We write $X \sim \mathcal{N}_D(\mu, \Sigma)$.

⁵Note that \mathcal{P} does not depend on $a \in A$, and recall that $\Omega_{\text{loc}} \ni (x_1, \dots, x_m) \mapsto \pi((x_1, \dots, x_m)) := (x_2, \dots, x_m) \in \mathbb{R}^{D \cdot (m-1)}$ denotes the projection onto the last $m - 1$ components.

4.3.2. *Data.* To train and test the performance of the portfolio optimization approach, we consider the price evolution of $d = 5$ constituents⁶ of the *S&P 500* between 2010 and 2021. We consider a *lookback period* of $m = 10$ days, i.e., the prediction of the optimal trading execution relies on the previously realized $m = 10$ returns. We split the data into a training period ranging from January 2010 until September 2018, and three different testing periods thereafter. The normalized evolution of the asset values is depicted in Figure 2.

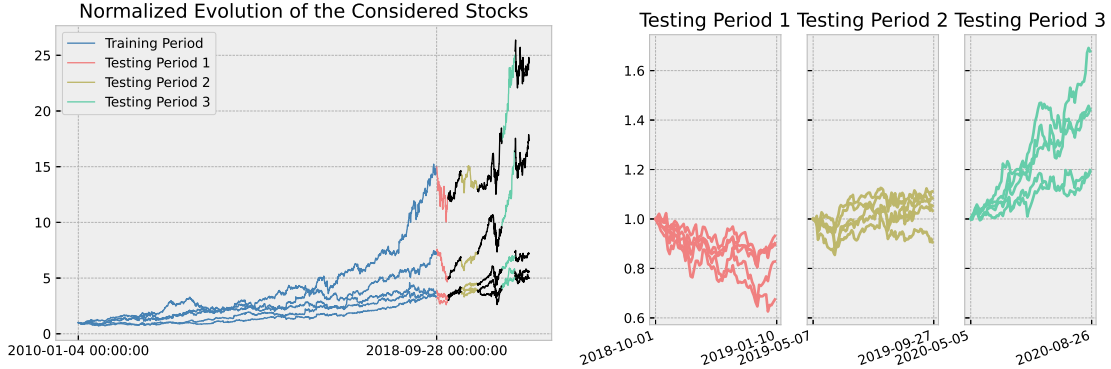


FIGURE 2. The left panel of the figure shows the normalized evolution (with initial value = 1) of the stocks of $d = 5$ constituents of the *S&P 500*. Further, we divide the data into a training phase (blue) and three testing periods thereafter that are highlighted with different colors.

The right panel of the figure shows the normalized evolution (with initial value = 1 for each of the assets) of the considered stocks in the three testing periods.

The testing periods are illustrated in detail in the right panel of Figure 2, and they comprise three different market scenarios. While the first testing period covers an overall declining market phase, the second testing period is a volatile period without a clear trend. Eventually, the third period is a bullish market with a strong upward trend.

4.3.3. *Results.* In Table 1, 2, and 3, respectively, we depict the results of the numerical method from Appendix A applied to the presented data in the three testing periods with different radii ε for both of the considered approaches to define ambiguity sets, see Section 4.1 and 4.2. Note also that we consider different ranges of values for ε for the two considered approaches. In Figure 3, 4, and 5, we depict the cumulated trading profits of the trained strategies in the respective testing periods.

Testing Period 1

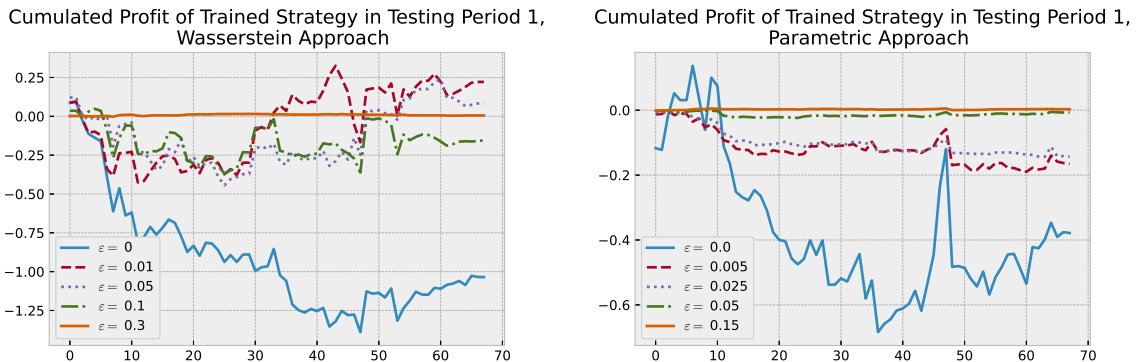


FIGURE 3. The figure shows the cumulated training profit of the trained strategies in testing period 1. The left panel illustrates the profit when applying a Wasserstein approach, whereas the right panel illustrates the profit under a parametric approach.

⁶The constituents are Apple, Microsoft, Google, Ebay, and Amazon.

	Overall Profit	Average Profit	% of profitable trades	Sharpe Ratio	Sortino Ratio
Wasserstein Approach					
$\varepsilon = 0$	-1.035019	-0.015221	47.06	-0.172211	-0.207249
$\varepsilon = 0.01$	0.221663	0.003260	58.82	0.035932	0.052785
$\varepsilon = 0.05$	0.080857	0.001189	55.88	0.015079	0.022100
$\varepsilon = 0.1$	-0.154961	-0.002279	44.12	-0.028555	-0.041425
$\varepsilon = 0.3$	0.005442	0.000080	48.53	0.035702	0.055826
Parametric Approach					
$\varepsilon = 0$	-0.378669	-0.005569	50.00	-0.067022	-0.083380
$\varepsilon = 0.005$	-0.165537	-0.002434	42.65	-0.127499	-0.146742
$\varepsilon = 0.025$	-0.143357	-0.002108	38.24	-0.198093	-0.225493
$\varepsilon = 0.05$	-0.006977	-0.000103	50.00	-0.035993	-0.043762
$\varepsilon = 0.15$	0.003136	0.000046	58.82	0.053274	0.066614

TABLE 1. The table shows the results of the Wasserstein-ball approach (Section 4.1) and the parametric approach (Section 4.2) in the first testing period.

Testing Period 2

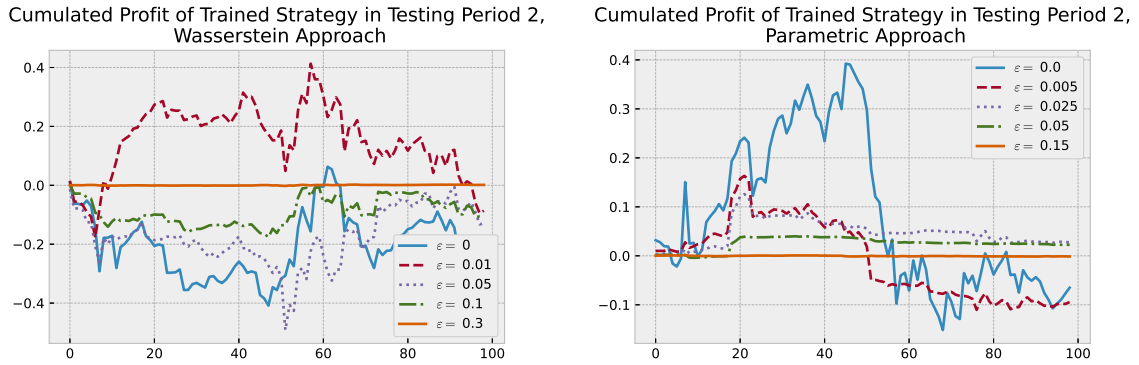


FIGURE 4. The figure shows the cumulated training profit of the trained strategies in testing period 2. The left panel illustrates the profit when applying a Wasserstein-ball approach, whereas the right panel illustrates the profit under a parametric approach.

	Overall Profit	Average Profit	% of profitable trades	Sharpe Ratio	Sortino Ratio
Wasserstein Approach					
$\varepsilon = 0$	-0.319652	-0.003229	51.52	-0.068955	-0.091312
$\varepsilon = 0.01$	-0.086229	-0.000871	56.57	-0.018503	-0.025662
$\varepsilon = 0.05$	-0.128121	-0.001294	52.53	-0.031648	-0.044057
$\varepsilon = 0.1$	-0.108500	-0.001096	50.51	-0.054067	-0.074651
$\varepsilon = 0.3$	0.001415	0.000014	52.53	0.027997	0.040511
Parametric Approach					
$\varepsilon = 0$	-0.065187	-0.000658	48.48	-0.014956	-0.020914
$\varepsilon = 0.005$	-0.094396	-0.000953	50.51	-0.063190	-0.079788
$\varepsilon = 0.025$	0.027990	0.000283	55.56	0.030187	0.053250
$\varepsilon = 0.05$	0.022454	0.000227	49.49	0.079555	0.277870
$\varepsilon = 0.15$	-0.001429	-0.000014	49.49	-0.048469	-0.061703

TABLE 2. The table shows the results of the Wasserstein-ball approach (Section 4.1) and the parametric approach (Section 4.2) in the second testing period.

Testing Period 3

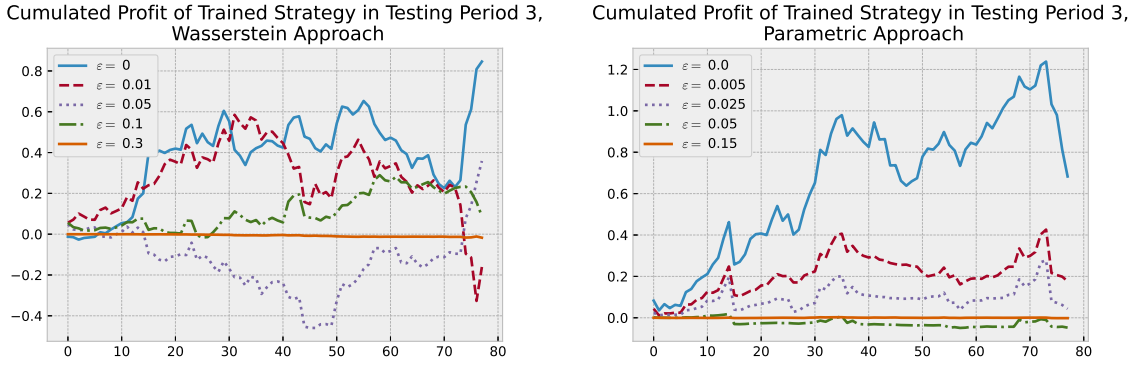


FIGURE 5. The figure shows the cumulated training profit of the trained strategies in testing period 3. The left panel illustrates the profit when applying a Wasserstein-ball approach, whereas the right panel illustrates the profit under a parametric approach.

	Overall Profit	Average Profit	% of profitable trades	Sharpe Ratio	Sortino Ratio
Wasserstein Approach					
$\varepsilon = 0$	0.845952	0.010846	52.56	0.166630	0.335301
$\varepsilon = 0.01$	-0.160733	-0.002061	48.72	-0.031032	-0.040478
$\varepsilon = 0.05$	0.358808	0.004600	52.56	0.096431	0.154537
$\varepsilon = 0.1$	0.087670	0.001124	52.56	0.037463	0.053469
$\varepsilon = 0.3$	-0.017214	-0.000221	43.59	-0.199209	-0.231336
Parametric Approach					
$\varepsilon = 0$	0.682794	0.008754	64.10	0.124028	0.170558
$\varepsilon = 0.005$	0.176339	0.002261	61.54	0.054989	0.069144
$\varepsilon = 0.025$	0.043324	0.000555	64.10	0.014557	0.017213
$\varepsilon = 0.05$	-0.047627	-0.000611	52.56	-0.064649	-0.075748
$\varepsilon = 0.15$	-0.001740	-0.000022	51.28	-0.039336	-0.049270

TABLE 3. The table shows the results of the Wasserstein-ball approach (Section 4.1) and the parametric approach (Section 4.2) in the third testing period.

4.3.4. *Discussion of the Results.* Note that in contrast to the third bullish testing period, the first and second testing periods comprise scenarios that did not occur in similar form during the training period. Hence, it cannot be guaranteed that a non-robust trading strategy that was trained in periods where such scenarios never occurred can trade profitably during testing period 1 and testing period 2. Indeed, as the numerical results reveal, applying a trained non-robust trading strategy not respecting any distributional ambiguity (i.e. $\varepsilon = 0$) leads to significant losses in the adverse scenarios considered in testing period 1 and 2, while robust trading strategies, which did encounter also adverse scenarios during training, clearly outperform the non-robust strategy for both considered approaches, namely the Wasserstein-ball approach and the parametric approach.

As the third testing period comprises a bullish market period, occurring in similar form in the training data, the non-robust approach turns out to be the most profitable approach in this period, since it is the approach that is best adjusted to this scenario. However, when choosing the right level of ε , the robust approach still can trade profitably in this period.

In general, it turns out to be important to correctly identify the appropriate *size* of the ambiguity set (here encoded by the radius ε). Due to the formulation of the robust optimization problem as a worst-case approach, respecting for more distributional ambiguity means to consider more *bad* scenarios, and therefore may eventually result in a more careful, less volatile trading behavior with smaller returns, as it clearly can be seen in the case $\varepsilon = 0.3$ for the Wasserstein-ball approach and

in the case $\varepsilon = 0.15$ for the parametric approach, respectively. In contrast, insufficiently accounting for uncertainty comes at the cost of not being well equipped when adverse scenarios occur, an observation that was already made in similar empirical studies that use different approaches to solve robust optimization problems, compare, e.g., [30] and [32]. Hence, choosing an *intermediate* level of ambiguity seems to be an appropriate choice. Supporting this rationale, our numerical results show that choosing a level of $\varepsilon = 0.05$ in the Wasserstein-ball approach and of $\varepsilon = 0.025$ in the parametric approach, respectively, lead to trading strategies that outperform in testing periods 1 and 2 the respective non-robust trading strategies and still lead to profits in testing period 3.

This provides evidence that taking into account distributional uncertainty may be of particular importance in volatile and crisis-like market scenarios which did not occur previously in a similar form.

5. PROOF OF THEOREM 2.7

Proof of Theorem 2.7 (i). Let $v \in C_p(\Omega_{\text{loc}}, \mathbb{R})$. We define the map

$$(5.1) \quad F : \text{Gr } \mathcal{P} = \{(x, a_0, \mathbb{P}_0) \mid x \in \Omega_{\text{loc}}, a_0 \in A, \mathbb{P}_0 \in \mathcal{P}(x, a_0)\} \rightarrow \mathbb{R}$$

$$(x, a_0, \mathbb{P}_0) \mapsto \int_{\Omega_{\text{loc}}} r(x, a_0, \omega_1) + \alpha v(\omega_1) \mathbb{P}_0(d\omega_1).$$

We claim that the map F is continuous. Indeed, to see this, let $(x^{(n)}, a_0^{(n)}, \mathbb{P}_0^{(n)}) \subseteq \text{Gr } \mathcal{P}$ with $(x^{(n)}, a_0^{(n)}, \mathbb{P}_0^{(n)}) \rightarrow (x, a_0, \mathbb{P}_0) \in \Omega_{\text{loc}} \times A \times \mathcal{M}_1(\Omega_{\text{loc}})$ for $n \rightarrow \infty$. Since $\Omega_{\text{loc}} \times A \ni (\tilde{x}, \tilde{a}) \mapsto \mathcal{P}(\tilde{x}, \tilde{a})$ is compact-valued and continuous (Assumption 2.2 (i)), we have $\mathbb{P}_0 \in \mathcal{P}(x, a)$. Moreover,

$$(5.2) \quad |F(x^{(n)}, a_0^{(n)}, \mathbb{P}_0^{(n)}) - F(x, a_0, \mathbb{P}_0)|$$

$$\leq |F(x^{(n)}, a_0^{(n)}, \mathbb{P}_0^{(n)}) - F(x, a_0, \mathbb{P}_0^{(n)})| + |F(x, a_0, \mathbb{P}_0^{(n)}) - F(x, a_0, \mathbb{P}_0)|.$$

The second summand $|F(x, a_0, \mathbb{P}_0^{(n)}) - F(x, a_0, \mathbb{P}_0)|$ vanishes for $n \rightarrow \infty$ since the integrand $\omega_1 \mapsto r(x, a_0, \omega_1) + \alpha v(\omega_1)$ is an element of $C_p(\Omega_{\text{loc}}, \mathbb{R})$. For the first summand we obtain, by using Assumption 2.4 (ii), that

$$\lim_{n \rightarrow \infty} |F(x^{(n)}, a_0^{(n)}, \mathbb{P}_0^{(n)}) - F(x, a_0, \mathbb{P}_0^{(n)})|$$

$$\leq \lim_{n \rightarrow \infty} \int_{\Omega_{\text{loc}}} |r(x^{(n)}, a_0^{(n)}, \omega_1) - r(x, a_0, \omega_1)| \mathbb{P}_0^{(n)}(d\omega_1)$$

$$\leq \lim_{n \rightarrow \infty} \int_{\Omega_{\text{loc}}} L \cdot \max\{1, \|\omega_1\|^p\} \cdot \left(\|x^{(n)} - x\| + \|a_0^{(n)} - a_0\| \right) \mathbb{P}_0^{(n)}(d\omega_1)$$

$$= \lim_{n \rightarrow \infty} L \cdot \left(\|x^{(n)} - x\| + \|a_0^{(n)} - a_0\| \right) \cdot \int_{\Omega_{\text{loc}}} \max\{1, \|\omega_1\|^p\} \mathbb{P}_0(d\omega_1) = 0,$$

where we use in the last equality that due to the convergence $\mathbb{P}_0^{(n)} \xrightarrow{\tau_p} \mathbb{P}_0$ also the p -th moments converge, see e.g. [38, Definition 6.8. (i), and Theorem 6.9]. Thus, F is continuous. Therefore, we may apply Berge's maximum theorem (Theorem B.2) and obtain that the map

$$(5.3) \quad G : \Omega_{\text{loc}} \times A \rightarrow \mathbb{R}$$

$$(x, a_0) \mapsto \inf_{\mathbb{P}_0 \in \mathcal{P}(x, a_0)} \int_{\Omega_{\text{loc}}} r(x, a_0, \omega_1) + \alpha v(\omega_1) \mathbb{P}_0(d\omega_1)$$

is continuous, and the set of minimizers is nonempty for each $(x, a_0) \in \Omega_{\text{loc}} \times A$. According to the measurable maximum theorem ([2, Theorem 18.19]) there exists a measurable selector $\Omega_{\text{loc}} \times A \ni (x, a_0) \mapsto \mathbb{P}_0^*(x, a_0) \in \mathcal{P}(x, a_0)$, where $\mathbb{P}_0^*(x, a_0)$ minimizes the integral in (5.3) for each $(x, a_0) \in \Omega_{\text{loc}} \times A$. This shows the assertion in (2.7). Next, we use that $\Omega_{\text{loc}} \times A \ni (x, a_0) \mapsto G(x, a_0)$ is continuous, and apply again Berge's maximum theorem to the constant compact-valued correspondence $\Omega_{\text{loc}} \ni x \rightarrow A \subset \mathbb{R}^m$ and to G . This yields that the map

$$(5.4) \quad H : \Omega_{\text{loc}} \rightarrow \mathbb{R}$$

$$x \mapsto \sup_{a_0 \in A} G(x, a_0).$$

is continuous and that the set of maximizers of $G(x, \cdot)$ is nonempty for each $x \in \Omega_{\text{loc}}$. Moreover, by the measurable maximum theorem ([2, Theorem 18.19]) there exists a measurable selector $\Omega_{\text{loc}} \ni x \mapsto a_{\text{loc}}^*(x) \in A$, where $a_{\text{loc}}^*(x)$ maximizes $G(x, \cdot)$ for each $x \in \Omega_{\text{loc}}$. This shows (2.8). Next, we define the map $\mathbb{P}_{\text{loc}}^* \in \mathbf{P}_{a_{\text{loc}}^*}$ by

$$\begin{aligned} \mathbb{P}_{\text{loc}}^* : \Omega_{\text{loc}} &\rightarrow \mathcal{M}_1(\Omega_{\text{loc}}) \\ x &\mapsto \mathbb{P}_0^*(x, a_{\text{loc}}^*(x)) \end{aligned}$$

and obtain for each $x \in \Omega_{\text{loc}}$ that

$$\begin{aligned} (5.5) \quad \mathcal{T}v(x) = H(x) &= \sup_{a_0 \in A} \inf_{\mathbb{P}_0 \in \mathcal{P}(x, a_0)} \int_{\Omega_{\text{loc}}} r(x, a_0, \omega_1) + \alpha v(\omega_1) \mathbb{P}_0(d\omega_1) \\ &= \inf_{\mathbb{P}_0 \in \mathcal{P}(x, a_0)} \int_{\Omega_{\text{loc}}} r(x, a_{\text{loc}}^*(x), \omega_1) + \alpha v(\omega_1) \mathbb{P}_0(d\omega_1) \\ &= \int_{\Omega_{\text{loc}}} r(x, a_{\text{loc}}^*(x), \omega_1) + \alpha v(\omega_1) \mathbb{P}_{\text{loc}}^*(x; d\omega_1) \\ &= \sup_{a_{\text{loc}} \in \mathcal{A}_{\text{loc}}} \inf_{\mathbb{P} \in \mathbf{P}_{a_{\text{loc}}}} \int_{\Omega_{\text{loc}}} r(x, a_{\text{loc}}(x), \omega_1) + \alpha v(\omega_1) \mathbb{P}(x; d\omega_1). \end{aligned}$$

This shows (2.9) and therefore completes the proof of (i). \square

Proof of Theorem 2.7 (ii). The continuity of $\Omega_{\text{loc}} \ni x \mapsto \mathcal{T}v(x)$ follows from the continuity of H defined in (5.4) and (5.5). By the growth conditions on r and X_1 (Assumption 2.4 (iii)), we obtain for all $x \in \Omega_{\text{loc}}$ that

$$\begin{aligned} \mathcal{T}v(x) &\leq \sup_{a_{\text{loc}} \in \mathcal{A}_{\text{loc}}} \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}}} \mathbb{E}_{\mathbb{P}_0} [C_r(1 + \|x\|^p + \|X_1\|^p) + \alpha \|v\|_{C_p} (1 + \|X_1\|^p)] \\ &\leq C_r (\|x\|^p + C_P(1 + \|x\|^p)) + \alpha \|v\|_{C_p} C_P(1 + \|x\|^p) \\ &\leq (C_r + C_r C_P + \alpha \|v\|_{C_p} C_P) (1 + \|x\|^p). \end{aligned}$$

Hence, $\mathcal{T}v \in C_p(\Omega_{\text{loc}}, \mathbb{R})$. Note that for every nonempty set \mathcal{Q} and for all $G, H : \mathcal{Q} \rightarrow \mathbb{R}$ we have that

$$(5.6) \quad \inf_{\mathcal{Q} \in \mathcal{Q}} G(\mathcal{Q}) - \inf_{\mathcal{Q} \in \mathcal{Q}} H(\mathcal{Q}) \leq \sup_{\mathcal{Q} \in \mathcal{Q}} |G(\mathcal{Q}) - H(\mathcal{Q})|.$$

Therefore, we obtain for every $v, w \in C_p(\Omega_{\text{loc}}, \mathbb{R}), x \in \Omega_{\text{loc}}$, by using (2.2), that

$$\begin{aligned} \mathcal{T}v(x) - \mathcal{T}w(x) &\leq \alpha \sup_{a \in A} \sup_{\mathbb{P}_0 \in \mathcal{P}(x, a)} \left| \mathbb{E}_{\mathbb{P}_0} [v(X_1) - w(X_1)] \right| \\ &\leq \alpha \sup_{a \in A} \sup_{\mathbb{P}_0 \in \mathcal{P}(x, a)} \mathbb{E}_{\mathbb{P}_0} [|v(X_1) - w(X_1)|] \\ &\leq \alpha \sup_{a \in A} \sup_{\mathbb{P}_0 \in \mathcal{P}(x, a)} \mathbb{E}_{\mathbb{P}_0} [\|v - w\|_{C_p} (1 + \|X_1\|^p)] \\ &\leq \alpha C_P \|v - w\|_{C_p} \cdot (1 + \|x\|^p). \end{aligned}$$

By interchanging the roles of v and w we obtain (2.10). Now, let $v_0 \in C_p(\Omega_{\text{loc}}, \mathbb{R})$. Then, since by Assumption 2.4 (iv) we have $0 < \alpha C_P < 1$, \mathcal{T} is a contraction on $C_p(\Omega_{\text{loc}}, \mathbb{R})$. Hence Banach's fix point theorem (Theorem B.1) implies existence and uniqueness of a fix point $v \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ such that $v = \mathcal{T}v = \lim_{n \rightarrow \infty} \mathcal{T}^n v_0$. \square

Proof of Theorem 2.7 (iii). We first show the inequality $v(x) \leq V(x)$.

Let $v \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ be the fix point of \mathcal{T} whose existence and uniqueness were proved in part (ii), i.e., we have that $v = \mathcal{T}v$. Let $a_{\text{loc}}^* \in \mathcal{A}_{\text{loc}}$ be the minimizer from part (i) and write $\mathbf{a}^* := (a_{\text{loc}}^*(X_0), a_{\text{loc}}^*(X_1), \dots) \in \mathcal{A}$, and let $x \in \Omega_{\text{loc}}$. First, we claim it holds for all $n \in \mathbb{N}$ with $n \geq 2$ that

$$(5.7) \quad \mathcal{T}^n v(x) \leq \inf_{\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}^*}} \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{n-1} \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) + \alpha^n v(X_n^x) \right].$$

We prove the claim (5.7) inductively. To that end, note that by part (i) we can write

$$\begin{aligned} \mathcal{T}v(x) &= \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}^*}} \mathbb{E}_{\mathbb{P}_0(x)} \left[r(x, a_{\text{loc}}^*(x), X_1) + \alpha v(X_1) \right] \\ &= \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}^*}} \int_{\Omega_{\text{loc}}} r(x, a_{\text{loc}}^*(x), \omega_1) + \alpha v(\omega_1) \mathbb{P}_0(x; d\omega_1). \end{aligned}$$

This implies for all $t \in \mathbb{N}_0$ and $\omega = (\omega_t)_{t \in \mathbb{N}_0} \in \Omega$ that

$$\begin{aligned} (5.8) \quad \mathcal{T}v(X_t(\omega)) &= \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}^*}} \int_{\Omega_{\text{loc}}} r(X_t(\omega), a_{\text{loc}}^*(X_t(\omega)), \omega_1) + \alpha v(\omega_1) \mathbb{P}_0(X_t(\omega); d\omega_1) \\ &= \inf_{\mathbb{P}_t \in \mathbf{P}_{a_{\text{loc}}^*}} \int_{\Omega_{\text{loc}}} r(\omega_t, a_{\text{loc}}^*(\omega_t), \omega_{t+1}) + \alpha v(\omega_{t+1}) \mathbb{P}_t(\omega_t; d\omega_{t+1}), \end{aligned}$$

where we just used the definition of the canonical process and relabelled the variables \mathbb{P}_0 and ω_1 . For $n = 2$, as $\mathcal{T}v = v$, we thus have

$$\begin{aligned} (5.9) \quad \mathcal{T}(\mathcal{T}v)(x) &= \mathcal{T}v(x) = \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}^*}} \int_{\Omega_{\text{loc}}} r(x, a_{\text{loc}}^*(x), \omega_1) + \alpha v(x) \mathbb{P}_0(x; d\omega_1) \\ &= \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}^*}} \int_{\Omega_{\text{loc}}} r(x, a_{\text{loc}}^*(x), \omega_1) + \alpha \mathcal{T}v(x) \mathbb{P}_0(x; d\omega_1) \\ &= \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}^*}} \int_{\Omega_{\text{loc}}} \left[r(x, a_{\text{loc}}^*(x), \omega_1) \right. \\ &\quad \left. + \alpha \sup_{a_{\text{loc}} \in \mathcal{A}_{\text{loc}}} \inf_{\mathbb{P}_1 \in \mathbf{P}_{a_{\text{loc}}}} \int_{\Omega_{\text{loc}}} \left\{ r(\omega_1, a_{\text{loc}}(\omega_1), \omega_2) + \alpha v(\omega_2) \right\} \mathbb{P}_1(\omega_1; d\omega_2) \right] \mathbb{P}_0(x; d\omega_1). \end{aligned}$$

This and part (i) ensures that

$$\begin{aligned} \mathcal{T}(\mathcal{T}v)(x) &= \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}^*}} \int_{\Omega_{\text{loc}}} \left[r(x, a_{\text{loc}}^*(x), \omega_1) \right. \\ &\quad \left. + \alpha \inf_{\mathbb{P}_1 \in \mathbf{P}_{a_{\text{loc}}^*}} \int_{\Omega_{\text{loc}}} \left\{ r(\omega_1, a_{\text{loc}}^*(\omega_1), \omega_2) + \alpha v(\omega_2) \right\} \mathbb{P}_1(\omega_1; d\omega_2) \right] \mathbb{P}_0(x; d\omega_1) \\ &\leq \inf_{\substack{\mathbb{P}_t \in \mathbf{P}_{a_{\text{loc}}^*}, \\ t=0,1}} \int_{\Omega_{\text{loc}}} \int_{\Omega_{\text{loc}}} \sum_{t=0}^1 \alpha^t r(\omega_t, a_{\text{loc}}^*(\omega_t), \omega_{t+1}) + \alpha^2 v(\omega_2) \mathbb{P}_1(\omega_1; d\omega_2) \mathbb{P}_0(x; d\omega_1) \\ &= \inf_{\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}^*}} \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^1 \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) + \alpha^2 v(X_2) \right], \end{aligned}$$

where we use (5.8) and the structure of the measures $\mathfrak{P}_{x, \mathbf{a}^*}$. The general case for arbitrary n follows with analogue arguments. Indeed, let the claim in (5.7) be true for $n - 1$, then it follows by the

same argument as in (5.9) and by the structure of every $\mathbb{P} \in \mathfrak{P}_{\omega_1, \mathbf{a}^*}$ that

$$\begin{aligned}
\mathcal{T}^n v(x) &= \mathcal{T}(\mathcal{T}^{n-1}v)(x) \\
&\leq \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}^*}} \int_{\Omega_{\text{loc}}} \left[r(x, a_{\text{loc}}^*(x), \omega_1) \right. \\
&\quad \left. + \alpha \inf_{\mathbb{P} \in \mathfrak{P}_{\omega_1, \mathbf{a}^*}} \int_{\Omega} \left\{ \sum_{t=1}^{n-1} \alpha^{t-1} r(\omega_t, a_{\text{loc}}^*(\omega_t), \omega_{t+1}) + \alpha^{n-1} v(\omega_n) \right\} \mathbb{P}(d\omega) \right] \mathbb{P}_0(x; d\omega_1) \\
&= \inf_{\mathbb{P}_0 \in \mathbf{P}_{a_{\text{loc}}^*}} \int_{\Omega_{\text{loc}}} \left[r(x, a_{\text{loc}}^*(x), \omega_1) \right. \\
&\quad \left. + \alpha \inf_{\substack{\mathbb{P}_t \in \mathbf{P}_{a_{\text{loc}}^*}, \\ t=1, \dots, n-1}} \int_{\Omega_{\text{loc}}} \cdots \int_{\Omega_{\text{loc}}} \left\{ \sum_{t=1}^{n-1} \alpha^{t-1} r(\omega_t, a_{\text{loc}}^*(\omega_t), \omega_{t+1}) \right. \right. \\
&\quad \left. \left. + \alpha^{n-1} v(\omega_n) \right\} \mathbb{P}_{n-1}(\omega_{n-1}; d\omega_n) \cdots \mathbb{P}_1(\omega_1; d\omega_2) \right] \mathbb{P}_0(x; d\omega_1) \\
&\leq \inf_{\substack{\mathbb{P}_t \in \mathbf{P}_{a_{\text{loc}}^*}, \\ t=0, \dots, n-1}} \int_{\Omega_{\text{loc}}} \int_{\Omega_{\text{loc}}} \cdots \int_{\Omega_{\text{loc}}} \left\{ \sum_{t=0}^{n-1} \alpha^t r(\omega_t, a_{\text{loc}}^*(\omega_t), \omega_{t+1}) \right. \\
&\quad \left. + \alpha^n v(\omega_n) \right\} \mathbb{P}_{n-1}(\omega_{n-1}; d\omega_n) \cdots \mathbb{P}_1(\omega_1; d\omega_2) \mathbb{P}_0(x; d\omega_1) \\
&= \inf_{\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}^*}} \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{n-1} \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) + \alpha^n v(X_n) \right].
\end{aligned}$$

According to (5.7) we have for all $n \in \mathbb{N}$ that

$$(5.10) \quad v(x) = \mathcal{T}v(x) = \mathcal{T}^n v(x) \leq \inf_{\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}^*}} \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{n-1} \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) + \alpha^n v(X_n) \right].$$

Further, we obtain for all $\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}^*}$ and $n \in \mathbb{N}$ by a repeated application of (2.2) that

$$\mathbb{E}_{\mathbb{P}} \left[|v(X_n)| \right] \leq \|v\|_{C_P} C_P^n (1 + \|x\|^p).$$

Note that Assumption 2.4 implies $0 < C_P \cdot \alpha < \frac{1}{C_r(C_P+1)} < 1$. When letting $n \rightarrow \infty$, we thus have for all $\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}}$ that

$$(5.11) \quad 0 \leq \limsup_{n \rightarrow \infty} \mathbb{E}_{\mathbb{P}} \left[\alpha^n |v(X_n)| \right] \leq \|v\|_{C_P} (1 + \|x\|^p) \cdot \limsup_{n \rightarrow \infty} (C_P \cdot \alpha)^n = 0.$$

Moreover, note that by the growth condition on r in Assumption 2.4 (iii) we have for each $n \in \mathbb{N}$ that

$$(5.12) \quad \sum_{t=0}^{n-1} \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) \leq \sum_{t=0}^{\infty} \alpha^t C_r (1 + \|X_t\|^p + \|X_{t+1}\|^p).$$

Let $\delta := \frac{1}{\alpha C_r(C_P+1)} - C_P$ which satisfies $\delta > 0$ by Assumption 2.4 (iv). Then for all $\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}^*}$, by using Assumption 2.4 (iv), (2.2), and Beppo Levi's theorem, we have that

$$\begin{aligned}
(5.13) \quad \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{\infty} \alpha^t C_r (1 + \|X_t\|^p + \|X_{t+1}\|^p) \right] &\leq \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{\infty} \alpha^t (C_r(1 + C_P)) (1 + \|X_t\|^p) \right] \\
&\leq \sum_{t=0}^{\infty} \alpha^t (C_r(1 + C_P))^{t+1} (1 + \|x\|^p) \\
&= C_r(1 + C_P) \left[\sum_{t=0}^{\infty} \left(\frac{1}{C_P + \delta} \right)^t \right] \cdot (1 + \|x\|^p) \\
&= C_r(1 + C_P) \frac{(C_P + \delta)(1 + \|x\|^p)}{C_P + \delta - 1} < \infty.
\end{aligned}$$

Hence the dominating function in (5.12) is integrable and we obtain, by using the dominated convergence theorem and (5.11), that

$$\begin{aligned}
(5.14) \quad v(x) &\leq \limsup_{n \rightarrow \infty} \inf_{\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}^*}} \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{n-1} \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) + \alpha^n v(X_n) \right] \\
&\leq \inf_{\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}^*}} \limsup_{n \rightarrow \infty} \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{n-1} \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) \right] + \limsup_{n \rightarrow \infty} \mathbb{E}_{\mathbb{P}} \left[\alpha^n |v(X_n)| \right] \\
&= \inf_{\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}^*}} \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{\infty} \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) \right] \leq V(x).
\end{aligned}$$

Next, we show the inequality $v(x) \geq V(x)$. To this end, let $\mathbb{P}_0^* : \Omega_{\text{loc}} \times A \rightarrow \mathcal{M}_1(\Omega_{\text{loc}})$ be defined as in part (i) with respect to the unique fix point $v \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ of \mathcal{T} . Moreover, for every $\mathbf{a} = (a_t)_{t \in \mathbb{N}_0} \in \mathcal{A}$ let $\mathbb{P}_{x, \mathbf{a}}^* := \delta_x \otimes \mathbb{P}_{a_0}^* \otimes \mathbb{P}_{a_1}^* \otimes \cdots$, where for $t \in \mathbb{N}$ we define $\mathbb{P}_{a_t}^* : \Omega_{\text{loc}} \ni \omega_t \mapsto \mathbb{P}_0^*(\omega_t, a_t(\omega_t)) \in \mathcal{P}(\omega_t, a_t(\omega_t))$. Thus, we have, by using the dominated convergence theorem with the same dominating function as in (5.12), that

$$\begin{aligned}
(5.15) \quad V(x) &= \sup_{\mathbf{a} \in \mathcal{A}} \inf_{\mathbb{P} \in \mathfrak{P}_{x, \mathbf{a}}} \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{\infty} \alpha^t r(X_t, a_t(X_t), X_{t+1}) \right] \\
&\leq \sup_{\mathbf{a} \in \mathcal{A}} \mathbb{E}_{\mathbb{P}_{x, \mathbf{a}}^*} \left[\sum_{t=0}^{\infty} \alpha^t r(X_t, a_t(X_t), X_{t+1}) \right] \\
&= \sup_{\mathbf{a} \in \mathcal{A}} \sum_{t=0}^{\infty} \mathbb{E}_{\mathbb{P}_{x, \mathbf{a}}^*} \left[\alpha^t r(X_t, a_t(X_t), X_{t+1}) \right] \\
&= \sup_{\mathbf{a} \in \mathcal{A}} \sum_{t=0}^{\infty} \left(\alpha^t \mathbb{E}_{\mathbb{P}_{x, \mathbf{a}}^*} \left[r(X_t, a_t(X_t), X_{t+1}) + \alpha v(X_{t+1}) \right] - \mathbb{E}_{\mathbb{P}_{x, \mathbf{a}}^*} \left[\alpha^{t+1} v(X_{t+1}) \right] \right) \\
&= \sup_{\mathbf{a} \in \mathcal{A}} \sum_{t=0}^{\infty} \left(\alpha^t \int_{\Omega_{\text{loc}}} \cdots \int_{\Omega_{\text{loc}}} r(\omega_t, a_t(\omega_t), \omega_{t+1}) + \alpha v(\omega_{t+1}) \mathbb{P}_0^*(\omega_t, a_t(\omega_t); d\omega_{t+1}) \cdots \mathbb{P}_0^*(x, a_0(x); d\omega_1) \right. \\
&\quad \left. - \mathbb{E}_{\mathbb{P}_{x, \mathbf{a}}^*} \left[\alpha^{t+1} v(X_{t+1}) \right] \right).
\end{aligned}$$

Moreover, by using the results from part (i) we have for all $\omega_t \in \Omega_{\text{loc}}$

$$\begin{aligned}
(5.16) \quad &\int_{\Omega_{\text{loc}}} r(\omega_t, a_t(\omega_t), \omega_{t+1}) + \alpha v(\omega_{t+1}) \mathbb{P}_0^*(\omega_t, a_t(\omega_t); d\omega_{t+1}) \\
&= \inf_{\mathbb{P}_0 \in \mathcal{P}(\omega_t, a_t(\omega_t))} \int_{\Omega_{\text{loc}}} r(\omega_t, a_t(\omega_t), \omega_{t+1}) + \alpha v(\omega_{t+1}) \mathbb{P}_0(\omega_t, a_t(\omega_t); d\omega_{t+1}) \\
&\leq \sup_{a_{\text{loc}} \in \mathcal{A}_{\text{loc}}} \inf_{\mathbb{P}_0 \in \mathcal{P}(\omega_t, a_{\text{loc}}(\omega_t))} \int_{\Omega_{\text{loc}}} r(\omega_t, a_{\text{loc}}(\omega_t), \omega_{t+1}) + \alpha v(\omega_{t+1}) \mathbb{P}_0(\omega_t, a_{\text{loc}}(\omega_t); d\omega_{t+1}) \\
&= \int_{\Omega_{\text{loc}}} r(\omega_t, a_{\text{loc}}^*(\omega_t), \omega_{t+1}) + \alpha v(\omega_{t+1}) \mathbb{P}_0^*(\omega_t, a_{\text{loc}}^*(\omega_t); d\omega_{t+1}) = \mathcal{T}v(\omega_t) = v(\omega_t).
\end{aligned}$$

Hence, we obtain with (5.15) and (5.16) that

$$\begin{aligned}
V(x) &\leq \sup_{\mathbf{a} \in \mathcal{A}} \sum_{t=0}^{\infty} \left(\alpha^t \int_{\Omega_{\text{loc}}} \cdots \int_{\Omega_{\text{loc}}} v(\omega_t) \mathbb{P}_0^*(\omega_{t-1}, a_{t-1}(\omega_{t-1}); d\omega_t) \cdots \mathbb{P}_0^*(x, a_0(x); d\omega_1) \right. \\
&\quad \left. - \mathbb{E}_{\mathbb{P}_{x, \mathbf{a}}^*} \left[\alpha^{t+1} v(X_{t+1}) \right] \right) \\
&= \sup_{\mathbf{a} \in \mathcal{A}} \sum_{t=0}^{\infty} \left(\alpha^t \mathbb{E}_{\mathbb{P}_{x, \mathbf{a}}^*} \left[v(X_t) \right] - \alpha^{t+1} \mathbb{E}_{\mathbb{P}_{x, \mathbf{a}}^*} \left[v(X_{t+1}) \right] \right) \\
&= \sup_{\mathbf{a} \in \mathcal{A}} v(x) = v(x).
\end{aligned}$$

This shows $V(x) = v(x)$. Eventually, to see that the first line of (2.11) holds, we compute by using (5.16), the definition $\mathbb{P}_x^* := \delta_x \otimes \mathbb{P}_{\text{loc}}^* \otimes \mathbb{P}_{\text{loc}}^* \otimes \cdots$ as well as the dominated convergence theorem that

$$\begin{aligned} v(x) &= \sum_{t=0}^{\infty} \left(\alpha^t \mathbb{E}_{\mathbb{P}_x^*} \left[v(X_t) \right] - \alpha^{t+1} \mathbb{E}_{\mathbb{P}_x^*} \left[v(X_{t+1}) \right] \right) \\ &= \sum_{t=0}^{\infty} \left(\alpha^t \mathbb{E}_{\mathbb{P}_x^*} \left[r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) + \alpha v(X_{t+1}) \right] - \alpha^{t+1} \mathbb{E}_{\mathbb{P}_x^*} \left[v(X_{t+1}) \right] \right) \\ &= \sum_{t=0}^{\infty} \mathbb{E}_{\mathbb{P}_x^*} \left[\alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) \right] = \mathbb{E}_{\mathbb{P}_x^*} \left[\sum_{t=0}^{\infty} \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) \right]. \end{aligned}$$

Moreover, by (5.14), as we have shown that $V = v$, we obtain that

$$(5.17) \quad V(x) = \inf_{\mathbb{P} \in \mathfrak{F}_{x, a^*}} \mathbb{E}_{\mathbb{P}} \left[\sum_{t=0}^{\infty} \alpha^t r(X_t, a_{\text{loc}}^*(X_t), X_{t+1}) \right].$$

□

6. PROOF OF RESULTS IN SECTION 3

6.1. Proof of Results in Section 3.1.

Proof of Proposition 3.1. Let $x \in \Omega_{\text{loc}}$, $a \in A$.

We see that $\mathcal{P}(x, a)$ is nonempty since the measure $\widehat{\mathbb{P}}(x, a)$ is contained in $\mathcal{P}(x, a)$ by definition of the q -Wasserstein-ball.

The compactness of $\mathcal{B}_{\varepsilon}^{(q)}(\widehat{\mathbb{P}}(x, a))$ with respect to τ_0 , which is the topology induced by the weak convergence of measures, follows from, e.g., [43, Theorem 1], where we use the assumption that $\widehat{\mathbb{P}}(x, a)$ has finite q -th moments.

To show the upper hemicontinuity of \mathcal{P} , we apply Lemma B.3. Let $(x^{(n)}, a^{(n)})_{n \in \mathbb{N}} \subseteq \Omega_{\text{loc}} \times A$ such that $(x^{(n)}, a^{(n)}) \rightarrow (x, a) \in \Omega_{\text{loc}} \times A$ for $n \rightarrow \infty$. Further, consider a sequence $(\mathbb{P}^{(n)})_{n \in \mathbb{N}}$ such that $\mathbb{P}^{(n)} \in \mathcal{B}_{\varepsilon}^{(q)}(\widehat{\mathbb{P}}(x^{(n)}, a^{(n)}))$ for all $n \in \mathbb{N}$, i.e., we have $((x^{(n)}, a^{(n)}), \mathbb{P}^{(n)})_{n \in \mathbb{N}} \subseteq \text{Gr } \mathcal{P}$.

Let $(\delta_n)_{n \in \mathbb{N}} \subseteq (0, 1)$ with $\lim_{n \rightarrow \infty} \delta_n = 0$. Note that, since $\Omega_{\text{loc}} \times A \ni (x, a) \mapsto \widehat{\mathbb{P}}(x, a)$ is, by assumption, continuous in τ_q we have $\lim_{n \rightarrow \infty} W_q(\widehat{\mathbb{P}}(x, a), \widehat{\mathbb{P}}(x^{(n)}, a^{(n)})) = 0$. Hence, there exists a subsequence $(\widehat{\mathbb{P}}(x^{(n_k)}, a^{(n_k)}))_{k \in \mathbb{N}}$ such that

$$(6.1) \quad W_q(\widehat{\mathbb{P}}(x, a), \widehat{\mathbb{P}}(x^{(n_k)}, a^{(n_k)})) < \delta_k \cdot \varepsilon \text{ for all } k \in \mathbb{N}.$$

This implies for each $\mathbb{P}^{(n_k)}$, $k \in \mathbb{N}$, that

$$W_q(\widehat{\mathbb{P}}(x, a), \mathbb{P}^{(n_k)}) \leq W_q(\widehat{\mathbb{P}}(x, a), \widehat{\mathbb{P}}(x^{(n_k)}, a^{(n_k)})) + W_q(\widehat{\mathbb{P}}(x^{(n_k)}, a^{(n_k)}), \mathbb{P}^{(n_k)}) \leq \delta_k \cdot \varepsilon + \varepsilon \leq 2\varepsilon.$$

Hence, $\mathbb{P}^{(n_k)} \in \mathcal{B}_{2\varepsilon}^{(q)}(\widehat{\mathbb{P}}(x, a))$ for all $k \in \mathbb{N}$. By the compactness of $\mathcal{B}_{2\varepsilon}^{(q)}(\widehat{\mathbb{P}}(x, a))$ in τ_0 , there exists a subsequence $(\mathbb{P}^{(n_{k_\ell})})_{\ell \in \mathbb{N}}$ such that $\mathbb{P}^{(n_{k_\ell})} \xrightarrow{\tau_0} \mathbb{P}$ as $\ell \rightarrow \infty$ for some $\mathbb{P} \in \mathcal{B}_{2\varepsilon}^{(q)}(\widehat{\mathbb{P}}(x, a))$. In particular, since by assumption $\widehat{\mathbb{P}}(x, a)$ possesses finite q -th moments, \mathbb{P} has also finite q -th moments, see [43, Lemma 1]. It remains to prove that $\mathbb{P} \in \mathcal{B}_{\varepsilon}^{(q)}(\widehat{\mathbb{P}}(x, a))$. To that end, define for each $k \in \mathbb{N}$

$$(6.2) \quad \widetilde{\mathbb{P}}^{(n_k)} := (1 - \delta_k) \cdot \mathbb{P}^{(n_k)} + \delta_k \cdot \widehat{\mathbb{P}}(x^{(n_k)}, a^{(n_k)}).$$

Then, for each $k \in \mathbb{N}$ we have

$$\begin{aligned} (6.3) \quad & W_q(\widehat{\mathbb{P}}(x^{(n_k)}, a^{(n_k)}), \widetilde{\mathbb{P}}^{(n_k)}) \\ &= W_q\left((1 - \delta_k) \cdot \widehat{\mathbb{P}}(x^{(n_k)}, a^{(n_k)}) + \delta_k \cdot \widehat{\mathbb{P}}(x^{(n_k)}, a^{(n_k)}), (1 - \delta_k) \cdot \mathbb{P}^{(n_k)} + \delta_k \cdot \widehat{\mathbb{P}}(x^{(n_k)}, a^{(n_k)})\right) \\ &= (1 - \delta_k) \cdot W_q(\widehat{\mathbb{P}}(x^{(n_k)}, a^{(n_k)}), \mathbb{P}^{(n_k)}) \leq (1 - \delta_k) \cdot \varepsilon. \end{aligned}$$

Therefore, by (6.1) and (6.3) we have for each $\ell \in \mathbb{N}$ that

$$(6.4) \quad \begin{aligned} W_q \left(\widehat{\mathbb{P}}(x, a), \widetilde{\mathbb{P}}^{(n_{k\ell})} \right) &\leq W_q \left(\widehat{\mathbb{P}}(x, a), \widehat{\mathbb{P}}(x^{(n_{k\ell})}, a^{(n_{k\ell})}) \right) + W_q \left(\widehat{\mathbb{P}}(x^{(n_{k\ell})}, a^{(n_{k\ell})}), \widetilde{\mathbb{P}}^{(n_{k\ell})} \right) \\ &\leq \delta_{k\ell} \cdot \varepsilon + (1 - \delta_{k\ell}) \cdot \varepsilon = \varepsilon. \end{aligned}$$

Furthermore, we have by (6.2) that

$$(6.5) \quad \lim_{\ell \rightarrow \infty} \widetilde{\mathbb{P}}^{(n_{k\ell})} = \lim_{\ell \rightarrow \infty} \mathbb{P}^{(n_{k\ell})} = \mathbb{P} \text{ in } \tau_0.$$

Since $\mu \mapsto W_q(\widehat{\mathbb{P}}(x, a), \mu)$ is lower semicontinuous in τ_0 , see [15, Corollary 5.3], we obtain from (6.4) and (6.5) that

$$W_q \left(\widehat{\mathbb{P}}(x, a), \mathbb{P} \right) \leq \liminf_{\ell \rightarrow \infty} W_q \left(\widehat{\mathbb{P}}(x, a), \widetilde{\mathbb{P}}^{(n_{k\ell})} \right) \leq \varepsilon,$$

and hence $\mathbb{P} \in \mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a))$. The assertion that \mathcal{P} is upper hemicontinuous follows now with the characterization of upper hemicontinuity provided in Lemma B.3.

To show the lower hemicontinuity of \mathcal{P} we first define the set-valued map

$$\overset{\circ}{\mathcal{P}} : \Omega_{\text{loc}} \times A \ni (x, a) \rightarrow \mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a)) := \left\{ \mathbb{P} \in \mathcal{M}_1(\Omega_{\text{loc}}) \mid W_q(\mathbb{P}, \widehat{\mathbb{P}}(x, a)) < \varepsilon \right\}$$

and conclude the lower hemicontinuity of $\overset{\circ}{\mathcal{P}}$ with Lemma B.4. To this end, we consider a sequence $(x^{(n)}, a^{(n)})_{n \in \mathbb{N}} \subset \Omega_{\text{loc}} \times A$ such that $(x^{(n)}, a^{(n)}) \rightarrow (x, a) \in \Omega_{\text{loc}} \times A$ for $n \rightarrow \infty$, and we consider some $\mathbb{P} \in \overset{\circ}{\mathcal{P}}((x, a)) = \mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a))$. Note that since $\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a))$ is defined as an open ball with respect to τ_q , there exists some $0 < \delta < \varepsilon$ such that $\mathbb{P} \in \mathcal{B}_{\varepsilon-\delta}^{(q)}(\widehat{\mathbb{P}}(x, a))$. We define for $n \in \mathbb{N}$ the measure

$$\mathbb{P}^{(n)} := \begin{cases} \widehat{\mathbb{P}}(x^{(n)}, a^{(n)}), & \text{if } W_q \left(\widehat{\mathbb{P}}(x^{(n)}, a^{(n)}), \widehat{\mathbb{P}}(x, a) \right) \geq \delta \\ \mathbb{P}, & \text{else.} \end{cases}$$

Then, we claim that $\mathbb{P}^{(n)} \in \overset{\circ}{\mathcal{P}}((x^{(n)}, a^{(n)}))$ for all $n \in \mathbb{N}$. Indeed, if $W_q \left(\widehat{\mathbb{P}}(x^{(n)}, a^{(n)}), \widehat{\mathbb{P}}(x, a) \right) \geq \delta$ this follows by definition of $\mathbb{P}^{(n)}$, whereas if $W_q \left(\widehat{\mathbb{P}}(x^{(n)}, a^{(n)}), \widehat{\mathbb{P}}(x, a) \right) < \delta$, then $\mathbb{P}^{(n)} = \mathbb{P}$, and hence by the triangle inequality

$$W_q \left(\mathbb{P}, \widehat{\mathbb{P}}(x^{(n)}, a^{(n)}) \right) \leq W_q \left(\mathbb{P}, \widehat{\mathbb{P}}(x, a) \right) + W_q \left(\widehat{\mathbb{P}}(x, a), \widehat{\mathbb{P}}(x^{(n)}, a^{(n)}) \right) < (\varepsilon - \delta) + \delta = \varepsilon.$$

By the continuity of $(x, a) \mapsto \widehat{\mathbb{P}}(x, a)$ in τ_q we have that $\widehat{\mathbb{P}}(x^{(n)}, a^{(n)}) \xrightarrow{\tau_q} \widehat{\mathbb{P}}(x, a)$ as $n \rightarrow \infty$. Thus, there exists some $N \in \mathbb{N}$ such that we have $\mathbb{P}^{(n)} = \mathbb{P}$ for all $n \geq N$ and thus, in particular $\mathbb{P}^{(n)} \rightarrow \mathbb{P}$ weakly for $n \rightarrow \infty$, which concludes the lower hemicontinuity of $\overset{\circ}{\mathcal{P}}$ with Lemma B.4. Next, we claim that the τ_0 -closure of $\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a))$, denoted by $\text{cl}_{\tau_0} \left(\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a)) \right)$, coincides with $\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a))$. Indeed, the inclusion $\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a)) \subseteq \text{cl}_{\tau_0} \left(\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a)) \right)$ follows, since $\text{cl}_{\tau_0} \left(\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a)) \right)$ is closed in τ_0 and hence also in τ_q . To show the reverse inclusion $\text{cl}_{\tau_0} \left(\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a)) \right) \subseteq \mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a))$ let $\mathbb{P} \in \text{cl}_{\tau_0} \left(\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a)) \right)$. Then, there exists a sequence $(\mathbb{P}^{(n)})_{n \in \mathbb{N}} \subseteq \mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a))$ with $\mathbb{P}^{(n)} \xrightarrow{\tau_0} \mathbb{P}$ as $n \rightarrow \infty$. Hence by using the lower semicontinuity of $\mu \mapsto W_q(\mu, \widehat{\mathbb{P}}(x, a))$ with respect to τ_0 we obtain

$$W_q \left(\mathbb{P}, \widehat{\mathbb{P}}(x, a) \right) \leq \liminf_{n \rightarrow \infty} W_q \left(\mathbb{P}^{(n)}, \widehat{\mathbb{P}}(x, a) \right) \leq \varepsilon.$$

Hence, $\text{cl}_{\tau_0} \left(\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a)) \right) = \mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a))$ and [2, Lemma 17.22] implies that the set-valued map $\mathcal{P} : \Omega_{\text{loc}} \times A \ni (x, a) \rightarrow \text{cl}_{\tau_0} \left(\mathcal{B}_\varepsilon^{(q)}(\widehat{\mathbb{P}}(x, a)) \right)$ is lower hemicontinuous.

Eventually, since $p = 0$, the growth constraint (2.2) is automatically fulfilled. \square

6.2. Proof of Results in Section 3.2.

Proof of Proposition 3.2. Let $(x, a) \in \Omega_{\text{loc}} \times A$.

The nonemptiness of $\mathcal{P}(x, a)$ follows directly since Θ is nonempty.

To show the compactness of $\mathcal{P}(x, a)$ let $(\mathbb{P}^{(n)})_{n \in \mathbb{N}} \subseteq \mathcal{P}(x, a)$, i.e., for all $n \in \mathbb{N}$ we have $\mathbb{P}^{(n)} = \widehat{\mathbb{P}}(x, a, \theta^{(n)})$ for some $\theta^{(n)} \in \Theta(x, a)$. The compactness of $\Theta(x, a)$ implies the existence

of a subsequence $(\theta^{(n_k)})_{k \in \mathbb{N}} \subseteq \Theta(x, a)$ such that $\theta^{(n_k)} \rightarrow \theta \in \Theta(x, a)$ for $k \rightarrow \infty$. Hence, since $\widehat{\mathbb{P}}$ is continuous, it follows $\widehat{\mathbb{P}}(x, a, \theta^{(n_k)}) \rightarrow \widehat{\mathbb{P}}(x, a, \theta) \in \mathcal{P}(x, a)$ in τ_p for $k \rightarrow \infty$.

We apply Lemma B.3 to show the upper hemicontinuity of \mathcal{P} . To this end, consider a sequence $(x^{(n)}, a^{(n)})_{n \in \mathbb{N}} \subseteq \Omega_{\text{loc}} \times A$ with $\lim_{n \rightarrow \infty} (x^{(n)}, a^{(n)}) = (x, a)$ and a sequence $(\mathbb{P}^{(n)})_{n \in \mathbb{N}}$ with $\mathbb{P}^{(n)} \in \mathcal{P}(x^{(n)}, a^{(n)})$ for all $n \in \mathbb{N}$. We have a representation $\mathbb{P}^{(n)} = \widehat{\mathbb{P}}(x^{(n)}, a^{(n)}, \theta^{(n)})$ for some $\theta^{(n)} \in \Theta(x^{(n)}, a^{(n)})$ for all $n \in \mathbb{N}$. Then, since Θ is upper hemicontinuous, there exists a subsequence $(\theta^{(n_k)})_{k \in \mathbb{N}}$ with $\theta^{(n_k)} \in \Theta(x^{(n_k)}, a^{(n_k)})$ for all $k \in \mathbb{N}$ such that $\theta^{(n_k)} \rightarrow \theta$ for $k \rightarrow \infty$ for some $\theta \in \Theta(x, a)$. Hence with the continuity of $\widehat{\mathbb{P}}$ it follows $\mathbb{P}^{(n_k)} \rightarrow \mathbb{P} := \widehat{\mathbb{P}}(x, a, \theta) \in \mathcal{P}(x, a)$ in τ_p for $k \rightarrow \infty$.

To show the lower hemicontinuity we let $(x^{(n)}, a^{(n)})_{n \in \mathbb{N}} \subseteq \Omega_{\text{loc}} \times A$ with $\lim_{n \rightarrow \infty} (x^{(n)}, a^{(n)}) = (x, a)$ and $\mathbb{P} := \widehat{\mathbb{P}}(x, a, \theta) \in \mathcal{P}(x, a)$ for some $\theta \in \Theta(x, a)$. Then, the lower hemicontinuity of Θ implies the existence of a subsequence $(x^{(n_k)}, a^{(n_k)})_{k \in \mathbb{N}}$ and of a sequence $(\theta^{(k)})_{k \in \mathbb{N}}$ with $\theta^{(k)} \in \Theta(x^{(n_k)}, a^{(n_k)})$ for all $k \in \mathbb{N}$ such that $\theta^{(k)} \rightarrow \theta$ for $k \rightarrow \infty$. Hence, it follows with the continuity of $\widehat{\mathbb{P}}$ that $\mathcal{P}(x^{(n_k)}, a^{(n_k)}) \ni \mathbb{P}^{(k)} := \widehat{\mathbb{P}}(x^{(n_k)}, a^{(n_k)}, \theta^{(k)}) \rightarrow \mathbb{P}$ for $k \rightarrow \infty$, implying with Lemma B.4 the lower hemicontinuity of \mathcal{P} . \square

6.3. Proof of Results in Section 3.3. Before reporting the proof of Proposition 3.3, we establish the following lemma.

Lemma 6.1. *Let $\mathcal{D} := \{\delta_x \mid x \in T^{m-1}\} \subseteq (\mathcal{M}_1(T^{m-1}), \tau_p)$ be the closed subset consisting of all Dirac measures on T^{m-1} . Then, for any $p \in \{0, 1\}$, the map*

$$\begin{aligned} \varphi : (\mathcal{D}, \tau_p) \times (\mathcal{M}_1(T), \tau_p) &\rightarrow (\mathcal{M}_1(T^m), \tau_p) \\ (\delta_x, \mathbb{P}) &\mapsto \delta_x \otimes \mathbb{P} \end{aligned}$$

is continuous.

Proof. First, we consider the case $p = 0$. Let $(\delta_{x^{(n)}})_{n \in \mathbb{N}} \subseteq \mathcal{D}$, and $(\mathbb{P}^{(n)})_{n \in \mathbb{N}} \subseteq \mathcal{M}_1(T)$ with $\delta_{x^{(n)}} \xrightarrow{\tau_0} \mu$ and $\mathbb{P}^{(n)} \xrightarrow{\tau_0} \mathbb{P} \in \mathcal{M}_1(T)$ for $n \rightarrow \infty$. Note that, as \mathcal{D} is closed, $\mu = \delta_x$ for some $x \in T^{m-1}$. Now, let $f : T^m \rightarrow \mathbb{R}$ be Lipschitz continuous with Lipschitz constant $L > 0$. Then we have

$$\begin{aligned} &\lim_{n \rightarrow \infty} \left| \int_{T^m} f(y, z) \delta_{x^{(n)}}(dy) \otimes \mathbb{P}^{(n)}(dz) - \int_{T^m} f(y, z) \delta_x(dy) \otimes \mathbb{P}(dz) \right| \\ &= \lim_{n \rightarrow \infty} \left| \int_T f(x^{(n)}, z) \mathbb{P}^{(n)}(dz) - \int_T f(x, z) \mathbb{P}(dz) \right| \\ &\leq \lim_{n \rightarrow \infty} \left(\int_T |f(x^{(n)}, z) - f(x, z)| \mathbb{P}^{(n)}(dz) + \left| \int_T f(x, z) \mathbb{P}^{(n)}(dz) - \int_T f(x, z) \mathbb{P}(dz) \right| \right) \\ (6.6) \quad &\leq \lim_{n \rightarrow \infty} L \cdot \|x^{(n)} - x\| + \lim_{n \rightarrow \infty} \left| \int_T f(x, z) \mathbb{P}^{(n)}(dz) - \int_T f(x, z) \mathbb{P}(dz) \right| = 0, \end{aligned}$$

where the second summand in (6.6) vanishes due to $\mathbb{P}^{(n)} \xrightarrow{\tau_0} \mathbb{P}$. By [25, Theorem 18.7] we conclude that φ is continuous.

Now, we consider the case $p = 1$. Let $(\delta_{x^{(n)}})_{n \in \mathbb{N}} \subseteq \mathcal{D}$, and $(\mathbb{P}^{(n)})_{n \in \mathbb{N}} \subseteq \mathcal{M}_1(T)$ with $\delta_{x^{(n)}} \xrightarrow{\tau_1} \delta_x \in \mathcal{D}$ and $\mathbb{P}^{(n)} \xrightarrow{\tau_1} \mathbb{P} \in \mathcal{M}_1(T)$ for $n \rightarrow \infty$. Since convergence in τ_1 implies convergence in τ_0 , we obtain, by the already considered case $p = 0$, that $\delta_{x^{(n)}} \otimes \mathbb{P}^{(n)} \xrightarrow{\tau_0} \delta_x \otimes \mathbb{P}$ for $n \rightarrow \infty$. It remains to show that the convergence also follows with respect to τ_1 . To conclude the convergence in τ_1 it suffices, by [38, Theorem 6.9], to show that

$$\lim_{n \rightarrow \infty} \int_{T^m} \|(y, z)\| \delta_{x^{(n)}}(dy) \otimes \mathbb{P}^{(n)}(dz) = \int_{T^m} \|(y, z)\| \delta_x(dy) \otimes \mathbb{P}(dz).$$

To see this, note that

$$\begin{aligned}
& \lim_{n \rightarrow \infty} \left| \int_{T^m} \|(y, z)\| \delta_{x^{(n)}}(dy) \otimes \mathbb{P}^{(n)}(dz) - \int_{T^m} \|(y, z)\| \delta_x(dy) \otimes \mathbb{P}(dz) \right| \\
&= \lim_{n \rightarrow \infty} \left| \int_T \|(x^{(n)}, z)\| \mathbb{P}^{(n)}(dz) - \int_T \|(x, z)\| \mathbb{P}(dz) \right| \\
&\leq \lim_{n \rightarrow \infty} \left(\int_T \left| \|(x^{(n)}, z)\| - \|(x, z)\| \right| \mathbb{P}^{(n)}(dz) + \left| \int_T \|(x, z)\| \mathbb{P}^{(n)}(dz) - \int_T \|(x, z)\| \mathbb{P}(dz) \right| \right) \\
&\leq \lim_{n \rightarrow \infty} \left(\|x^{(n)} - x\| + \left| \int_T \|(x, z)\| \mathbb{P}^{(n)}(dz) - \int_T \|(x, z)\| \mathbb{P}(dz) \right| \right) = 0,
\end{aligned}$$

where we use that $T \ni y \mapsto \|(x, y)\| \in C_1(T, \mathbb{R})$ and $\mathbb{P}^{(n)} \xrightarrow{\tau_1} \mathbb{P}$ for $n \rightarrow \infty$. \square

Proof of Proposition 3.3. Let $(x, a) \in \Omega_{\text{loc}} \times A$.

It is immediate that $\mathcal{P}(x, a) \neq \emptyset$, since $\tilde{\mathcal{P}}(x, a) \neq \emptyset$ by assumption.

To show the compactness of $\mathcal{P}(x, a)$ we consider a sequence $(\mathbb{P}^{(n)})_{n \in \mathbb{N}} \subseteq \mathcal{P}(x, a)$, where for all $n \in \mathbb{N}$ we have $\mathbb{P}^{(n)} = \delta_{\pi(x)} \otimes \tilde{\mathbb{P}}^{(n)}$ for some $\tilde{\mathbb{P}}^{(n)} \in \tilde{\mathcal{P}}(x, a)$. Then, by the compactness of $\tilde{\mathcal{P}}(x, a)$ there exists a subsequence $(\tilde{\mathbb{P}}^{(n_k)})_{k \in \mathbb{N}}$ such that $\tilde{\mathbb{P}}^{(n_k)} \rightarrow \tilde{\mathbb{P}} \in \tilde{\mathcal{P}}(x, a)$ in τ_p as $k \rightarrow \infty$. Now, let $g \in C_p(\Omega_{\text{loc}}, \mathbb{R})$. Then the map $T \ni y \mapsto g(\pi(x), y)$ is contained in $C_p(T, \mathbb{R})$, and hence

$$\begin{aligned}
\lim_{k \rightarrow \infty} \int_{\Omega_{\text{loc}}} g(z) \mathbb{P}^{(n_k)}(dz) &= \lim_{k \rightarrow \infty} \int_T g(\pi(x), y) \tilde{\mathbb{P}}^{(n_k)}(dy) \\
&= \int_T g(\pi(x), y) \tilde{\mathbb{P}}(dy) = \int_{\Omega_{\text{loc}}} g(z) \mathbb{P}(dz)
\end{aligned}$$

for $\mathbb{P} := \delta_{\pi(x)} \otimes \tilde{\mathbb{P}} \in \mathcal{P}(x, a)$, which proves the compactness of $\mathcal{P}(x, a)$.

To show the upper hemicontinuity of \mathcal{P} , let $(x^{(n)}, a^{(n)}) \subseteq \Omega_{\text{loc}} \times A$ with $(x^{(n)}, a^{(n)}) \rightarrow (x, a)$ for $n \rightarrow \infty$, and let $\mathbb{P}^{(n)} \in \mathcal{P}(x^{(n)}, a^{(n)})$ for all $n \in \mathbb{N}$. Then, we have a representation $\mathbb{P}^{(n)} = \delta_{\pi(x^{(n)})} \otimes \tilde{\mathbb{P}}^{(n)}$ with $\tilde{\mathbb{P}}^{(n)} \in \tilde{\mathcal{P}}(x^{(n)}, a^{(n)})$ for all $n \in \mathbb{N}$. By the upper hemicontinuity of $\tilde{\mathcal{P}}$, there exists, according to Lemma B.3, a subsequence $(\tilde{\mathbb{P}}^{(n_k)})_{k \in \mathbb{N}}$ with $\tilde{\mathbb{P}}^{(n_k)} \rightarrow \tilde{\mathbb{P}} \in \tilde{\mathcal{P}}(x, a)$ in τ_p as $k \rightarrow \infty$. Moreover $\delta_{\pi(x^{(n)})} \rightarrow \delta_{\pi(x)}$ in τ_1 as $n \rightarrow \infty$.

We apply Lemma 6.1 and obtain that $\delta_{\pi(x^{(n_k)})} \otimes \tilde{\mathbb{P}}^{(n_k)} \rightarrow \delta_{\pi(x)} \otimes \tilde{\mathbb{P}} \in \mathcal{P}(x, a)$, and hence the upper hemicontinuity follows with Lemma B.3.

To prove the lower hemicontinuity of \mathcal{P} we consider again a sequence $(x^{(n)}, a^{(n)}) \subseteq \Omega_{\text{loc}} \times A$ with $(x^{(n)}, a^{(n)}) \rightarrow (x, a)$ for $n \rightarrow \infty$, and some $\mathbb{P} \in \mathcal{P}(x, a)$ with a representation $\mathbb{P} = \delta_{\pi(x)} \otimes \tilde{\mathbb{P}}$ for $\tilde{\mathbb{P}} \in \tilde{\mathcal{P}}(x, a)$. By the lower hemicontinuity of $\tilde{\mathcal{P}}$ there exists, according to Lemma B.4, a subsequence $(x^{(n_k)}, a^{(n_k)})_{k \in \mathbb{N}}$ and $\tilde{\mathbb{P}}^{(n_k)} \in \tilde{\mathcal{P}}(x^{(n_k)}, a^{(n_k)})$ for all $k \in \mathbb{N}$ such that $\tilde{\mathbb{P}}^{(n_k)} \rightarrow \tilde{\mathbb{P}}$ in τ_p . Then, we set $\mathbb{P}^{(n_k)} := \delta_{\pi(x^{(n_k)})} \otimes \tilde{\mathbb{P}}^{(n_k)} \in \mathcal{P}(x^{(n_k)}, a^{(n_k)})$ for all $k \in \mathbb{N}$, and we conclude $\mathbb{P}^{(n_k)} \rightarrow \mathbb{P}$ in τ_p for $k \rightarrow \infty$ with Lemma 6.1. Hence, the lower hemicontinuity follows with Lemma B.4. \square

7. PROOF OF RESULTS IN SECTION 4

7.1. Proof of Results in Section 4.1.

Proof of Proposition 4.1. Since Assumption 2.2 (ii) is automatically fulfilled for $p = 0$, the fulfilment of Assumption 2.2 follows from Proposition 3.1 and Proposition 3.3, once we have shown that $\Omega_{\text{loc}} \ni x \mapsto \hat{\mathbb{P}}(x) \in \mathcal{M}_1(T)$ is continuous in τ_q and possesses finite q -th moments.

To show the continuity of $\hat{\mathbb{P}}$, we consider a sequence $(X_t^{(n)})_{n \in \mathbb{N}} \subseteq \Omega_{\text{loc}}$ with $X_t^{(n)} \rightarrow X_t \in \Omega_{\text{loc}}$ for $n \rightarrow \infty$. By construction $\Omega_{\text{loc}} \ni x \mapsto \pi_s(x) \in [0, 1]$ is continuous for all $s = m, \dots, N - 1$, which

implies for all $g \in C_q(T, \mathbb{R})$ that

$$\begin{aligned} \lim_{n \rightarrow \infty} \int_T g(y) \widehat{\mathbb{P}}(X_t^{(n)}; dy) &= \lim_{n \rightarrow \infty} \sum_{s=m}^{N-1} \pi_s(X_t^{(n)}) g(\mathcal{R}_{s+1}) \\ &= \sum_{s=m}^{N-1} \pi_s(X_t) g(\mathcal{R}_{s+1}) = \int_T g(y) \widehat{\mathbb{P}}(X_t; dy). \end{aligned}$$

Moreover, the existence of the q -th moment follows by

$$\int_T \|y\|^q \widehat{\mathbb{P}}(X_t; dy) = \sum_{s=m}^{N-1} \pi_s(X_t) \cdot \|\mathcal{R}_{s+1}\|^q < \infty.$$

Now, to verify Assumption 2.4 note that r is continuous and that the compactness of T and of A imply that r is bounded, and thus Assumption 2.4 (i) and (iii) are fulfilled. Next, let $X_t, X'_t \in \Omega_{\text{loc}}$, $X_{t+1} = (\mathcal{R}_{t-m+2}, \dots, \mathcal{R}_{t+1}) \in \Omega_{\text{loc}}$, and let $a_t, a'_t \in A$. Then, by the Cauchy–Schwarz inequality we see that

$$|r(X_t, a_t, X_{t+1}) - r(X'_t, a'_t, X_{t+1})| = \left| \sum_{i=1}^D (a_t^i - a'_t{}^i) \mathcal{R}_{t+1}^i \right| \leq \|\mathcal{R}_{t+1}\| \cdot \|a_t - a'_t\| \leq \max_{x \in T} \|x\| \cdot \|a_t - a'_t\|$$

and hence Assumption 2.4 (ii) is fulfilled. \square

7.2. Proof of Results in Section 4.2.

Proof of Proposition 4.2. To show Assumption 2.2 (ii), let

$$C_P := 1 + \sqrt{\varepsilon^2 + \frac{1}{m} + 4 \cdot \frac{\varepsilon^2 + 1}{m-1}}$$

Now we consider some $(x, a) \in \Omega_{\text{loc}} \times A$ and some $\mathbb{P} \in \mathcal{P}(x, a)$. Then, we have a representation of the form $\mathbb{P} = \delta_{\pi(x)} \otimes \widetilde{\mathbb{P}}$ for some $\widetilde{\mathbb{P}} \in \widetilde{\mathcal{P}}(x)$, where $\widetilde{\mathbb{P}} \sim \mathcal{N}_D(\mu, \Sigma)$ with $(\mu, \Sigma) \in \mathbb{R}^D \times \mathbb{R}^{D \times D}$ fulfilling $\|\mu - \mathbf{m}(x)\| \leq \varepsilon$ and $\Sigma = \mathbf{c}(y)$ for some $y \in \Omega_{\text{loc}}$ with $\|y - x\| \leq \varepsilon$. Therefore, we have with Jensen's inequality that

$$\begin{aligned} \int_{\Omega_{\text{loc}}} 1 + \|y\| \mathbb{P}(dy) &= 1 + \int_{\mathbb{R}^D} \|(\pi(x), z)\| \widetilde{\mathbb{P}}(dz) \\ (7.1) \qquad \qquad \qquad &\leq 1 + \int_{\mathbb{R}^D} \|\pi(x)\| + \|z\| \widetilde{\mathbb{P}}(dz) \\ &\leq 1 + \|x\| + \sqrt{\int_{\mathbb{R}^D} \|z\|^2 \widetilde{\mathbb{P}}(dz)}. \end{aligned}$$

Moreover, for $Z = (Z_1, \dots, Z_D) \sim \mathcal{N}_D(\mu, \Sigma)$ we have

$$\begin{aligned} \int_{\mathbb{R}^D} \|z\|^2 \widetilde{\mathbb{P}}(dz) &= \mathbb{E}[\|Z\|^2] = \mathbb{E}\left[\sum_{i=1}^D Z_i^2\right] \\ (7.2) \qquad \qquad \qquad &= \sum_{i=1}^D \mathbb{E}[Z_i^2] + \sum_{i=1}^D \text{Var}(Z_i) \\ &= \|\mu\|^2 + \text{trace}(\Sigma). \end{aligned}$$

In the next step, we write $x = (x_i^{(j)})_{i=1, \dots, m}^{j=1, \dots, D}$, use the Cauchy–Schwarz inequality, and compute

$$\begin{aligned}
\|\mu\|^2 &\leq (\|\mu - \mathbf{m}(x)\| + \|\mathbf{m}(x)\|)^2 \\
&\leq \left(\varepsilon + \sqrt{\sum_{j=1}^D \left(\frac{1}{m} \sum_{i=1}^m x_i^{(j)} \right)^2} \right)^2 \\
(7.3) \quad &\leq \left(\varepsilon + \sqrt{\sum_{j=1}^D \left(\frac{1}{m} \sum_{i=1}^m (x_i^{(j)})^2 \right)} \right)^2 \\
&= \left(\varepsilon + \frac{1}{\sqrt{m}} \|x\| \right)^2 \leq \left(\varepsilon^2 + \frac{1}{m} \right) (1 + \|x\|^2).
\end{aligned}$$

Further, we write $y = (y_i^{(j)})_{i=1, \dots, m}^{j=1, \dots, D}$, and obtain with the Cauchy–Schwarz inequality that

$$\begin{aligned}
\|\mathbf{m}(y) - \mathbf{m}(x)\|^2 &= \frac{1}{m^2} \sum_{j=1}^D \left(\sum_{i=1}^m (y_i^{(j)} - x_i^{(j)}) \right)^2 \\
(7.4) \quad &\leq \frac{1}{m} \sum_{j=1}^D \sum_{i=1}^m (y_i^{(j)} - x_i^{(j)})^2 = \frac{\|y - x\|^2}{m} \leq \frac{\varepsilon^2}{m}.
\end{aligned}$$

The above inequality (7.4), and $\|\mathbf{m}(x)\| \leq \frac{1}{\sqrt{m}} \|x\|$ (see also (7.3)) imply together with the Cauchy–Schwarz inequality that

$$\begin{aligned}
\text{trace}(\Sigma) &= \frac{1}{m-1} \sum_{i=1}^m \text{trace}((y_i - \mathbf{m}(y))(y_i - \mathbf{m}(y))^T) \\
&= \frac{1}{m-1} \sum_{i=1}^m \sum_{j=1}^D (y_i^{(j)} - \mathbf{m}(y)^{(j)})^2 \\
&\leq \frac{2}{m-1} \sum_{i=1}^m \sum_{j=1}^D \left[(y_i^{(j)})^2 + (\mathbf{m}(y)^{(j)})^2 \right] \\
(7.5) \quad &= \frac{2}{m-1} \|y\|^2 + \frac{2m}{m-1} \|\mathbf{m}(y)\|^2 \\
&\leq \frac{2}{m-1} (\|y - x\| + \|x\|)^2 + \frac{2m}{m-1} (\|\mathbf{m}(y) - \mathbf{m}(x)\| + \|\mathbf{m}(x)\|)^2 \\
&\leq \frac{2}{m-1} (\varepsilon + \|x\|)^2 + \frac{2m}{m-1} \left(\frac{\varepsilon}{\sqrt{m}} + \frac{1}{\sqrt{m}} \|x\| \right)^2 \\
&\leq \frac{2}{m-1} (\varepsilon^2 + 1) (1 + \|x\|^2) + \frac{2m}{m-1} \cdot \frac{\varepsilon^2 + 1}{m} \cdot (1 + \|x\|^2) \\
&= 4 \cdot \frac{\varepsilon^2 + 1}{m-1} \cdot (1 + \|x\|^2).
\end{aligned}$$

Hence, by combining (7.1), (7.2), (7.3), and (7.5) we have

$$\begin{aligned}
\int_{\Omega_{\text{loc}}} 1 + \|y\| \mathbb{P}(dy) &\leq 1 + \|x\| + \sqrt{\left(\varepsilon^2 + \frac{1}{m} \right) (1 + \|x\|^2) + 4 \cdot \frac{\varepsilon^2 + 1}{m-1} \cdot (1 + \|x\|^2)} \\
&\leq \left(1 + \sqrt{\varepsilon^2 + \frac{1}{m} + 4 \cdot \frac{\varepsilon^2 + 1}{m-1}} \right) \cdot (1 + \|x\|) \\
&= C_P \cdot (1 + \|x\|),
\end{aligned}$$

as required in Assumption 2.2 (ii).

Since Assumption 2.2 (ii) is fulfilled, the fulfilment of Assumption 2.2 follows now with an application of Proposition 3.3. Thus, to verify the assumptions of Proposition 3.3, we need to show that $\Omega_{\text{loc}} \ni x \mapsto \tilde{\mathcal{P}}(x) \rightarrow (\mathcal{M}_1(\mathbb{R}^D), \tau_1)$ is nonempty, compact-valued, and continuous. This, in turn

follows from Proposition 3.2 once we have shown that $\Omega_{\text{loc}} \ni x \rightarrow \Theta(x) \subseteq \mathbb{R}^D \times \mathbb{R}^{D \times D}$ is nonempty, compact-valued, and continuous and that

$$(7.6) \quad \begin{aligned} \{(x, \mu, \Sigma) \mid x \in \Omega_{\text{loc}}, (\mu, \Sigma) \in \Theta(x)\} &\rightarrow (\mathcal{M}_1(\mathbb{R}^D), \tau_1) \\ (x, \mu, \Sigma) &\mapsto \widehat{\mathbb{P}}(x, \mu, \Sigma) := \mathcal{N}_D(\mu, \Sigma). \end{aligned}$$

is continuous.

To that end, let $x \in \Omega_{\text{loc}}$.

The non-emptiness of $\Theta(x)$ follows by definition.

To show the compactness of $\Theta(x)$, let $(\mu^{(n)}, \Sigma^{(n)})_{n \in \mathbb{N}} \subseteq \Theta(x)$. Then, we have $\|\mu^{(n)} - \mathbf{m}(x)\| \leq \varepsilon$ for all $n \in \mathbb{N}$ as well as $\Sigma^{(n)} = \mathbf{c}(y^{(n)})$ for some $y^{(n)} \in \Omega_{\text{loc}}$ with $\|y^{(n)} - x\| \leq \varepsilon$. Then, according to the Bolzano–Weierstrass theorem there exists a subsequence $(\mu^{(n_k)}, y^{(n_k)})_{k \in \mathbb{N}} \subseteq \mathbb{R}^D \times \Omega_{\text{loc}}$ such that $y^{(n_k)} \rightarrow y \in \Omega_{\text{loc}}$ with $\|y - x\| \leq \varepsilon$, and $\mu^{(n_k)} \rightarrow \mu \in \mathbb{R}^D$ with $\|\mu - \mathbf{m}(x)\| \leq \varepsilon$ for $k \rightarrow \infty$. Since \mathbf{c} is continuous we obtain that $(\mu^{(n_k)}, \Sigma^{(n_k)}) \rightarrow (\mu, \Sigma) := (\mu, \mathbf{c}(y)) \in \Theta(x)$ for $k \rightarrow \infty$.

To show the upper hemicontinuity of Θ , let $(x^{(n)})_{n \in \mathbb{N}} \subseteq \Omega_{\text{loc}}$ with $x^{(n)} \rightarrow x \in \Omega_{\text{loc}}$ for $n \rightarrow \infty$ as well as $(\mu^{(n)}, \Sigma^{(n)})_{n \in \mathbb{N}}$ with $(\mu^{(n)}, \Sigma^{(n)}) \in \Theta(x^{(n)})$ for all $n \in \mathbb{N}$. We have for all $n \in \mathbb{N}$ that $\|\mu^{(n)} - \mathbf{m}(x^{(n)})\| \leq \varepsilon$ and that $\Sigma^{(n)} = \mathbf{c}(y^{(n)})$ with $\|y^{(n)} - x^{(n)}\| \leq \varepsilon$ for some $y^{(n)} \in \Omega_{\text{loc}}$. Therefore, since $\|\mu^{(n)} - \mathbf{m}(x)\| \leq \|\mu^{(n)} - \mathbf{m}(x^{(n)})\| + \|\mathbf{m}(x^{(n)}) - \mathbf{m}(x)\|$, the continuity of \mathbf{m} ensures for every n large enough that $\|\mu^{(n)} - \mathbf{m}(x)\| \leq 2\varepsilon$. Hence, there exists according to the Bolzano–Weierstrass theorem a subsequence $(\mu^{(n_k)})_{k \in \mathbb{N}}$ with $\mu^{(n_k)} \rightarrow \mu$ for $k \rightarrow \infty$ for some $\mu \in \mathbb{R}^D$. Therefore, since \mathbf{m} is continuous, we obtain

$$(7.7) \quad \|\mu - \mathbf{m}(x)\| = \lim_{k \rightarrow \infty} \left\| \mu^{(n_k)} - \mathbf{m}(x^{(n_k)}) \right\| \leq \varepsilon.$$

Analogously, we have that $\|y^{(n)} - x\| \leq \|y^{(n)} - x^{(n)}\| + \|x^{(n)} - x\| < 2\varepsilon$ for every n large enough. This implies the existence of a subsequence $(y^{(n_k)})_{k \in \mathbb{N}}$ converging against some $y \in \Omega_{\text{loc}}$ with $\|y - x\| = \lim_{k \rightarrow \infty} \|y^{(n_k)} - x^{(n_k)}\| \leq \varepsilon$. Then, for $\Sigma := \mathbf{c}(y)$ we have $(\mu, \Sigma) \in \Theta(x)$ and $(\mu^{(n_k)}, \Sigma^{(n_k)}) \rightarrow (\mu, \Sigma)$ for $k \rightarrow \infty$. Thus, the upper hemicontinuity follows with Lemma B.3.

To show the lower hemicontinuity of Θ we consider a sequence $(x^{(n)})_{n \in \mathbb{N}} \subseteq \Omega_{\text{loc}}$ with $x^{(n)} \rightarrow x \in \Omega_{\text{loc}}$ for $n \rightarrow \infty$ and some $(\mu, \Sigma) \in \Theta(x)$. We have by definition $\|\mu - \mathbf{m}(x)\| \leq \varepsilon$ as well as $\Sigma = \mathbf{c}(y)$ for some $y \in \mathbb{R}^D$ with $\|y - x\| \leq \varepsilon$. We define for every $n \in \mathbb{N}$

$$\mu^{(n)} := \left(1 - \frac{1}{n}\right) \mu + \frac{1}{n} \mathbf{m}(x^{(n)}).$$

Then, due to the convergence $\mathbf{m}(x^{(n)}) \rightarrow \mathbf{m}(x)$, there exists a subsequence $(x^{(n_k)})_{k \in \mathbb{N}}$ such that for every $k \in \mathbb{N}$ we have $\|\mathbf{m}(x^{(n_k)}) - \mathbf{m}(x)\| < \varepsilon/(n_k - 1)$. This implies for all $k \in \mathbb{N}$ that

$$(7.8) \quad \begin{aligned} \|\mu^{(n_k)} - \mathbf{m}(x^{(n_k)})\| &= \left(1 - \frac{1}{n_k}\right) \left\| \mu - \mathbf{m}(x^{(n_k)}) \right\| \\ &\leq \left(1 - \frac{1}{n_k}\right) \left(\|\mu - \mathbf{m}(x)\| + \left\| \mathbf{m}(x) - \mathbf{m}(x^{(n_k)}) \right\| \right) \leq \left(1 - \frac{1}{n_k}\right) \left(\varepsilon + \frac{\varepsilon}{n_k - 1} \right) = \varepsilon. \end{aligned}$$

Next, we define for all $n \in \mathbb{N}$

$$(7.9) \quad y^{(n)} := \left(1 - \frac{1}{n}\right) y + \frac{1}{n} x^{(n)}, \quad \Sigma^{(n)} := \mathbf{c}(y^{(n)}).$$

We obtain by the convergence $x^{(n)} \rightarrow x$ for $n \rightarrow \infty$ the existence of a subsequence $x^{(n_{k_l})}$ such that $\|x^{(n_{k_l})} - x\| < \varepsilon/(n_{k_l} - 1)$ for all $l \in \mathbb{N}$. This implies

$$(7.10) \quad \begin{aligned} \|y^{(n_{k_l})} - x^{(n_{k_l})}\| &= \left(1 - \frac{1}{n_{k_l}}\right) \cdot \|y - x^{(n_{k_l})}\| \\ &\leq \left(1 - \frac{1}{n_{k_l}}\right) \cdot \left(\|y - x\| + \|x - x^{(n_{k_l})}\| \right) \\ &\leq \left(1 - \frac{1}{n_{k_l}}\right) \left(\varepsilon + \frac{\varepsilon}{n_{k_l} - 1} \right) = \varepsilon. \end{aligned}$$

Hence, with (7.8), (7.9), and (7.10), we have shown the existence of a subsequence $(\mu^{(n_{kl})}, \Sigma^{(n_{kl})})_{l \in \mathbb{N}}$ with $(\mu^{(n_{kl})}, \Sigma^{(n_{kl})}) \in \Theta(x^{(n_{kl})})$ for all $l \in \mathbb{N}$ and such that, by the continuity of r , $(\mu^{(n_{kl})}, \Sigma^{(n_{kl})}) \rightarrow (\mu, \Sigma)$ for $l \rightarrow \infty$. This implies the lower hemicontinuity of Θ by Lemma B.4.

It remains to show that the map defined in (7.6) is continuous with respect to τ_1 .

To that end, consider a sequence $(x^{(n)})_{n \in \mathbb{N}} \subseteq \Omega_{\text{loc}}$ as well as a sequence $(\mu^{(n)}, \Sigma^{(n)})_{n \in \mathbb{N}}$ with $(\mu^{(n)}, \Sigma^{(n)}) \in \Theta(x^{(n)})$ for all $n \in \mathbb{N}$ and such that $(x^{(n)}, \mu^{(n)}, \Sigma^{(n)}) \rightarrow (x, \mu, \Sigma) \in \Omega_{\text{loc}} \times \Theta(x)$ for $n \rightarrow \infty$. Then we write $\mathbb{P}^{(n)} := \mathcal{N}_D(\mu^{(n)}, \Sigma^{(n)}) \in \mathcal{M}_1(\mathbb{R}^D)$ for $n \in \mathbb{N}$ as well as $\mathbb{P} := \mathcal{N}_D(\mu, \Sigma) \in \mathcal{M}_1(\mathbb{R}^D)$. The characteristic function of $\mathbb{P}^{(n)}$, denoted by

$$\mathbb{R}^D \ni u \mapsto \varphi_{\mathbb{P}^{(n)}}(u) := \exp\left(iu^T \mu^{(n)} - \frac{1}{2}u^T \Sigma^{(n)}u\right)$$

converges for $n \rightarrow \infty$ pointwise against

$$\mathbb{R}^D \ni u \mapsto \varphi_{\mathbb{P}}(u) := \exp\left(iu^T \mu - \frac{1}{2}u^T \Sigma u\right),$$

which is the characteristic function of \mathbb{P} , and hence by Lévy's continuity theorem (see, e.g., [25, Theorem 19.1]) we have $\mathbb{P}^{(n)} \rightarrow \mathbb{P}$ weakly, i.e., in τ_0 for $n \rightarrow \infty$.

The convergence of $\mathbb{P}^{(n)} \rightarrow \mathbb{P}$ with respect to τ_1 now follows with, e.g., [10, Example 3.8.15], since $(\mathbb{P}^{(n)})_{n \in \mathbb{N}}$, and \mathbb{P} are Gaussian.

To verify Assumption 2.4 first note that r is continuous, and hence Assumption 2.4 (i) is fulfilled. Let $X_t, X'_t \in \Omega_{\text{loc}}$, $X_{t+1} = (\mathcal{R}_{t-m+2}, \dots, \mathcal{R}_{t+1}) \in \Omega_{\text{loc}}$, and let $a_t, a'_t \in A$. Then, the Cauchy–Schwarz inequality implies

$$\left| r(X_t, a_t, X_{t+1}) - r(X'_t, a'_t, X_{t+1}) \right| = \left| \sum_{i=1}^D (a_t^i - a'_t{}^i) \mathcal{R}_{t+1}^i \right| \leq \|\mathcal{R}_{t+1}\| \cdot \|a_t - a'_t\| \leq \|X_{t+1}\| \cdot \|a_t - a'_t\|,$$

implying Assumption 2.4 (ii). Moreover, we have by using the Cauchy–Schwarz inequality that

$$\left| r(X_t, a_t, X_{t+1}) \right| = \left| \sum_{i=1}^D a_t^i \mathcal{R}_{t+1}^i \right| \leq \|a_t\| \cdot \|\mathcal{R}_{t+1}\| \leq \max_{a \in A} \|a\| \cdot \|X_{t+1}\|,$$

as required in Assumption 2.4 (iii). \square

ACKNOWLEDGMENTS

Financial support by the MOE AcRF Tier 1 Grant *RG74/21* and by the Nanyang Assistant Professorship Grant (NAP Grant) *Machine Learning based Algorithms in Finance and Insurance* is gratefully acknowledged.

APPENDIX A. NUMERICS

We present an explicit numerical algorithm that can be applied to compute the optimal value function V and to determine an optimal policy $\mathbf{a}^* \in \mathcal{A}$.

A.1. Value Iteration. Theorem 2.7 directly provides an algorithm for the computation of the optimal value $V(x)$ to which we refer as the *value iteration algorithm*.

In this algorithm we start with an arbitrary $V^{(0)} \in C_p(\Omega_{\text{loc}}, \mathbb{R})$ and then compute recursively $V^{(n+1)} := \mathcal{T}V^{(n)}$ for all $n \in \mathbb{N}_0$. According to Theorem 2.7 we then have

$$(A.1) \quad \lim_{n \rightarrow \infty} \mathcal{T}V^{(n)} = V.$$

Compare also, e.g., [7, Section 7], where this algorithm (in a non-robust setting) is discussed in detail.

A.2. Numerical Algorithm. With Algorithm 1 we present a pseudocode of the methodology that can be applied to compute both the optimal value function and the optimal policy numerically. The algorithm relies on the value iteration principle. This means we solve (A.1) by approximating the value function V through neural networks⁷ and by repeatedly applying the recursion $V^{(n+1)} := \mathcal{T}V^{(n)}$. Note that Algorithm 1 approximates an optimal one-step action

⁷For a general introduction to neural networks we refer the reader to [28], for applications of neural networks in finance see [16], for a proof of the universal approximation property of neural networks compare [23].

$a_{\text{loc}}^* \in \mathcal{A}_{\text{loc}}$. According to Theorem 2.7 an approximation of the optimal policy can then be obtained as $\mathbf{a}^* := (a_{\text{loc}}^*(X_0), a_{\text{loc}}^*(X_1), \dots) \in \mathcal{A}$.

Algorithm 1: Value Iteration

Input : Batch Size $B \in \mathbb{N}$; Hyperparameters for the neural networks; Number of epochs E ; Number of iterations Iter_v for the improvement of the value function; Number of iterations Iter_a for the improvement of the action function; Number of measures $N_{\mathcal{P}}$; Number of Monte-Carlo simulations N_{MC} ; State space Ω_{loc} ; Action space A ; Reward function r ;

Initialize a neural network V^0 ;

Initialize a neural network a_{loc} ;

for epoch = 1, ..., E **do**

Set $V^{\text{epoch}} = V^{\text{epoch}-1}$ and freeze the weights of $V^{\text{epoch}-1}$;

for iteration = 1, ..., Iter_a **do**

// We maximize $\inf_{\mathbb{P}_0 \in \mathcal{P}(x, a_{\text{loc}}(x))} \mathbb{E}_{\mathbb{P}_0}[r(x, a_{\text{loc}}(x), X_1) + \alpha V^{\text{epoch}-1}(X_1)]$ with respect to $a_{\text{loc}} \in \mathcal{A}_{\text{loc}}$.

Sample a batch of states $(x_i)_{i=1, \dots, B} \subseteq \Omega_{\text{loc}}$;

for $i = 1, \dots, B$ **do**

 | Pick measures $\mathbb{P}_1^{(i)}, \dots, \mathbb{P}_{N_{\mathcal{P}}}^{(i)} \in \mathcal{P}(x_i, a_{\text{loc}}(x_i))$;

end

Denote by $X_{1, \mathbb{P}}^{(j)}$ a random variable that is sampled according to a measure $\mathbb{P} \in \mathcal{M}_1(\Omega_{\text{loc}})$ for $j = 1, \dots, N_{\text{MC}}$;

Sample $X_{1, \mathbb{P}}^{(j)}$ for all $\mathbb{P} \in \{\mathbb{P}_1^{(i)}, \dots, \mathbb{P}_{N_{\mathcal{P}}}^{(i)}\}$, $j \in \{1, \dots, N_{\text{MC}}\}$, $i \in \{1, \dots, B\}$;

Define for $i = 1, \dots, B$:

$$\widehat{\mathcal{T}V}(x_i) := \min_{\mathbb{P} \in \{\mathbb{P}_1^{(i)}, \dots, \mathbb{P}_{N_{\mathcal{P}}}^{(i)}\}} \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} r(x_i, a_{\text{loc}}(x_i), X_{1, \mathbb{P}}^{(j)}) + \alpha V^{\text{epoch}-1}(X_{1, \mathbb{P}}^{(j)});$$

Maximize

$$\sum_{i=1}^B \widehat{\mathcal{T}V}(x_i)$$

with respect to parameters from the neural network a_{loc} (e.g. back-propagation with a stochastic gradient descent algorithm).

end

for iteration = 1, ..., Iter_v **do**

// We minimize the quadratic error between $V^{\text{epoch}}(x)$ and the approximation of

$\sup_{a_{\text{loc}} \in \mathcal{A}_{\text{loc}}} \inf_{\mathbb{P}_0 \in \mathcal{P}(x, a_{\text{loc}}(x))} \mathbb{E}_{\mathbb{P}_0}[r(x, a_{\text{loc}}(x), X_1) + \alpha V^{\text{epoch}-1}(X_1)] = \mathcal{T}V^{\text{epoch}-1}(x)$, that was computed in the previous step, for all states x from a batch of sampled values.

Sample Batch of states $(x_i)_{i=1, \dots, B} \subseteq \Omega_{\text{loc}}$;

Minimize

$$\sum_{i=1}^B \left(V^{\text{epoch}}(x_i) - \widehat{\mathcal{T}V}(x_i) \right)^2$$

with respect to parameters from V^{epoch} .

end

end

Output: Neural network V^E approximating the optimal value function;

Neural network a_{loc} approximating the optimal one-step policy;

APPENDIX B. SUPPLEMENTARY RESULTS

The first auxiliary result is Banach's fix point theorem, compare, e.g., [7, Theorem A 3.5], or any standard monograph on analysis or functional analysis.

Theorem B.1 (Banach's Fix Point Theorem). *Let M be a complete metric space with metric $d(x, y)$ and let $\mathcal{T} : M \rightarrow M$ be an operator such that there exists a number $\beta \in (0, 1)$ such that $d(\mathcal{T}v, \mathcal{T}w) \leq \beta d(v, w)$ for all $v, w \in M$. Then, we have that*

- (i) \mathcal{T} has a unique fix point v^* in M , i.e., $\mathcal{T}v^* = v^*$.
- (ii) $\lim_{n \rightarrow \infty} \mathcal{T}^n v = v^*$ for all $v \in M$.
- (iii) For $v \in M$ we obtain

$$d(v^*, \mathcal{T}^n v) \leq \frac{\beta^n}{1 - \beta} d(\mathcal{T}v, v).$$

The following result, Berge's Maximum Theorem, can for example be found in [2, Theorem 17.31].

Theorem B.2 (Berge's Maximum Theorem). *Let $\varphi : X \rightrightarrows Y$ be an upper and lower hemicontinuous correspondence between topological spaces with nonempty compact values, and suppose that $f : \{(x, y) \in X \times Y \mid y \in \varphi(x)\} \rightarrow \mathbb{R}$ is continuous. Then the following holds.*

- (i) The function

$$\begin{aligned} m : X &\rightarrow \mathbb{R} \\ x &\mapsto \max_{y \in \varphi(x)} f(x, y) \end{aligned}$$

is continuous.

- (ii) The correspondence

$$\begin{aligned} c : X &\rightrightarrows Y \\ x &\mapsto \{y \in \varphi(x) \mid f(x, y) = m(x)\} \end{aligned}$$

has nonempty, compact values.

- (iii) If Y is Hausdorff, then c is upper hemicontinuous.

The following two lemmas provide characterizations of upper and lower hemicontinuity, respectively. The results can be found, e.g., in [2, Theorem 17.20], and [2, Theorem 17.21].

Lemma B.3 (Upper Hemicontinuity). *Assume that the topological space X is first countable and that Y is metrizable. Then, for a correspondence $\varphi : X \rightrightarrows Y$ the following statements are equivalent.*

- (i) The correspondence φ is upper hemicontinuous and $\varphi(x)$ is compact for all $x \in X$.
- (ii) For any $x \in X$, if a sequence $((x^{(n)}, y^{(n)}))_{n \in \mathbb{N}} \subseteq \text{Gr } \varphi$ satisfies $x^{(n)} \rightarrow x$ for $n \rightarrow \infty$, then there exists a subsequence $(y^{(n_k)})_{k \in \mathbb{N}}$ with $y^{(n_k)} \rightarrow y \in \varphi(x)$ for $k \rightarrow \infty$.

Lemma B.4 (Lower Hemicontinuity). *For a correspondence $\varphi : X \rightrightarrows Y$ between first countable topological spaces the following statements are equivalent.*

- (i) The correspondence φ is lower hemicontinuous.
- (ii) For any $x \in X$, if $x^{(n)} \rightarrow x$ for $n \rightarrow \infty$, then for each $y \in \varphi(x)$ there exists a subsequence $(x^{(n_k)})_{k \in \mathbb{N}}$ and elements $y^{(k)} \in \varphi(x^{(n_k)})$ for each $k \in \mathbb{N}$ such that $y^{(k)} \rightarrow y$ for $k \rightarrow \infty$.

REFERENCES

- [1] Victor Aguirregabiria and Pedro Mira. Swapping the nested fixed point algorithm: A class of estimators for discrete Markov decision models. *Econometrica*, 70(4):1519–1543, 2002.
- [2] Charalambos D. Aliprantis and Kim C. Border. *Infinite dimensional analysis*. Springer, Berlin, third edition, 2006. A hitchhiker's guide.
- [3] Andrea Angiuli, Nils Detering, Jean-Pierre Fouque, and Jimin Lin. Reinforcement learning algorithm for mixed mean field control games. *arXiv preprint arXiv:2205.02330*, 2022.
- [4] Andrea Angiuli, Jean-Pierre Fouque, and Mathieu Lauriere. Reinforcement learning for mean field games, with applications to economics. *arXiv preprint arXiv:2106.13755*, 2021.
- [5] Nicole Bäuerle and Alexander Glauner. Q-learning for distributionally robust Markov decision processes. In *Modern Trends in Controlled Stochastic Processes*, pages 108–128. Springer, 2021.
- [6] Nicole Bäuerle and Ulrich Rieder. MDP algorithms for portfolio optimization problems in pure jump markets. *Finance and Stochastics*, 13(4):591–611, 2009.

- [7] Nicole Bäuerle and Ulrich Rieder. *Markov decision processes with applications to finance*. Springer Science & Business Media, 2011.
- [8] Claude Berge. *Espaces topologiques: Fonctions multivoques*. Collection Universitaire de Mathématiques, Vol. III. Dunod, Paris, 1959.
- [9] Francesco Bertoluzzo and Marco Corazza. Reinforcement learning for automatic financial trading: Introduction and some applications. *University Ca'Foscari of Venice, Dept. of Economics Research Paper Series No*, 33, 2012.
- [10] Vladimir Igorevich Bogachev. *Gaussian measures*. Number 62. American Mathematical Soc., 1998.
- [11] Stephen Boyd, Enzo Busseti, Steve Diamond, Ronald N Kahn, Kwangmoo Koh, Peter Nystrup, and Jan Speth. Multi-period trading via convex optimization. *Foundations and Trends® in Optimization*, 3(1):1–76, 2017.
- [12] Jay Cao, Jacky Chen, John Hull, and Zissis Poulos. Deep hedging of derivatives using reinforcement learning. *The Journal of Financial Data Science*, 3(1):10–27, 2021.
- [13] Ying-Hua Chang and Ming-Sheng Lee. Incorporating Markov decision process on genetic algorithms to formulate trading strategies for stock markets. *Applied Soft Computing*, 52:1143–1153, 2017.
- [14] Zhi Chen, Pengqian Yu, and William B Haskell. Distributionally robust optimization for sequential decision-making. *Optimization*, 68(12):2397–2426, 2019.
- [15] Philippe Clément and Wolfgang Desch. Wasserstein metric and subordination. *Studia Mathematica*, 1(189):35–52, 2008.
- [16] Matthew F Dixon, Igor Halperin, and Paul Bilokon. *Machine Learning in Finance*. Springer, 2020.
- [17] Jiayi Du, Muyang Jin, Petter N Kolm, Gordon Ritter, Yixuan Wang, and Bofei Zhang. Deep reinforcement learning for option replication and hedging. *The Journal of Financial Data Science*, 2(4):44–57, 2020.
- [18] Angelos Filos. *Reinforcement learning for portfolio management*. PhD thesis, Imperial College London, 2019.
- [19] Carl Gold. FX trading via recurrent reinforcement learning. In *2003 IEEE International Conference on Computational Intelligence for Financial Engineering, 2003. Proceedings.*, pages 363–370. IEEE, 2003.
- [20] Allan Gut. The multivariate normal distribution. In *An Intermediate Course in Probability*, pages 117–145. Springer, 2009.
- [21] Igor Halperin. QLBS: Q-learner in the Black-Scholes (-Merton) worlds. *The Journal of Derivatives*, 28(1):99–122, 2020.
- [22] Ben Hambly, Renyuan Xu, and Huining Yang. Recent advances in reinforcement learning in finance. *arXiv preprint arXiv:2112.04553*, 2021.
- [23] Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural networks*, 4(2):251–257, 1991.
- [24] Yuh-Jong Hu and Shang-Jen Lin. Deep reinforcement learning for optimizing finance portfolio management. In *2019 Amity International Conference on Artificial Intelligence (AICAI)*, pages 14–20. IEEE, 2019.
- [25] Jean Jacod and Philip Protter. *Probability essentials*. Universitext. Springer-Verlag, Berlin, second edition, 2003.
- [26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [27] Frank Hyneman Knight. *Risk, uncertainty and profit*, volume 31. Houghton Mifflin, 1921.
- [28] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- [29] Yuxi Li, Csaba Szepesvari, and Dale Schuurmans. Learning exercise policies for American options. In *Artificial Intelligence and Statistics*, pages 352–359. PMLR, 2009.
- [30] Eva Lütkebohmert, Thorsten Schmidt, and Julian Sester. Robust deep hedging. *Quantitative Finance*, 2022.
- [31] John Moody, Lizhong Wu, Yuansong Liao, and Matthew Saffell. Performance functions and reinforcement learning for trading systems and portfolios. *Journal of Forecasting*, 17(5-6):441–470, 1998.
- [32] Ariel Neufeld, Julian Sester, and Daiying Yin. Detecting data-driven robust statistical arbitrage strategies with deep neural networks. *arXiv preprint arXiv:2203.03179*, 2022.
- [33] John Rust. Structural estimation of Markov decision processes. *Handbook of econometrics*, 4:3081–3143, 1994.
- [34] Manfred Schäl. Markov decision processes in finance and dynamic options. In *Handbook of Markov decision processes*, pages 461–487. Springer, 2002.
- [35] Salman Sadiq Shuvo, Yasin Yilmaz, Alan Bush, and Mark Hafen. A Markov decision process model for socio-economic systems impacted by climate change. In *International Conference on Machine Learning*, pages 8872–8883. PMLR, 2020.
- [36] Sorawoot Srisuma and Oliver Linton. Semiparametric estimation of Markov decision processes with continuous state space. *Journal of Econometrics*, 166(2):320–341, 2012.
- [37] Kerem Uğurlu. Robust optimal control using conditional risk mappings in infinite horizon. *Journal of Computational and Applied Mathematics*, 344:275–287, 2018.
- [38] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer, 2009.
- [39] Douglas J White. A survey of applications of markov decision processes. *Journal of the operational research society*, 44(11):1073–1096, 1993.
- [40] Zhuoran Xiong, Xiao-Yang Liu, Shan Zhong, Hongyang Yang, and Anwar Walid. Practical deep reinforcement learning approach for stock trading. *arXiv preprint arXiv:1811.07522*, 2018.
- [41] Huan Xu and Shie Mannor. Distributionally robust markov decision processes. *Mathematics of Operations Research*, 37(2):288–300, 2012.
- [42] Pengqian Yu, Joon Sern Lee, Ilya Kulyatin, Zekun Shi, and Sakyasingha Dasgupta. Model-based deep reinforcement learning for dynamic portfolio optimization. *arXiv preprint arXiv:1901.08740*, 2019.
- [43] Man-Chung Yue, Daniel Kuhn, and Wolfram Wiesemann. On linear optimization over wasserstein balls. *arXiv preprint arXiv:2004.07162*, 2020.

- [44] Zihao Zhang, Stefan Zohren, and Stephen Roberts. Deep reinforcement learning for trading. *The Journal of Financial Data Science*, 2(2):25–40, 2020.